# Bounded Version Vectors

José Bacelar Almeida       Paulo Sérgio Almeida       Carlos Baquero

Departamento de Informática, Universidade do Minho
{jba,psa,cbm}@di.uminho.pt

## Abstract

Version vectors play a central role in update tracking under optimistic distributed systems, allowing the detection of obsolete or inconsistent versions of replicated data. Version vectors do not have a bounded representation; they are based on integer counters that grow indefinitely as updates occur. Existing approaches to this problem are scarce; the mechanisms proposed are either unbounded or operate only under specific settings. This paper examines version vectors as a mechanism for data causality tracking and clarifies their role with respect to vector clocks. Then, it introduces bounded stamps and proves them to be a correct alternative to integer counters in version vectors. The resulting mechanism, bounded version vectors, represents the first bounded solution to data causality tracking between replicas subject to local updates and pairwise symmetrical synchronization.

**Keywords:** Replication, causality, version vectors, update tracking, bounded state.

## 1   Introduction

Optimistic replication is a critical technology in distributed systems, in particular when improving availability of database systems and adding support to mobility and partitioned operation [17]. Under optimistic replication, data replicas can evolve autonomously by incorporation new updates into their state. Thus, when contact can be established between two or more replicas, mutual consistency must be evaluated and potential divergence detected.

The classic mechanism for assessing divergence between mutable replicas is provided by *version vectors* which, since their introduction by Parker et al [13], have been one of the cornerstones of optimistic data management. Version vectors associate to each replica a vector of integer counters that keeps track of the last update that is known to have been originated in every other replica and in the replica itself. The mechanism is simple and intuitive but requires a state of unbounded size, since each counter in the vector can grow indefinitely.

The potential existence of a bounded substitute to version vectors has been overlooked by the community. A possible cause is a frequent confusion of the roles played by *version vectors* and *vector clocks* (e.g. [16, 17]), that have the same representation [13, 4, 12], together with the existence of a minimality result by Charron-Bost [3], stating that vector clocks are the most concise characterization of causality among process events.

In this article we show that a bounded solution is possible for the problem addressed by version vectors: the detection of mutual inconsistency between replicas subject to local updates and pairwise symmetrical synchronization. We present a mechanism, *bounded stamps*, that can be used to replace integer counters in version vectors, stressing that the minimality result that precludes bounded vector clocks does not apply to version vectors.

### 1.1   On version vectors and vector clocks

Asynchronous distributed systems track causality and logical time among communicating processes by means of several mechanisms [11, 18], in particular vector clocks [4, 12].

While being structurally equivalent to version vectors, vector clocks serve a very distinct purpose. Vector clocks track causality by establishing a strict partial order on the events of processes that communicate by message passing, and are known to be the most concise solution to this problem. Vector clocks, being a vector of integer counters, are unbounded in size, but so is the number of events that must be ordered and timestamped by them. In short, *vector clocks order an unlimited number of events occurring in a*

*given number of processes.*

If we consider the role of version vectors, data causality, there is always a limit to the number of possible relations that can be established on the set of replicas. This limit is independent on the number of update events that are considered on any given run. For example, in a two replica system $\{r_a, r_b\}$ only four cases can occur: $r_a = r_b$, $r_a < r_b$, $r_b > r_a$ and $r_a \parallel r_b$. If the two replicas are already divergent the inclusion of *new* update events on any of the replicas does not change their mutual divergence and the corresponding relation between them. In short, *version vectors order a given number of replicas, according to an unlimited number of update events.*

The existence of a limited number of relations is a necessary but not sufficient condition for the existence of a bounded characterization mechanism. A relation, which is a global abstraction, must be encoded and computed through local operations on replica pairs without the need for a global view. This is one of the important properties of version vectors.

## 2 Data causality and version vectors

Data causality on a set of replicas can be assessed via set inclusion of the sets of update events known to each replica. Data causality is the pre-order defined by:

$$r_a \leq r_b \quad \text{iff} \quad U_a \subseteq U_b$$

being $U_a$ and $U_b$ the sets of update events (globally unique events), known to replicas $r_a$ and $r_b$.

When tracking data causality with version vectors in a $N$ replica system, one associates to each replica $r_i \in \{r_0, \ldots, r_{N-1}\}$ a vector $V_i$ of $N$ integer counters. The order on version vectors is the standard pointwise (coordinatewise) order:

$$V_a \leq_{\mathsf{V}} V_b \quad \text{iff} \quad \forall k.\, V_a^k \leq V_b^k$$

where $V_i^k$ denotes component $k$ of vector $V_i$.

The operations on version vectors, formally presented in Figure 1, are as follows:

**Initialization** (I) establishes the initial system state. All vectors are initialized with zeroes.

**Update** ($\mathsf{U}^a$) an update event in replica $r_a$ increments $V_a^a$.

Operation I:

$$(V_i^k)' \quad = \quad 0.$$

Operation $\mathsf{U}^a$:

$$(V_i^k)' \quad = \quad \begin{cases} V_i^k + 1 & \text{if } i = k = a; \\ V_i^k & \text{otherwise.} \end{cases}$$

Operation $\mathsf{S}^{ab}$:

$$(V_a^k)' = (V_b^k)' \quad = \quad V_a^k \sqcup V_b^k.$$

Figure 1: Semantics of version vector operations.

**Synchronization** ($\mathsf{S}^{ab}$) synchronization of $r_a$ and $r_b$ is achieved by taking the pointwise join (greatest element) of $V_a$ and $V_b$.

This classic mechanism encodes data causality because comparing version vectors gives the same result as comparing sets of known update events. For all runs and replicas $r_a$ and $r_b$:

$$r_a \leq r_b \quad \text{iff} \quad U_a \subseteq U_b \quad \text{iff} \quad V_a \leq_{\mathsf{V}} V_b.$$

Figure 2 shows a run with version vectors in a four replica system. Updates are depicted by a "•" and synchronization by two "○" connected by a line.

### 2.1 Version vector slices

All operations over version vectors exhibit a pointwise nature: a given vector position is only compared or updated to the same position in other vectors, resulting from all information about updates originated in replica $r_k$ being stored in component $k$ of each version vector. This allows a decomposition of the replicated system into $N$ *slices*, where each slice represents the updates that were originated in a given replica. Slice $i$ for a $N$ replica system is made up of the $i^{\text{th}}$ component of each version vector:

$$\langle V_0^i, \ldots, V_{N-1}^i \rangle.$$

This means that data causality in $N$ replicas can be encoded by the concatenation of the representation for each of the $N$ slices. It also means that it is enough to concentrate on a subproblem: encoding the distributed knowledge about a single source of updates, and the corresponding version vector slice (VVS). The source of updates increments
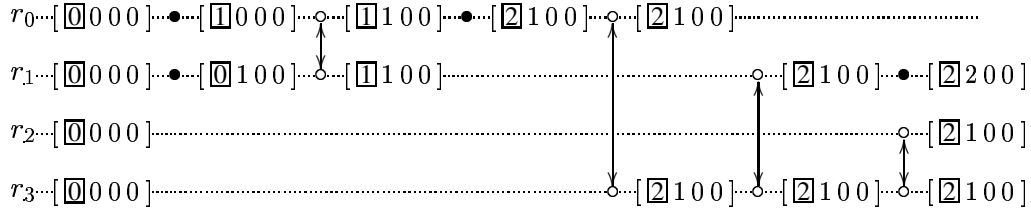
2

$r_0 \cdots [\boxed{0}\,0\,0\,0] \cdots \bullet \cdots [\boxed{1}\,0\,0\,0] \cdots \circ \cdots [\boxed{1}\,1\,0\,0] \cdots \bullet \cdots [\boxed{2}\,1\,0\,0] \cdots \circ \cdots [\boxed{2}\,1\,0\,0] \cdots$

$r_1 \cdots [\boxed{0}\,0\,0\,0] \cdots \bullet \cdots [\boxed{0}\,1\,0\,0] \cdots \circ \cdots [\boxed{1}\,1\,0\,0] \cdots\cdots\cdots [\boxed{2}\,1\,0\,0] \cdots \bullet \cdots [\boxed{2}\,2\,0\,0]$

$r_2 \cdots [\boxed{0}\,0\,0\,0] \cdots\cdots\cdots\cdots\cdots\cdots\cdots \circ \cdots [\boxed{2}\,1\,0\,0]$

$r_3 \cdots [\boxed{0}\,0\,0\,0] \cdots\cdots\cdots\cdots \circ \cdots [\boxed{2}\,1\,0\,0] \cdots \circ \cdots [\boxed{2}\,1\,0\,0] \cdots \circ \cdots [\boxed{2}\,1\,0\,0]$

Figure 2: Version Vectors: example run, depicting slice 0 counters by a boxed digit.

Operation I:

$$(S_i)' = 0.$$

Operation $\mathsf{U}^0$:

$$(S_i)' = \begin{cases} S_i + 1 & \text{if } i = 0; \\ S_i & \text{otherwise.} \end{cases}$$

Operation $\mathsf{S}^{ab}$:

$$(S_a)' = (S_b)' = S_a \sqcup S_b.$$

Figure 3: VVS semantics for slice 0.

its counter and all other replicas keep potentially outdated copies of that counter; this subproblem amounts to storing a distributed representation of a total order.

For the remainder of the paper we will concentrate, for notational convenience and without loss of generality, on finding a bounded representation for slice 0. Figure 3 presents the semantics of version vectors restricted to slice 0; in the run presented in Figure 2 this slice is shown using boxed counters.

# 3  Informal presentation

We now give an informal presentation of the mechanism and give some intuition of how it works and how it accomplishes its purpose. Having shown that it is enough to concentrate on a subproblem (a single source of updates) and the corresponding slice of version vectors, we now present the stamp that will replace, in each replica, the integer counter of the corresponding version vector.

For problem size $N$, i.e. assuming $N$ replicas, with $r_0$ the "primary" where updates take place and $r_1, \ldots, r_{N-1}$ the "secondary" replicas, we represent a stamp by something like

| | | |
|---|---|---|
| **c** | **b** | a |
| **c** | a | |
| **a** | | |
| **c** | a | |

It has a representation of bounded size, as it consists of $N$ rows, each with at most $N$ symbols (letters here), taken from a finite set $\mathcal{S}_N$. An example run consisting of four replicas is presented in Figure 4.

A stamp is, in abstract, a vector of totally ordered sets. Each of the $N$ components (rows in our notation) represents a total order, with the greatest element on the left (the first row above means $c > b > a$). In a stamp for replica $r_i$, row $i$ ($i \in \{0, \ldots N - 1\}$) is what we call the *principal order* (displayed with a gray background), while the other rows are the *cached orders*. (Thus, the stamp above would belong to replica $r_3$.) The cached order in row $j$ represents the principal order of replica $j$ at some point in time, propagated to replica $i$ (either directly or indirectly through several synchronizations).

The greatest element of the principal order (on the left, depicted in bold over gray) is what we call the *principal element*. It represents the most recent update (in the primary) known by the replica. In a representation using an infinite total ordered set instead of $\mathcal{S}_N$ nothing more would be needed. This element can be thought of as "corresponding" to the value of the integer counter in version vectors.

The left column in a stamp (depicted in bold) is what we call the *principal vector*; it is made up of the greatest element of each order (row). It represents the most recent local knowledge about the principal element of each replica (including itself).

In a stamp, there is a relationship between the principal order and the principal vector: the elements in the principal vector are the same ones as in the principal order. In other words, the set of elements in the principal vector is ordered according to the principal order.

## 3.1  Comparison and synchronization as well defined local operations

As we will show below, the mechanism is able to compare two stamps by a local operation on the respective principal orders. No global knowledge is used: not even a global
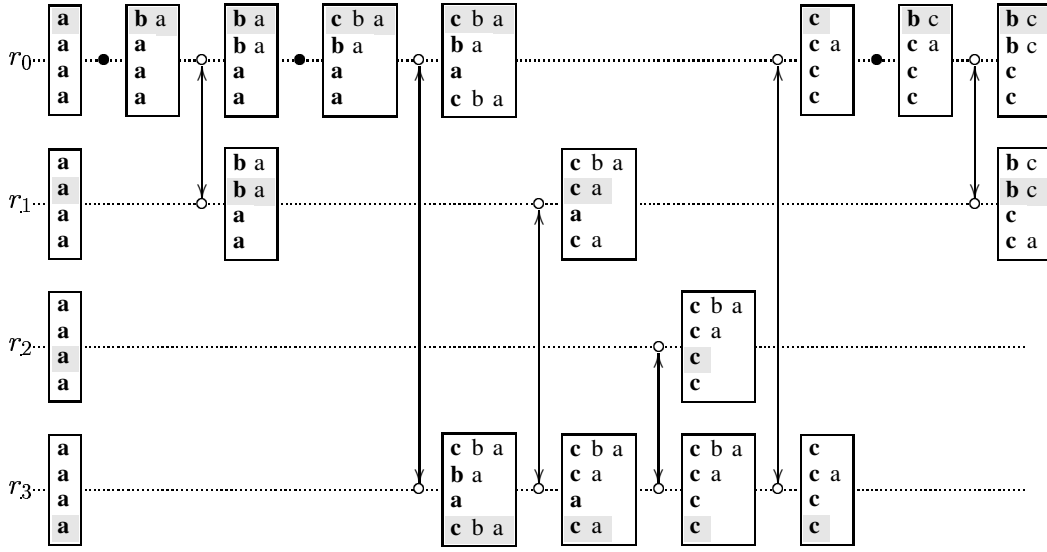
Figure 4: Bounded stamps: example run.

order on the set of symbols $\mathcal{S}_N$ is assumed. For comparison purposes $\mathcal{S}_N$ is simply an unordered set, with elements that are ordered differently in different stamps. As an example, the comparison of

$$r_0 = \begin{array}{|c|} \hline \mathbf{b}\ \text{c} \\ \mathbf{c}\ \text{a} \\ \mathbf{c} \\ \mathbf{c} \\ \hline \end{array} \quad \text{with} \quad r_1 = \begin{array}{|c|} \hline \mathbf{c}\ \text{b}\ \text{a} \\ \mathbf{c}\ \text{a} \\ \mathbf{a} \\ \text{c}\ \text{a} \\ \hline \end{array}$$

involves looking at **b** c and **c** a , and gives $r_0 > r_1$.

When synchronizing two stamps, in the positions of the two principal elements, the resulting value will be the maximum of the two principal elements; the rest of the resulting principal vector will be the pointwise maximum of the respective values. The comparisons are performed according to the principal orders of the two stamps involved.

Is is important to notice that, in general, it is not possible to take two arbitrary total orders and merge them into a new total order. As such, it could be thought that computing the maximum as mentioned above is ill defined. As we will show, several properties of the model can be explored that make these operations indeed possible and well defined. We will also show that it is possible to totally order the elements in the resulting principal vector, i.e. to obtain a new principal order.

## 3.2  Garbage collection for symbol reuse

The boundedness of the mechanism is only possible through symbol reuse. When an update operation is per-formed, instead of incrementing an integer counter, some symbol is chosen to become the new principal element. By using a finite set of symbols $\mathcal{S}_N$, an update will eventually reuse a symbol that was already used in the past to represent some previous update that has been synchronized with other replicas.

However, by reusing symbols, an obvious problem arises that needs to be addressed: the symbol reuse cannot compromise the well-definedness of the comparison operations described above. As an example, it would not be acceptable that, due to reuse, the principal orders of two stamps end up being **a** b c and **c** a , as it would not be possible to overcome the ambiguity between $a > b > c$ and $c > a$ and to infer which one is the greatest stamp.

To address the problem, the mechanism implements a distributed "garbage collection" of symbols. This is accomplished through the extra information in the cached orders. As we will show, any element in the principal order/vector of any replica is also present in the primary replica (in some principal or cached order). This is the key property towards symbol reuse: when an update is performed, any symbol which is not present in the primary replica is considered "garbage" and can be (re)used for the new principal element.

As an example, in Figure 4, when the final update occurs, symbol $b$ can be used for the new principal element because
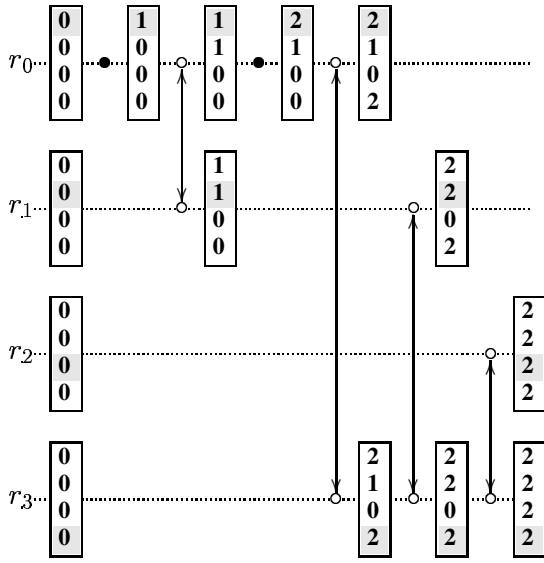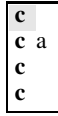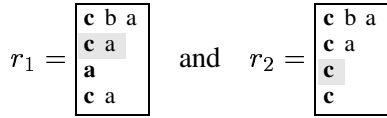
4

Figure 5: Counter mode principal vectors.

it is not present in the primary replica:

$$\begin{array}{l} \mathbf{c} \\ \mathbf{c}\ a \\ \mathbf{c} \\ \mathbf{c} \end{array}$$

Notice that the scheme only assures that $b$ does not occur in the principal orders/vectors. In this example $b$ occurs in some cached orders of replicas

$$r_1 = \begin{array}{l} \mathbf{c}\ b\ a \\ \mathbf{c}\ a \\ \mathbf{a} \\ \mathbf{c}\ a \end{array} \quad \text{and} \quad r_2 = \begin{array}{l} \mathbf{c}\ b\ a \\ \mathbf{c}\ a \\ \mathbf{c} \\ \mathbf{c} \end{array}$$

but this is not a problem because those elements will not be used in comparisons; the "old" $b$ will not be confused with the "new" $b$.

## 3.3 Synopsis of formal presentation

The formal presentation and proof of correctness will make use of an unbounded mechanism which we call the *counter mode principal vectors* (CMPV). This auxiliary mechanism represents what the evolution of the principal vector would be if we could afford to use integer counters. The mechanism makes use of the total order on natural numbers and does not encode orders locally. In Figure 5 we present part of the run in Figure 4 using the counter mode mechanism.

The bulk of the proof consists in establishing several properties of the CMPV model that allow the relevant comparison operations to be computed in a well-defined way

Operation I:
$$(\mathsf{a}^k)' = 0.$$

Operation $\mathsf{U}^0$:
$$(\mathsf{a}^k)' = \begin{cases} \mathsf{a}^k + 1 & \text{if } a = k = 0; \\ \mathsf{a}^k & \text{otherwise.} \end{cases}$$

Operation $\mathsf{S}^{ab}$:
$$(\mathsf{a}^k)' = (\mathsf{b}^k)' = \begin{cases} \mathsf{a}^a \sqcup \mathsf{b}^b & \text{if } k \in \{a, b\}; \\ \mathsf{a}^k \sqcup \mathsf{b}^k & \text{otherwise.} \end{cases}$$

Figure 6: Semantics of operations in CMPV.

using only local information. The key idea is that, exploiting these properties, bounded stamps can be seen as an encoding of CMPV using a finite set $\mathcal{S}_N$, where the principal orders are used to encode the relevant order information.

## 4 Counter Mode Principal Vectors

Version Vector Slices (VVS) rely on an unbounded totally ordered set — the integers. Their unbounded nature is actually a consequence of adopting a predetermined order relation (and hence globally known) to capture data causality among replicas. To overcome this, we enrich VVS in a way that order judgments become, in a sense, local to each replica. In this way, it will be possible to dynamically encode the causality order and open the perspective of bounding the "counters" domain.

For a replica index $a$, its local state in the CMPV model is denoted by $\mathsf{C}_a$ and defined as the tuple $\langle a, \mathsf{a} \rangle$ where $\mathsf{a}$ is a vector of integers with size $N$ — the *principal vector* for $\mathsf{C}_a$ (see Figure 5). The value in position $k$ of vector $\mathsf{a}$ is denoted by $\mathsf{a}^k$ and represents the knowledge of stamp $\mathsf{C}_a$ concerning the most recent update known by stamp $\mathsf{C}_k$. The element $\mathsf{a}^a$ plays a central role since it holds $\mathsf{C}_a$'s view about the more recent update — this is essentially the information contained in VVS counters and we call it the *principal element* for stamp $\mathsf{C}_a$.

Figure 6 defines the semantics of the operations in the CMPV model. Symbol $\sqcup$ denotes the join operation under integer ordering (i.e. taking the maximum element). Notice that the order information is only required to perform the synchronization operation. Moreover, comparisons are always between principal elements or pointwise (between the same position in two principal vector). Occasionally, it will be convenient to write $\mathsf{a} \sqcup \mathsf{b}$ for the result of the synchro-

5

nization on stamps $C_a$ and $C_b$ (i.e. the principal vector of one of these stamps after synchronization).

A *trace* consists of a sequence of operations starting with | and followed by an arbitrary number of updates and synchronizations. In the remainder, when stating properties in the CMPV, we will leave implicit that they only refer to reachable states, i.e. states that result from some trace of operations. Induction over the traces is the fundamental tool to prove invariance properties, as the following simple facts about CMPV.

**Proposition 1.** *For every replica $C_a$, $C_b$ and index $k$,*

*1.* $a^b \leq b^b$,

*2.* $a^a \leq 0^0$,

*3.* $a^k \leq a^a$.

*Proof.* Simple induction on the length of traces. $\square$

Given stamps $C_a$ and $C_b$ we define their *data causality order under CMPV* ($\leq_C$) as the comparison of their principal elements:

$$C_a \leq_C C_b \quad \text{iff} \quad a^a \leq b^b.$$

By Figure 6 it can be seen that the computation of principal elements only depends upon principal elements. Moreover, if we restrict the impact of the operations to the principal element we recover the VVS semantics (Figure 3). This observation leads immediately to the correctness of CMPV as a data causality encoding for slice 0:

$$C_a \leq_C C_b \quad \text{iff} \quad V_a^0 \leq_V V_b^0.$$

This result is not surprising since CMPV was defined as a semantics preserving extension of VVS.

Next we will show that the additional information contained in the CMPV model makes it possible to avoid relying on the integer order, and to replace it with a locally encoded order. For this, we will use a non-trivial invariant on the global state given by the following lemma. Its proof is presented in the appendix since it requires an auxiliary definition and some additional lemmata.

**Lemma 2.** *For every stamp $C_a$ and $C_b$ and index $k$,*

$$a^a \leq b^b \ \text{ and } \ b^k \leq a^k \quad \textit{implies} \quad a^k \in b.$$

*Proof.* See appendix A. $\square$

Recall that the order information is only required to perform the synchronization operation. Moreover, comparisons are always between principal elements or pointwise (between the same position in two principal vector). In the following we will show that these comparisons can be performed without relying on integer order as long as we can order the elements in the principal vector of each stamp individually.

Comparison between principal elements reduces to a membership testing.

**Proposition 3.** *For every stamp $C_a$, $C_b$,*

$$a^a \leq b^b \quad \textit{iff} \quad a^a \in b.$$

*Proof.* $\implies$ If $a^a \leq b^b$ then, by Proposition 1(1) we have that $b^a \leq a^a$ and so, by Lemma 2, $a^a \in b$.

$\impliedby$ If $a^a \in b$ then, by Proposition 1(3) we have that $a^a \leq b^b$. $\square$

For a stamp $C_a$, let us denote by $\leq^a$ the restriction of the intrinsic integer order to the values contained in the principal vector $a$:

$$x \leq^a y \quad \text{iff} \quad x \leq y \text{ and } x \in a \text{ and } y \in a.$$

Using these orderings, we define new ones that are appropriate to perform the required comparisons. For stamps $C_a$ and $C_b$, let their combined order $\leq^{ab}$ be defined as:

$$x \leq^{ab} y \quad \text{iff} \quad (b^b \in a \text{ and } (x \in a \Rightarrow x \leq^a y)) \text{ or } (a^a \in b \text{ and } (x \in b \Rightarrow x \leq^b y)).$$

For convenience, we also define the corresponding join operation $\underset{ab}{\sqcup}$ as:

$$x \underset{ab}{\sqcup} y = \begin{cases} y & \text{if } x \leq^{ab} y, \\ x & \text{otherwise.} \end{cases}$$

The following proposition establishes the claimed properties for this ordering.

**Proposition 4.** *For every stamp $C_a$ and $C_b$ and index $k$,*

*1.* $a^a \leq b^b \quad \textit{iff} \quad a^a \leq^{ab} b^b$,

*2.* $a^k \leq b^k \quad \textit{iff} \quad a^k \leq^{ab} b^k$.

6

*Proof.* (1) Follows directly from Propositions 1 and 3.

(2) $\Longrightarrow$ Let $\mathsf{a}^k \leq \mathsf{b}^k$. When $\mathsf{b}^b \leq \mathsf{a}^a$ Proposition 3 guaranties that $\mathsf{b}^b \in \mathsf{a}$ and, by Lemma 2, we have $\mathsf{b}^k \in \mathsf{a}$ and then $\mathsf{a}^k \leq^{\mathsf{a}} \mathsf{b}^k$, which establishes $\mathsf{a}^k \leq^{\mathsf{ab}} \mathsf{b}^k$. The case $\mathsf{a}^a < \mathsf{b}^b$ is trivial since, either $\mathsf{a}^k \in \mathsf{b}$ (in which case $\mathsf{a}^k \leq^{\mathsf{b}} \mathsf{b}^k$), or $\mathsf{a}^k \notin \mathsf{b}$ and so $\mathsf{a}^k \leq^{\mathsf{ab}} \mathsf{b}^k$. $\Longleftarrow$ Let $\mathsf{a}^k \not\leq \mathsf{b}^k$ (that is, $\mathsf{b}^k < \mathsf{a}^k$). The proof proceeds as in the previous implication.

$\square$

Restricted orders can be explicitly encoded (e.g. by a sequence) and can be easily manipulated. We now show that when a synchronization is performed, all the elements in the resulting principal vector were already present in the more up-to-date stamp. This means that the restricted order that results is a restriction of the one from the more up-to-date stamp.

**Proposition 5.** *Let* $\mathsf{C}_a$ *and* $\mathsf{C}_b$ *be stamps and* $\mathsf{C}_x = \mathsf{C}_a \sqcup \mathsf{C}_b$. *If* $\mathsf{a}^a \leq \mathsf{b}^b$ *then, for all* $k$,

$$\mathsf{x}^k \in \mathsf{b}.$$

*Proof.* For the pointwise join $\mathsf{x}^k = \mathsf{a}^k \sqcup \mathsf{b}^k$: if $\mathsf{a}^k \leq \mathsf{b}^k$ then $\mathsf{x}^k = \mathsf{b}^k \in \mathsf{b}$; if $\mathsf{b}^k \leq \mathsf{a}^k$ then, by Lemma 2, $\mathsf{a}^k \in \mathsf{b}$. Otherwise, note that the resulting principal element $(\mathsf{b}^b)$ is already in $\mathsf{b}$. $\square$

These observations together with the fact that the global state can only retain a bounded amount of integer values (an obvious limit is $N^2$) opens the way for a change in the domain from the integers in the CMPV model to a finite set.

# 5 Bounded Stamps

A migration from the domain of integer counters in CMPV to a finite set $\mathcal{S}_N$ is faced with the following difficulty: the update operation should be able to choose a value, that is not present in any principal vector, for the new principal element in the primary.

Adopting a set $\mathcal{S}_N$ sufficiently large (e.g. with $N^2$ elements) guaranties that such a choice exists under a global view. The problem lies in making that choice using only the information in the state of the primary. To overcome this problem we make a new extension of the model that allows the primary to keep track of all the values in use in the principal vectors of all stamps.

We will present this new model parameterized by a set $\mathcal{S}_N$ (the symbol domain), a distinguished element $\mathbf{0} \in$ $\mathcal{S}_N$ (the initial element), and an oracle for new symbols $\text{new}(-)$ (satisfying an axiom described bellow). For each replica index $a$, its local state in the bounded stamps model is denoted by $\mathsf{B}_a$ and defined as $\langle a, \underline{\mathsf{a}}, \boxed{\mathsf{a}} \rangle$ where:

- $a$ is the replica index;

- $\underline{\mathsf{a}}$ is a vector of values from $\mathcal{S}_N$ with size $N$ — the principal vector;

- $\boxed{\mathsf{a}}$ is a vector of $N$ total orders, encoded as sequences, representing the full bounded stamp.

This last component contains all the information in the principal vector, the principal order and the cached orders. Although the principle vector $\underline{\mathsf{a}}$ is redundant (as each component $\underline{\mathsf{a}}^k$ is also present in the first position of each $\boxed{\mathsf{a}}^k$), it is kept in the model for notational convenience in describing the operations and in establishing the correspondence between the models.

The intuitive idea is that the state for each stamp keeps an explicit representation of the restricted orders. More precisely, for stamp $\mathsf{B}_a$, the sequence $\boxed{\mathsf{a}}^a$ contains precisely the elements of $\underline{\mathsf{a}}$ ordered downward (first element is $\underline{\mathsf{a}}^a$). From that sequence one easily defines the restricted order for stamp $\mathsf{B}_a$, what we call *principal order* to emphasize its explicit nature.

$$x \leq_{\mathsf{B}}{}^{\mathsf{a}} y \quad \text{iff} \quad x = y \text{ or } \langle y, x \rangle = \boxed{\mathsf{a}}^a_{|\{x,y\}}$$

where $l_{|m}$ denotes the sequence $l$ restricted to the elements in $m$, i.e. $\langle x \mid x \in l \text{ and } x \in m \rangle$. The combined order $\leq^{\mathsf{ab}}$ and associated join are defined precisely as in counter mode, that is

$$x \leq^{\mathsf{ab}} y \quad \text{iff} \quad \begin{aligned}&(\underline{\mathsf{b}}^b \in \underline{\mathsf{a}} \wedge (x \in \underline{\mathsf{a}} \Rightarrow x \leq_{\mathsf{B}}{}^{\mathsf{a}} y)) \text{ or}\\&(\underline{\mathsf{a}}^a \in \underline{\mathsf{b}} \wedge (x \in \underline{\mathsf{b}} \Rightarrow x \leq_{\mathsf{B}}{}^{\mathsf{b}} y)).\end{aligned}$$

The other sequences in $\boxed{\mathsf{a}}$ keep information about (potentially outdated) principal orders of other stamps — these are called the *cached orders*.

Figure 7 gives the semantics for the operations in this model. The oracle for new symbols $\text{new}(-)$ is a function that gives an element of $\mathcal{S}_N$ satisfying the following axiom:

$$\text{For every stamp } \mathsf{B}_a, \qquad \text{new}(\boxed{\mathsf{0}}) \notin \underline{\mathsf{a}}.$$

The argument $\boxed{\mathsf{0}}$ in the oracle $\text{new}(-)$ intends to emphasize that the choice of the new symbol should be made based on the primary local state.

Operation l:
$$(\underline{a}^k)' = \mathbf{0},$$
$$(\boxed{\underline{a}}^k)' = \langle \mathbf{0} \rangle.$$

Operation $U^0$:
$$(\underline{0}^0)' = \mathrm{new}(\boxed{0}),$$
$$(\boxed{0}^0)' = \mathrm{new}(\boxed{0}) \cdot \boxed{0}^0_{|(\underline{0})'}.$$

Operation $S^{ab}$:
$$(\underline{a}^k)' = (\underline{b}^k)' = \begin{cases} \underline{a}^a \underset{ab}{\sqcup} \underline{b}^b & \text{if } k \in \{a,b\}, \\ \underline{a}^k \underset{ab}{\sqcup} \underline{b}^k & \text{otherwise,} \end{cases}$$

if $k \in \{a, b\}$:
$$(\boxed{\underline{a}}^k)' = (\boxed{\underline{b}}^k)' = \begin{cases} \boxed{\underline{b}}^b_{|(\underline{b})'} & \text{if } \underline{a}^a \in \underline{b}, \\ \boxed{\underline{a}}^a_{|(\underline{a})'} & \text{otherwise,} \end{cases}$$

if $k \neq a$ and $k \neq b$:
$$(\boxed{\underline{a}}^k)' = \begin{cases} \boxed{\underline{b}}^k & \text{if } (\underline{a}^k)' \neq \underline{a}^k, \\ \boxed{\underline{a}}^k & \text{otherwise,} \end{cases}$$
$$(\boxed{\underline{b}}^k)' = \begin{cases} \boxed{\underline{a}}^k & \text{if } (\underline{b}^k)' \neq \underline{b}^k, \\ \boxed{\underline{b}}^k & \text{otherwise.} \end{cases}$$

Figure 7: Semantics of operations on BS model.

Data causality ordering under the Bounded Stamps model is defined by

$$\mathsf{B}_a \leq_\mathsf{B} \mathsf{B}_b \quad \text{iff} \quad \underline{a}^a \in \underline{b}.$$

The correctness of the proposed model follows from the observation that, apart from the cached orders used for the symbol reuse mechanism, it is actually an encoding of the CMPV model. To formalize the correspondence between both models, we introduce an encoding function $[\![ - ]\!]_{-}$ that maps each integer in the CMPV model into the corresponding symbol (in $\mathcal{S}_N$) in the state resulting from a given trace. This map is defined recursively on the traces.

$$[\![ n ]\!]_\mathsf{l} = \mathbf{0},$$
$$[\![ n ]\!]_{\alpha \cdot \mathsf{U}^0} = \begin{cases} \mathrm{new}(\boxed{0}_\alpha) & \text{if } n = |\alpha_{|\mathsf{U}^0}| + 1, \\ [\![ n ]\!]_\alpha & \text{otherwise,} \end{cases}$$
$$[\![ n ]\!]_{\alpha \cdot \mathsf{S}^{xy}} = [\![ n ]\!]_\alpha.$$

Where $|\alpha_{|\mathsf{U}^0}|$ is the number of update events in $\alpha$, $\boxed{0}_\alpha$ is the bounded stamp for the primary after trace $\alpha$, and $\mathrm{new}(\boxed{0}_\alpha)$ gives a canonical choice for the new principal element on the primary after the update. When we discard the cached orders, the semantics of operations given in Figure 7 are precisely the ones in CMPV (Figure 6) affected by the encoding map. Moreover, the principal orders are encodings for the restricted orders presented in the previous section.

**Lemma 6.** *For an arbitrary trace $\alpha$, replicas index $a$ and $b$:*

1. $\underline{a}^k = [\![ \mathsf{a}^k ]\!]_\alpha$,

2. $[\![ \mathsf{a}^i ]\!]_\alpha = [\![ \mathsf{a}^j ]\!]_\alpha$ *implies* $\mathsf{a}^i = \mathsf{a}^j$,

3. $x \leq^\mathsf{a} y$ *iff* $[\![ x ]\!]_\alpha \leq_\mathsf{B}^\mathsf{a} [\![ y ]\!]_\alpha$.

*Proof.* This results from a simple induction on the length of traces. When the last operation was l it is trivial. When it was $\mathsf{U}^0$, the result follows from the induction hypothesis and the axiom for the oracle $\mathrm{new}(-)$. When it was $\mathsf{S}^{xy}$, the result follows from induction hypothesis, the fact that, since $\leq^{\mathsf{ab}}$ computes the required joins (Proposition 4), the definitions of both models are the same, and the correctness of the new restricted orders (Proposition 5). $\square$

As a simple consequence of the previous result, we can state the following correctness result.

**Proposition 7.** *For any arbitrary trace $\alpha$ and replica indexes $a$ and $b$ we have*

$$\mathsf{B}_a \leq_\mathsf{B} \mathsf{B}_b \quad \text{iff} \quad \mathsf{C}_a \leq_\mathsf{C} \mathsf{C}_b.$$

*Proof.* Immediate from Lemma 6 and the definitions of $\leq_\mathsf{B}$ and $\leq_\mathsf{C}$. $\square$

It remains to instantiate the parameters of the model. A trivial but unbounded instantiation would be: set $\mathcal{S}_N$ as the integers, $\mathbf{0}$ as value 0 and $\mathrm{new}(\boxed{0}) = \underline{0}^0 + 1$. In this setting, principal orders would be an explicit representation of counter mode restricted orders. Obviously, we are interested in bounded instantiations of $\mathcal{S}_N$. To show that such instantiations exists, we introduce the following lemma that puts in evidence the role of cached orders. Once again we will postpone its proof to the appendix since it uses a similar technique as the proof of lemma 2.

**Lemma 8.** *For every stamp $\mathsf{B}_a$ there exists an $i$ such that*

$$\boxed{\underline{a}}^a \subseteq \boxed{0}^i.$$

*Proof.* See appendix B. $\square$

8

We are now able to present a bounded instantiation for the model. Let $\mathcal{S}_N$ be a totally ordered set with $N^2$ elements (the total order is here only to avoid making non-deterministic choices). We define:

$$
\begin{aligned}
\mathbf{0} &= \sqcap \mathcal{S}_N, \\
\text{new}(\boxed{a}) &= \sqcap \{x \mid x \in \mathcal{S}_N \text{ and } x \notin \boxed{a}\}.
\end{aligned}
$$

Lemma 8 guaranties that new($\boxed{0}$) satisfies the axiom. It follows then that it acts as an encoding of counter mode model (Proposition 7). Thus we have constructed a bounded model for the data causality problem in a slice, which generalizes, by concatenating slices, to the full data causality problem addressed by version vectors.

# 6 Related Work

On what concerns bounded replacements for version vectors there is, up to our knowledge, no previous solution to the problem. The possible existence of a bounded substitute to version vectors was referred in [1] while introducing the version stamps concept. Version stamps allow the characterization of data causality in settings where version vectors cannot operate, namely when replicas can be created and terminated autonomously.

There have been several approaches to version vector compression. Update coalescing [14] takes advantage of the fact that several consecutive updates issued in isolation in a single replica can be made equivalent to a single large update. Update coalescing is intrinsic in bounded stamps since sequence restriction in the update operation discards non-propagated symbols. Dynamic compression [14] can effectively reduce the size of version vectors by removing a common minimum from all entries (along each slice). However, this technique requires distributed consensus on all replicas and therefore cannot progress if one or more replicas are unreachable. Unilateral version vector pruning [16] avoids distributed consensus by allowing unilateral deletion of inactive version vectors entries, but relays on some timing assumptions on the physical-clock's skew.

Lightweight version vectors [8] develop an integer encoding technique that allows a gradual increase of integer storage as counters increase. This technique is used in conjunction with update coalescing to provide a dynamic size representation. Hash histories [9] track data causality by collecting hash fingerprints of contents. This representation is independent of the number of replicas but grows in proportion to the number of updates.

The minimality of vectors clocks as a characterization of Lamport causality [11], presented by Charron-Bost [3] and recently re-addressed in [6], indicates particular runs where the full expressiveness of vectors clocks is required. However there are cases in which smaller representations can operate: Plausible Clocks [19] offer a bounded substitute to vectors clocks that are accurate in a large percentage of situations and may be used in settings were deviations only impacts performance and not correctness; Resettable Vector Clocks [2] allow a bounded implementation of vector clocks under a specific communication pattern between processes.

The collection of cached copies of the knowledge in other replicas has been explored before in [5, 20] and used for optimization of message passing strategies. This concept is sometimes referred to as matrix clocks [15]. These clocks are based on integer counters and are similar to our intermediate "counter mode principal vector" representation.

# 7 Conclusions

Version vectors are the key mechanism in the detection of inconsistency and obsolescence among optimistically replicated data. This mechanism has been used extensively in the design of distributed file systems [10, 7], in particular for data causality tracking among file copies. It is well known that version vectors are unbounded due to their use of counters; some approaches in the literature have tried to address this problem.

We have brought the attention to the fact that causally ordering a limited number of replicas does not require the full expressive power of version vectors. Due to the limited number of configurations among replicas, data causality tracking does not necessarily imply the use of unbounded mechanisms. As a consequence, Charron-Bost's minimality of vector clocks cannot be transposed to version vectors.

We have noted that to find a bounded alternative to version vectors, it was enough to concentrate on a sub-problem: keeping distributed knowledge about a total order generated by a single entity.

The key to bounded stamps was defining an intermediate unbounded mechanism and showing that it was possible to perform comparisons without requiring a global total order; this was the bulk of the proof correctness; bounded stamps were then derived as an encoding into a finite set of symbols. This required the definition of a non-trivial symbol

9

reuse mechanism that is able to progress even if an arbitrary number of replicas ceases to participate in the exchanges. This mechanism may have a broader applicability beyond its current use (e.g. log dissemination and pruning) and become a building block in other mechanisms for distributed systems.

The construction of the mechanism was supported by a simulator[1], which was used in the proof of correctness so as to probe (and discard) tentative hypotheses. The simulator was also turned into a model checker which proved the correctness up to $N = 4$, giving some confidence before the full proof of correctness was attempted.

Bounded version vectors are obtained by substituting integer counters on version vectors by bounded stamps. It represents the first bounded mechanism for detection of obsolescence and mutual inconsistency in distributed systems.

# References

[1] Paulo Sérgio Almeida, Carlos Baquero, and Victor Fonte. Version stamps – decentralized version vectors. In *Proceedings of the 22nd International Conference on Distributed Computing Systems (ICDCS)*, pages 544–551. IEEE Computer Society, 2002.

[2] A. Arora, S. S .Kulkarni, and M. Demirbas. Resettable vector clocks. In *19th Symposium on Principles of Distributed Computing (PODC'2000), Portland, 2000*. ACM, 2000.

[3] Bernadette Charron-Bost. Concerning the size of logical clocks in distributed systems. *Information Processing Letters*, 39:11–16, 1991.

[4] Colin Fidge. Timestamps in message-passing systems that preserve the partial ordering. In *11th Australian Computer Science Conference*, pages 55–66, 1989.

[5] Michael J. Fischer and A. Michael. Sacrificing serializability to attain high availability of data. In *Proceedings of the ACM Symposium on Principles of Database Systems*, pages 70–75. ACM, 1982.

[6] V. K. Garg and C. Skawratananond. String realizers of posets with applications to distributed computing. In *Proceedings of the ACM Symposium on Principles of Distributed Computing (PODC'01)*, pages 72–80. ACM, 2001.

[7] Richard G. Guy, John S. Heidemann, Wai Mak, Thomas W. Page, Gerald J. Popek, and Dieter Rothmeier. Implementation of the ficus replicated file system. In *USENIX Conference Proceedings*, pages 63–71. USENIX, June 1990.

[8] Yun-Wu Huang and Philip Yu. Lightweight version vectors for pervasive computing devices. In *Proceedings of the 2000 International Workshops on Parallel Processing*, pages 43–48. IEEE Computer Society, 2000.

[9] Brent ByungHoon Kang, Robert Wilensky, and John Kubiatowicz. The hash history approach for reconciling mutual inconsistency. In *Proceedings of the 23nd International Conference on Distributed Computing Systems (ICDCS)*, pages 670–677. IEEE Computer Society, 2003.

[10] James Kistler and M. Satyanarayanan. Disconnected operation in the coda file system. *ACM Transaction on Computer Systems*, 10(1):3–25, February 1992.

[11] Leslie Lamport. Time, clocks and the ordering of events in a distributed system. *Communications of the ACM*, 21(7):558–565, July 1978.

[12] Friedemann Mattern. Virtual time and global clocks in distributed systems. In *Workshop on Parallel and Distributed Algorithms*, pages 215–226, 1989.

[13] D. Stott Parker, Gerald Popek, Gerard Rudisin, Allen Stoughton, Bruce Walker, Evelyn Walton, Johanna Chow, David Edwards, Stephen Kiser, and Charles Kline. Detection of mutual inconsistency in distributed systems. *Transactions on Software Engineering*, 9(3):240–246, 1983.

[14] David Howard Ratner. *Roam: A Scalable Replication System for Mobile and Distributed Computing*. PhD thesis, 1998. UCLA-CSD-970044.

[15] Frédéric Ruget. Cheaper matrix clocks. In *Proceedings of the 8th International Workshop on Distributed Algorithms*, pages 355–369. Springer Verlag, LNCS, 1994.

[16] Yasushi Saito. Unilateral version vector pruning using loosely synchronized clocks. Technical Report HPL-2002-51, HP Labs, 2002.

[17] Yasushi Saito and Marc Shapiro. Optimistic replication. Technical Report MSR-TR-2003-60, Microsoft Research, 2003.

[18] R. Schwarz and F. Mattern. Detecting causal relationships in distributed computations: In search of the holy grail. *Distributed Computing*, 3(7):149–174, 1994.

[19] FJ Torres-Rojas and M. Ahamad. Plausible clocks: constant size logical clocks for distributed systems. *Distributed Computing*, 12(4):179–196, 1999.

[20] G. T. J. Wuu and A. J. Bernstein. Efficient solutions to the replicated log and dictionary problems. In *Proceedings of the ACM Symposium on Principles of Distributed Computing (PODC'84)*, pages 232–242. ACM, 1984.

---

[1]http://gsd.di.uminho.pt/bvv/bvv-simulator.py

# A    Proof of Lemma 2

The hypothesis of lemma 2 concern two stamps (say a and b) in which we can identify some sort of conflict between each stamp knowledge: Stamp b has a better knowledge concerning the primary state ($a^a \leq b^b$) but has an outdated vision concerning some other stamp (say k), i.e. $b^k \leq a^k$. Lemma 2 states that when this happens stamp b already attributes the value of $a^k$ to some other stamp (say $l$ — that is, $b^l = a^k$). In order to prove this result, it will be necessary to reinforce this statement: not only $b^l = a^k$ but it is possible to identify a flow of information between stamp l and k. Moreover, this flow of information (a sequence of synchronization operations starting from l to k) can be traced in stamp b's local state as a sequence of indexes enjoying some properties. These sequence of indexes are called *delay paths* and are defined as follows.

**9 Definition (Delay Path).** A *delay path* between $a^k$ and b is a non-empty sequence of indexes $\langle i_0, \ldots, i_n \rangle$ such that, for any stamp c,

1. $i_0 = k$,

2. $b^{i_n} = a^k$,

3. $b^{i_p} < a^k$          for all $0 \leq p < n$,

4. $a^k \leq i_p{}^{i_{p-1}}$       for all $0 < p \leq n$,

5. $a^k < c^{i_p} \Rightarrow a^k \leq c^{i_{p-1}}$     for all $0 < p \leq n$.

Some simple facts concerning delay paths.

**Proposition 10.** *Let $\langle i_0, \ldots, i_n \rangle$ be a delay path between $a^k$ and b. The following facts hold:*

1. $b^k \leq a^k$,

2. $a^k \in b$,

3. $a^k \leq i_p{}^{i_p}$          *for all $0 \leq p \leq n$,*

4. $b \in \langle i_0, \ldots, i_n \rangle \Rightarrow n = 0$.

*Proof.* The first three facts are immediate consequences from the definition and Proposition 1. Regarding the last fact, if $b$ occurred in a position $i_p$, being $n > 0$, by condition (4) of delay paths we have $a^k \leq b^{i_{p-1}}$; but this contradicts condition (3). Thus, $b$ only occurs in a singleton delay path. $\square$

Some of the conditions on delay paths impose global constrains on them that will allow to reason about global state changes and their impact on the local states. The following Lemma exposes the use of such global constrains.

**Lemma 11 (Pointwise-join Lemma).** *Let $\langle i_0, \ldots, i_n \rangle$ be a non-empty sequence of indexes. If for some $x$,*

1. $b^{i_n} = x$,

2. *for all $0 \leq p < n$, $b^{i_p} < x$,*

3. *for all $0 < p \leq n$ and any stamp c, if $x < c^{i_p}$ then $x \leq c^{i_{p-1}}$.*

*Then, for any stamp d for which $d^{i_0} \leq x$, there exists $0 \leq q \leq n$ such that $b^{i_q} \sqcup d^{i_q} = x$ and, for all $0 \leq p < q$, $b^{i_p} \sqcup d^{i_p} < x$.*

*Proof.* By induction on the length of the sequence $\langle i_0, \ldots, i_n \rangle$. For the base case (singular sequence) we have that $b^{i_0} = x$. Since $d^{i_0} \leq x$ we have $b^{i_0} \sqcup d^{i_0} = x$ and the remaining condition is vacuous. For the induction step, we consider the following cases: If $d^{i_0} = x$ then we set $q = 0$ since $b^{i_0} \sqcup d^{i_0} = x$. Otherwise, we know that $b^{i_0} \sqcup d^{i_0} < x$ and, by (4), that $d^{i_1} \leq x$. So we apply the induction hypothesis to the sequence $\langle i_1, \ldots, i_n \rangle$ and set $q$ to the resulting index plus 1. $\square$

We now show that the conditions in Lemma 2 are sufficient to establish the existence of delay paths.

**Lemma 12.** *If a and b are two stamps and $k$ a position such that*
$$a^a \leq b^b \quad \text{and} \quad b^k \leq a^k,$$
*then there exists a delay path between $a^k$ and b.*

*Proof.* We prove by induction on the length of the trace. If the last operation was l we use the singleton sequence $\langle k \rangle$ for the delay path and the conditions hold trivially. If the last operation was $U^0$ consider the following cases:

$a = b = 0$: we pick the sequence $\langle k \rangle$ that satisfies trivially all conditions;

$a = 0 \neq b$: after the update $a^a \not\leq b^b$, which contradicts the hypothesis;

$b = 0 \neq a$: if $k = 0$ then $b^k \not\leq a^k$, which contradicts the hypothesis. If $k \neq 0$ we use the same delay path that comes from the IH, which is still valid after the update because it does not contain position 0, since $b = 0 \neq k$ (Proposition 10).

11

$0 \notin \{a, b\}$**:** we use the same delay path from the IH, which is still valid: (1,2,3) because $a$ and $b$ are not affected by the update; (4) because only $0^0$ changes; (5) because even if for some $p$ we have $i_p = 0$, if $a^k < 0^0$, then $a^k \leq 0^{i_{p-1}}$ due to (4).

If the last operation was $\mathsf{S}^{xy}$ (and lets assume, without loss of generality, that $y$ is the more up-to-date stamp, i.e. $x^x \leq y^y$) we need to distinguish the following cases:

$\{x, y\} \cap \{a, b\} = \emptyset$**:** we use the same delay path from the IH, which is still valid: (1,2,3) because $a$ and $b$ are not affected; (4) because $i_p^{i_{p-1}}$ can only increase; (5) because for every $c = x \sqcup y$, if $a^k < c^{i_p}$, then either $c^{i_p}$ is computed pointwise and $a^k \leq c^{i_{p-1}}$ follows from the IH, or $i_p$ is either $x$ or $y$ and (by 4) $a^k \leq i_p^{i_{p-1}} \leq c^{i_{p-1}}$.

$\{x, y\} = \{a, b\}$**:** stamps $a$ and $b$ become equal after the synchronization and we pick the sequence $\langle k \rangle$ for the delay path;

$\{x, y\} \cap \{a, b\} = \{a\} \neq \{b\}$**:** in this case the stamp $a$ results from the synchronization of $x$ and $y$ and we have $x^x \leq y^y = a^a \leq b^b$. Consider the following two cases:

When $k = x$ and $x^x < y^y = a^k$. First, given that $y^y \leq b^b$ and $b^y \leq y^y$, we can apply the IH to $y$ and $b$ on index $y$ and establish the existence of a delay path $\langle i_0 = y, \ldots, i_n \rangle$ for $y^y$ in $b$. Then we prefix it by $k$, obtaining $\langle k, y, \ldots, i_n \rangle$, which is a suitable delay path between $a^k$ and $b$, given that: (1) holds by construction, (2) from the IH, (3) from the IH and $b^k < a^k$ (since $b^k \leq x^x < y^y = a^k$); (4) from the IH and $a^k = y^y = y^x = y^k$; (5) from the IH and because for every stamp $c$, $c^y \leq y^y = a^a = a^k$.

Otherwise, then either $a^k = x^k$ or $a^k = y^k$; applying the IH to either $x^k$ or $y^k$ and $b$ in position $k$ gives us a valid delay path for the resulting configuration (all conditions hold, including (5) as shown for the case $\{x, y\} \cap \{a, b\} = \emptyset$).

$\{x, y\} \cap \{a, b\} = \{b\} \neq \{a\}$**:** in this case the stamp $b$ results from the synchronization of $x$ and $y$.

When $k$ is either $x$ or $y$, we have $b^k = b^b = y^y$; but this means (as $a^a \leq b^b$ and $b^k \leq a^k$) that $a^k = b^k$; therefore $\langle k \rangle$ is a delay path.

Otherwise, $b^k = x^k \sqcup y^k$; this means that $y^k \leq b^k \leq a^k$ and by the IH there exists a delay path $P$ between $a^k$ and $y$. Given that also $x^k \leq a^k$, Lemma 11 establishes the existence of a sequence $Q = \langle i_0, \ldots, i_q \rangle$ (prefix of $P$) that is a delay path between $a^k$ and $b$ for the following reasons. Positions $x$ and $y$ do not appear in $Q$ — $x, y \neq i_0$ because we are assuming $k \neq x, y$, and $x, y \neq i_p$ for $p > 0$, otherwise we would have $a^k \leq x^{i_{p-1}}, y^{i_{p-1}}$ (condition (4) of delay paths of which $Q$ is a prefix) and then $a^k \leq x^{i_{p-1}} \sqcup y^{i_{p-1}}$, which contradicts Lemma 11. Thus, all elements $b^j$, with $j \in Q$ are computed pointwise (i.e $b^j = x^j \sqcup y^j$), making conditions (2,3 and 5) immediate consequences of Lemma 11. Condition (1) is trivially observed ($Q$ is a prefix of $P$); and condition (4) from the IH and because upon a join values can only increase.

$\square$

We can finally state Lemma 2.

**Lemma (2).** *For every stamp $\mathsf{C}_a$ and $\mathsf{C}_b$, and every index $k$,*

$$a^a \leq b^b \text{ and } b^k \leq a^k \quad \text{implies} \quad a^k \in b.$$

*Proof.* Direct from Lemma 12. $\square$

# B  Proof of Lemma 8

Lemma 8 says that each principal order is already contained in some cached order on the primary. Note that Lemma 2 already states that every principal element $a^a$ belongs to the primary principal vector, and delay paths were used to show where it can be found. Now, we will show that it is precisely in the primary cached order located in the position pointed out by the delay path between $a^a$ and $0$ that we can find all the elements in $\boxed{a}^a$. To prove this we need to reason about cached orders along delay paths. This suggests an extension of these to what we call *principal delay paths*.

**13 Definition.** A *principal delay path* for stamp $\mathsf{B}_a$ is a delay path $\langle i_0, \ldots, i_n \rangle$ between $a^a$ and $0$ that additionally satisfies the following condition: for every $0 \leq p \leq n$ and any stamp $\mathsf{B}_c$,

$$a^a = c^{i_p} \quad \text{implies} \quad \boxed{a}^a \subseteq \boxed{c}^{i_p} \text{ or} \\ (p > 0 \text{ and } a^a \leq c^{i_{p-1}}).$$

We now prove the existence of principal delay paths by extending the proof of existence in Lemma 12. Here we only go through the cases that are relevant for the additional condition.

**Lemma 14.** *For every stamp* $\mathsf{B}_a$ *there exists a principal delay path.*

*Proof.* (Sketch)

Consider the following additional arguments to the proof of Lemma 12. If the last operation was $\mathsf{S}^{xy}$ (assume $\mathsf{x}^x \leq \mathsf{y}^y$):

$\{x, y\} \cap \{a, 0\} = \emptyset$: let $\mathsf{c} = \mathsf{x} \sqcup \mathsf{y}$. If $i_p$ is either $x$ or $y$, we know that $p > 0$ (since $a \notin \{x, y\}$). Let $\mathsf{c}^{i_p} = \mathsf{x}^x \sqcup \mathsf{y}^y = \mathsf{y}^y$. When $\mathsf{a}^a = \mathsf{c}^{i_p}$, by condition (4), we have $\mathsf{a}^a \leq \mathsf{x}^{i_{p-1}}$ or $\mathsf{a}^a \leq \mathsf{y}^{i_{p-1}}$ which determines that $\mathsf{a}^a \leq \mathsf{c}^{i_{p-1}}$. When $\mathsf{c}^{i_p}$ is computed pointwise, the new condition follows by the induction hypothesis.

$\{x, y\} \cap \{a, 0\} = \{a\} \neq \{0\}$: when $a = x$ and $\mathsf{x}^x < \mathsf{y}^y$, let $\langle i_0 = y, \ldots, i_n \rangle$ be the principal delay path for $\mathsf{y}$. The new condition if verified for $\langle a, y, \ldots, i_n \rangle$ since, the case $c \neq y$ is trivial (because $\mathsf{c}^a = \mathsf{c}^x \leq \mathsf{x}^x < \mathsf{y}^y = \mathsf{a}^a$). For $c = y$, the new condition is satisfied since $\boxed{a}^a \subseteq \boxed{y}^y$ (Proposition 5).

$\{x, y\} \cap \{a, 0\} = \{0\} \neq \{a\}$ in this case the primary results from the synchronization of $\mathsf{x}$ and $\mathsf{y}$ (i.e. $\mathsf{y}$ is the primary before synchronization). Since $x \neq a$, then $0^a$ is computed pointwise. By IH we get a principal delay path $P$ to which we apply Lemma 11 to get a new sequence $Q$ where $x$ and $y$ never occur (c.f. proof of Lemma 12). The new condition follows by the induction hypothesis.

$\square$

**Lemma (8).** *For every stamp* $\mathsf{B}_a$ *there exists a position* $i$ *such that*
$$\boxed{a}^a \subseteq \boxed{0}^i.$$

*Proof.* Let $\langle i_0, \ldots, i_n \rangle$ be the principal delay path for $\mathsf{a}$ (given by Lemma 14). Instantiating the new condition for $0$ on $i_n$ we get that
$$\boxed{a}^a \subseteq \boxed{0}^{i_n}.$$

$\square$

13