

DSON: JSON CRDT Using Delta-Mutations For Document Stores

Arik Rinberg*
Technion
Haifa, Israel
ArikRinberg@campus.technion.ac.il

Tomer Solomon
IBM
Tel-Aviv, Israel
tomersolomon@ibm.com

Roe Shlomo
IBM
Tel-Aviv, Israel
Roe.Shlomo@ibm.com

Guy Khazma
IBM
Tel-Aviv, Israel
Guy.Khazma@ibm.com

Gal Lushi
IBM
Tel-Aviv, Israel
gal.lushi@ibm.com

Idit Keidar
Technion
Haifa, Israel
idish@technion.ac.il

Paula Ta-Shma
IBM
Tel-Aviv, Israel
paula@il.ibm.com

ABSTRACT

We propose DSON, a space efficient δ -based CRDT approach for distributed JSON document stores, enabling high availability at a global scale, while providing **strong** eventual consistency guarantees. We define the semantics of our CRDT based approach formally, and prove its correctness and convergence. Previous approaches optimize for collaborative document editing and store metadata proportional to the number of updates to a document, which is not acceptable for long lived document management. The metadata stored with our approach is bounded by $O(k^2D + n \log n)$, where n is the number of replicas, D is the number of document elements, and $k \leq n$ is the number of concurrent document updates. We also implement our approach[8] and demonstrate its space efficiency empirically. Experimental analysis shows that the metadata stored is typically significantly less than the worst case. This provides the basis for robust highly available distributed document stores with well defined semantics and safety guarantees, relieving application developers from the burden of conflict resolution.

PVLDB Reference Format:

Arik Rinberg¹, Tomer Solomon, Roe Shlomo, Guy Khazma, Gal Lushi, Idit Keidar, and Paula Ta-Shma. DSON: JSON CRDT Using Delta-Mutations For Document Stores. PVLDB, 14(1): XXX-XXX, 2020.
doi:XX.XX/XXX.XX

PVLDB Artifact Availability:

The source code and benchmarks have been made available at <https://github.com/crdt-ibm-research/json-delta-crdt>.

¹This work was done during an internship at IBM Research.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 14, No. 1 ISSN 2150-8097.
doi:XX.XX/XXX.XX

1 INTRODUCTION

NoSQL systems, and in particular document stores, have significantly grown in popularity over the last decade [7], and those in most widespread use adopt a JSON-based data model [6]. Such a data model maps directly to common programming language data structures and relaxes schema restrictions, making it suitable for developers' needs. Applications are diverse but include managing metadata for digital media such as photos and videos, social media posts, product information, and machine-generated data/logs. Distributed document stores are often applied in mobile and edge scenarios where both global scale and high availability under disconnected operation are paramount.

Distributed document stores often default to eventual consistency to enable high availability for both reads and writes on a global scale e.g. CouchDB [5], Couchbase [3], DynamoDB [9]. However, eventual consistency only guarantees that a document's state will converge at some unspecified future point in time if there are no additional updates. Reasoning without substantial safety guarantees adds a burden to application developers and is not suitable for all applications. Conflict resolution policies are often arbitrary [19] and have complex, non-intuitive semantics, e.g., Couchbase offers conflict resolution where the replica with the most mutations wins [4]. To achieve more control over consistency, some document stores provide an API for application developers to manually specify how to merge conflicts when they arise [19]. This approach is time consuming, error prone, and difficult to reason about.

We adopt *Conflict-free Replicated Data Types (CRDTs)* [42] to achieve a distributed document store design with **strong** eventual consistency, meaning that any two nodes that have received the same set of updates will be in the same state. In addition we also provide causal consistency and read-your-writes guarantees. The semantics we propose is well defined and intuitive, enabling programmers to reason about application correctness. We leverage and extend existing theoretical results for CRDTs and apply them to the document store context, where documents can have arbitrarily long lifetimes, but where the amount of metadata used to support CRDTs must be reasonably bounded.

The CRDT approach relies on predetermined conflict resolution rules which dictate their semantics. These rules are typically data structure specific, and CRDTs have been defined for counters, sets, directed graphs etc. [42]. In this paper, we focus on JSON document stores, and treat JSON[12] as a composition of nested maps and arrays (lists) over base data types. Note that since we handle both maps and arrays our work is applicable to other kinds of document stores as well. We define the semantics of our approach formally and prove its convergence and correctness. Existing CRDT approaches for JSON [2, 15] address collaborative document editing and collect metadata proportional to the number of updates. This metadata cannot easily be bounded or cleaned up over time, and therefore these approaches are not suited to our context. Our work enables global high availability for distributed JSON document stores with reasonable overheads and proven semantics.

CRDTs are often split into two main families: op-based [20] and state-based [32]. Op-based approaches are simpler and require smaller messages, but they assume reliable, exactly-once ordered messaging which is hard to maintain even when TCP is used [25]. On the other hand messages in state-based systems include the entire state, which can be prohibitively large. Almeida et al. introduced δ -state-based CRDTs [18] to achieve the best of both worlds – a δ -based CRDT sends a delta instead of the entire state, reducing communication overheads. We chose this approach because it naturally enables bounding metadata growth. However, it requires theoretical guarantees of causal consistency and designing causal δ -based CRDTs is non trivial. The original paper defines several causal δ -CRDTs including multi-value registers, add-wins and remove-wins sets and an enable-wins flag. It also defines a composable (i.e. arbitrarily nestable) Observed-Remove Map which can contain other causal δ -based CRDTs including maps themselves. However, maps on their own are not sufficient to handle JSON[12].

In this paper we define a composable causal δ -based Array CRDT. In addition to arrays (lists) being an important datatype in their own right, this results in a causal δ -based CRDT library powerful enough to capture JSON data. To date various array CRDTs have been proposed, such as Logoot [44], Treedoc [37], WOOT [36] and others [33, 40], as well as a technique to support concurrent move operations on array elements [27]. However none of these achieve composable, causal δ -based array CRDTs. In this paper we define causal composable δ -based array CRDTs supporting arbitrary levels of nesting and concurrent move operations without data duplication.

We base our approach on that of Almeida et. al. [18] and adopt *Observed-Remove (OR)* semantics for Arrays and combine them with their OR-Maps. OR semantics are intuitive because only observed state can be removed. For example, under OR semantics, if there is an update concurrent to a remove, the update wins and the element isn't removed. A shopping cart example is given in Figure 1, where OR semantics is contrasted with Remove-Wins (RW) semantics.

Our main contributions are as follows. We define a δ -based CRDT for Observed-Remove JSON documents which is built on a novel composable δ -based array CRDT. While ORMap semantics was previously defined [42], we formalize the δ -based approach with respect to its execution. We provide a formal definition of the array CRDT, and prove both its convergence and the correctness of its semantics. We show that the amount of metadata stored is bounded

by $O(k^2D + n \log n)$, where n is the number of replicas, D is the number of document elements, and $k \leq n$ is the number of concurrent operations. We implemented a prototype in JavaScript[8] and verify that metadata stays constant with the number of updates, unlike the case for existing JSON CRDTs [2, 15]. Some initial results were presented as an extended abstract[39].

The rest of the paper is structured as follows. Section 2 covers related work, section 3 covers preliminaries, and Section 4 covers our core theoretical results. Section 5 presents the algorithm for the novel Observed-Remove Array CRDT, and Section 6 presents the JSON CRDT. Section 7 surveys our implementation and Section 8 provides an experimental evaluation of metadata growth. Section 9 presents some ideas for future work and Section 10 concludes.

2 RELATED WORK

As described in the introduction, there is no known JSON CRDT that is implemented using delta mutations. The key missing piece is a delta-based array implementation that supports moving of elements and arbitrary nesting, without using tombstones. However, state-based and operation-based constructions do exist. In this section we describe existing CRDTs for arrays and JSON documents as well as describing usage of CRDTs in datastores.

The problem of collaborative text editing has motivated many proposals for array CRDTs such as WOOT [35], Treedoc [38], Logoot [44], and RGA [41]. However, text editing deals with a linear structure and does not consider arbitrary nesting of CRDTs as required for sequences in JSON documents. Higher level constructions, such as the ORArray, are built upon such base constructions.

Moreover, while all mentioned constructions support the insert and delete operations, only RGA supports update and none of them support moving elements. In Moving elements in list CRDTs A generic black box approach to adding move support to existing CRDTs has been proposed [27]. However, the described construction cannot easily be applied in a delta-based setting, e.g., it relies on AWSet that in the delta-based framework is only defined for primitive types and lacks support for an update operation.

A JSON CRDT has been proposed [28] and is implemented in the Automerge library [2]. Automerge is an operation-based CRDT that implements many different data types including composable maps and arrays. It provides a good solution for use cases that require observing all history and supporting undo operations. The array implementation in Automerge is based on RGA, an array CRDT that is implemented similar to a linked list. Each node contains an identifier, an identifier of the next item, a value, and a flag to indicate whether the node is a tombstone. While RGA was originally designed for linear data, nesting is enabled in Automerge by cleverly holding the identifiers of nested CRDTs in the value fields. However, the RGA construction inherently requires tombstones and has an ever-growing metadata problem. Garbage collection cannot be fully implemented without compromising consistency, leading to a significant drawback in use cases such as document stores. In contrast, ORArray and the JSON CRDT built on top of it are naturally implemented without tombstones. Another design choice in Automerge is the handling of conflicts that cannot be resolved automatically. When users concurrently update the same key in the same map (or the same index in the same list) then one of

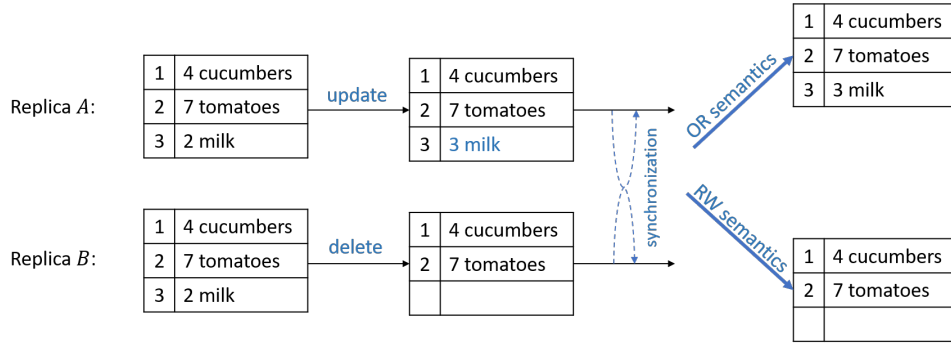


Figure 1: Observed-remove vs. remove-wins semantics via a shopping list example. With OR semantics, only observed state can be removed, therefore the delete at replica B cannot override the update at replica A which had not been observed.

the concurrently written values is picked arbitrarily and the others are stored in a conflicts structure as raw values (i.e., not CRDTs). In contrast, the JSON CRDT presented in this paper preserves all subtrees as valid CRDTs and gives more control to the user, similar to the semantics of a multi-value register, when desired.

Another JSON CRDT has been proposed in YATA [34] and is implemented in the Yjs library [15]. The core of the YATA CRDT is an Array CRDT implemented as a doubly linked list with insert and delete operation. YATA supports update operations and complex data types such as maps by storing *abstract operation types* as the array elements. For example, to support the update operation a *Replace Manager* operation type inserts the updated value to the beginning of an array such that eventually all replicas observe the same first array element as the updated value. While YATA is an operation based CRDT, delete operations are handled using a state based CRDT for storing tombstones. YATA describes a time based garbage collection mechanism and some ad-hoc optimizations are implemented to try minimizing the amount of state stored due to tombstones. The δ -based CRDTs are naturally built without the need for garbage collection, due to the use of a causal context.

δ -based CRDTs have been adopted in industry as part of Akka toolkit [1] for building applications using the actor model. Akka supports delta-based CRDTs for counters, sets, maps and registers but lacks support for arrays. We believe that the addition of a full JSON CRDT can both open the door for new use cases and greatly simplify application development.

δ -based CRDTs are also used to build real-time collaborative applications on top of IPFS [13, 21]. The array implementation is based on RGA with tombstones and without any support for nesting. That is, you can't store an array inside a map nor store any CRDT type inside the array.

Distributed databases such as Redis, Riak, Cosmos DB and Anti-dote DB include support for CRDT data types such as sets and maps. However, for data stores such as CouchDB that use JSON as the data model a full JSON CRDT is required. AutoCouch [24] leverages both the AutoMerge library and CouchDB to enable developing collaborative applications. AutoCouch synchronizes the entries of the AutoMerge history using the database's revision tree. As stated by the authors the main issue is the ever growing metadata size, both due to tombstones and due to the inherent history keeping in

AutoMerge. We believe that our approach provides an alternative that preserves all of the functionality required by document stores while keeping the synchronization cost low due to the use of delta mutations and a design that does not require garbage collection. Our approach may also be suitable for other use cases that use the JSON model, such as decentralized Kubernetes control planes [26, 31].

3 PRELIMINARIES

Section 3.1 presents our model, and in Section 3.2 we define our notations. In Section 3.3 we present background δ -based CRDTs.

3.1 Model

Our system consists of a set of *nodes* I hosting *replicas* of objects; we use node and replica interchangeably. The object exposes API *methods* which are used to alter the underlying state of the object. Nodes communicate with each other by *broadcasting* messages over unreliable, unordered links. The network provides the following two assumptions: (1) If a node delivers message m , then another node previously broadcast message m , and (2) if a message is broadcast infinitely many times it is eventually delivered by every node. We denote a broadcast of message m by node i as $bcast_i(m)$, and deliver of message m broadcast by node i at node j as $deliver_j(m, i)$. An *event* occurring at a node is either an API method invocation, a message broadcast, or a deliver.

An *execution* σ of the system is an interleaving sequence of events and states. A *trace* of an execution σ is the sub-sequence of events in σ . We denote the trace of an execution σ as $Tr(\sigma)$.

We make use of the *happens-before* relation as defined by Lamport [30]. An event $e_1 \in Tr(\sigma)$ happens-before or *precedes* an event $e_2 \in Tr(\sigma)$ in an execution σ , denoted $e_1 <_\sigma e_2$ if: (1) e_1 occurs before e_2 in the same replica in $Tr(\sigma)$, or (2) e_1 is a broadcast of message m and e_2 is a deliver of that message, or (3) there exists some event $e \in Tr(\sigma)$ such that $e_1 <_\sigma e$ and $e <_\sigma e_2$. We define the set of operations in σ as $Ops(\sigma) = \{e \in Tr(\sigma) | e \text{ is an API method invocation}\}$. We say that two events e_1 and e_2 are *concurrent* if neither precedes the other in σ , denoted $e_1 ||_\sigma e_2$.

Previous work defines CRDT object semantics in terms of the poset $(Ops(\sigma), <_\sigma)$ [32], i.e., the set of operations and the corresponding happens-before partial order. As an example, we present

the *Multi-Value Register (MVReg)*, a CRDT register. The register exposes the following API methods: (1) `WRITE(x)`, and (2) `READ()`. Given an execution σ , the `READ` method returns:

$$\{v \mid \exists \text{WRITE}(v) \in \text{Ops}(\sigma) \wedge \nexists \text{WRITE}(u) : \text{WRITE}(v) <_{\sigma} \text{WRITE}(u)\}.$$

For example, if node i executes `WRITE(v1)` and concurrently node j executes `WRITE(v2)`, a `READ` operation on node k that is preceded by both these `WRITE` operations returns $\{v_1, v_2\}$. If node k executes a `WRITE(v3)` after the `READ`, then an ensuing `READ` returns $\{v_3\}$.

3.2 Notations

We consider standard notations for sets and maps. The sets we consider in this paper are partially ordered (namely, posets), with a least element \perp . A *map* is a set of key-value pairs $\{k \mapsto v\}$, where each key k is from some set \mathcal{K} , and is associated with a single value v from some set \mathcal{V} . Given a map m and a key $k \in \mathcal{K}$, we denote its associated value by $m[k]$. The domain of a map m is denoted $\text{dom } m$ and the range of a map m as $\text{ran } m$, i.e., $\text{dom } m = \{k \mid \{k \mapsto v\} \in m\}$ and $\text{ran } m = \{v \mid \{k \mapsto v\} \in m\}$. For a map m and a key k , if $k \notin \text{dom } m$, we say that $m[k] = \perp$, i.e., all keys not present in the map are mapped to \perp . For a pair $p = (a, b)$, we use $\text{fst}(p)$ and $\text{snd}(p)$ to denote its first (a) and second (b), respectively. Finally, for a set S we denote its power set by $P(S)$.

A *join semilattice* [22], or simply a *lattice*, is a three tuple $(S, <, \sqcup)$, where S is a set, $<$ is a partial order on S , and \sqcup is a function $S \times S \mapsto S$ that is associative, commutative, and idempotent.

Notation	Explanation
\mathcal{I}	Set of nodes
$\text{bcst}_i(m)$	Broadcast event of message m by node $i \in \mathcal{I}$
$\text{deliver}_j(m, i)$	Deliver event of message m broadcast by node $i \in \mathcal{I}$, at node $j \in \mathcal{I}$
σ	An execution of the system, including events and states.
$\text{Tr}(\sigma)$	Sub-sequence of events (trace) of execution σ .
$e_1 <_{\sigma} e_2$	Event e_1 precedes event e_2 in execution σ
$e_1 \parallel_{\sigma} e_2$	Event e_1 and e_2 are concurrent in execution σ
$\text{Ops}(\sigma)$	The set of API methods invoked in σ .

Table 1: Table of notations.

3.3 δ -Based CRDTs

As mentioned in the introduction, δ -based CRDTs were proposed by Almeida et al. [18] to achieve the best of both state-based CRDTs and op-based CRDTs. This is achieved by using a join semilattice. The state of the CRDT and the broadcasted messages are both elements of the join semilattice. We denote the broadcasted messages as *deltas*.

The API methods use the current object state X to generate a delta X^{δ} that represents the incremental movement from the old state (before the method invocation) to the new one (after), i.e., X is the old state and $X \sqcup X^{\delta}$ is the new one. After generating X^{δ} , the node joins its state with X^{δ} to get the new state, and broadcasts X^{δ} to all nodes.

Nodes needn't keep track of which messages have been received, as the join is idempotent, and don't need to deliver messages in the same order, as the join is commutative and associative. Finally, convergence is achieved whenever all nodes deliver all messages from all other nodes. An *anti-entropy algorithm* [17] is used to enforce convergence and preserve causality: it retransmits outstanding deltas, enforces the causal order of delivery, and merges outstanding deltas in order to send fewer messages. The full algorithm is described in Appendix B.

Commonly, δ -based CRDTs assign unique identifiers to operations. The simplest way to generate these identifiers is for replica $i \in \mathcal{I}$ to generate the sequence of pairs $(i, 1), (i, 2), \dots$ and assign a pair per event. We denote such a pair a *dot*. An operation's dot is broadcast along with the delta reflecting its mutation. We say that a dot has been *observed* by replica i if: (1) i generated the dot, or (2) i delivered a message containing the dot. The collection of observed events at a node up to a given point in an execution is called a *causal context*, and it is maintained locally at every replica as part of its state. The causal context c local to i exposes three functions: $\text{max}_i(c)$, which gives the maximum sequence number in the causal context, $\text{next}_i(c)$, which produces the next dot for replica i , and $\text{insert}_i(c, d)$, which adds dot d to the causal context. The causal context is defined in Equation 1.

$$\begin{aligned} \text{CAUSALCONTEXT} &\in P(\mathcal{I} \times \mathbb{N}) \\ \text{max}_i(c) &= \max(\{n \mid (i, n) \in c\} \cup \{0\}) \\ \text{next}_i(c) &= (i, \text{max}_i(c) + 1) \\ \text{insert}_i(c, d) &= c \cup \{d\} \end{aligned} \quad (1)$$

As the causal context is a grow-only set, its size is unbounded. In practice, it can be efficiently compressed without loss of information. As dots are created sequentially, instead of maintaining a list of dots $(i, 1), (i, 2), \dots, (i, n)$ for every replica i , we can maintain a mapping *prefix* : $\mathcal{I} \mapsto \mathbb{N}$, such that if $\text{prefix}[i] = n$ then all dots (i, m) for $m \leq n$ are in the causal context. Due to the unreliable network, we may have a gap in the sequence, e.g., replica j may observe $(i, 1)$ and $(i, 3)$ but not $(i, 2)$. We therefore also maintain a separate (uncompressed) set of dots exceeding the longest prefix.

A `DOTSTORE` is a container for data-type specific information and is used to store the state of a δ -based CRDT. It provides the operation `dots` for querying the set of event identifiers (i.e., dots) currently stored in the dot store. We now present the dot stores defined by Almeida et al. [18]. The `DOTFUN` instantiated with lattice V is a map from dots in $\mathcal{I} \times \mathbb{N}$ to V . Its `dots` operation returns the domain of the map. It is presented formally in Equation 2.

$$\begin{aligned} \text{DOTFUN}(V : \text{Lattice}) : \mathcal{I} \times \mathbb{N} &\mapsto V \\ \text{dots}(s) &= \text{dom } s \end{aligned} \quad (2)$$

The `DOTMAP` instantiated with a set K and a `DOTSTORE` V is a map from K to V . Its `dots` operation returns the union of the return values of invoking the `dots` operation on each `DOTSTORE` in the range. Equation 3 presents it formally.

$$\begin{aligned} \text{DOTMAP}(K, V : \text{DotStore}) : K &\mapsto V \\ \text{dots}(m) &= \bigcup_{k \in \text{dom } m} \text{dots}(m[k]) \end{aligned} \quad (3)$$

To form the join semilattice that makes up the state and deltas of δ -based CRDTs, dot stores are paired with causal contexts. We denote the pair $\text{DOTSTORE} \times \text{CAUSALCONTEXT}$ as $\text{CAUSAL}\langle \text{DOTSTORE} \rangle$. For each type of $\text{CAUSAL}\langle \text{DOTSTORE} \rangle$ we define the join function (\sqcup). The semilattice's partial ordering is defined with respect to the join function: For any $X \in \text{CAUSAL}\langle \text{DOTSTORE} \rangle$, $X \sqcup \perp = X$. For two elements $X_1, X_2 \in \text{CAUSAL}\langle \text{DOTSTORE} \rangle$ we say that $X_1 < X_2$ iff $\exists X \neq \perp \in \text{CAUSAL}\langle \text{DOTSTORE} \rangle$ such that $X_1 \sqcup X = X_2$.

The merge in the $\text{CAUSAL}\langle \text{DOTFUN} \rangle$ keeps values that exist in both of the mappings and merges their respective values, or exist in either one of the mappings and are “new” to the other, in the sense that they are not in its causal history. It is formally defined in Equation 4.

$$(m, c) \sqcup (m', c') = \begin{aligned} & \{(d \mapsto m[d] \sqcup m'[d] \mid d \in \text{dom } m \cap \text{dom } m') \\ & \cup \{(d, v) \in m \mid d \notin c'\} \cup \{(d, v) \in m' \mid d \notin c\}, c \cup c'\} \end{aligned} \quad (4)$$

The merge in the $\text{CAUSAL}\langle \text{DOTMAP} \rangle$ applies the merge recursively on each of the keys in either domains, and keeps all none \perp values. It is presented formally in Equation 5. An example of a \perp value is a merge between two elements of the $\text{CAUSAL}\langle \text{DOTFUN} \rangle$ semilattice, where the domains are disjoint but all mappings are in the others causal history. Consider for example a write w_1 that precedes a write w_2 , i.e., $w_1 <_{\sigma} w_2$, then the dot generated by w_1 is in the causal context of the delta generated by w_2 . By the definition of join, the mapping doesn't “survive” the join, and therefore the old value (written by w_1) is overwritten – it isn't present in the range of the map after w_2 .

$$(m, c) \sqcup (m', c') = \begin{aligned} & \{(k \mapsto v(k) \mid k \in \text{dom } m \cup \text{dom } m') \\ & \wedge v(k) \neq \perp\}, c \cup c' \end{aligned} \quad (5) \\ \textbf{where } v(k) = \text{fst}((m(k), c) \sqcup (m'(k), c'))$$

We next give an example of a δ -based CRDT, via the MVReg [18]. Its pseudo-code is presented in Algorithm 1. The state is an element (m, c) in the semilattice $\text{CAUSAL}\langle \text{DOTFUN} \rangle$. Namely, a mapping from dots to values in V . The $\text{WRITE}(v)$ operation creates a delta $X^{\delta} = (m^{\delta}, c^{\delta}) \in \text{CAUSAL}\langle \text{DOTFUN} \rangle$ as follows: c^{δ} contains the current domain of the DOTFUN , which is the set of dots written but not over-written, and m^{δ} is a mapping from a new dot to the written value v . Note that any preceding $\text{WRITE}(u)$ that is contained m before the operation has its dot added to c^{δ} , therefore, when joining X^{δ} with the current state, u is removed from m . A READ operation returns all values that are not removed. For example, consider two replicas i and j of a new MVReg. Replica i executes $\text{WRITE}_i(v)$, generates delta $(\{(i, 1) \mapsto v\}, \{(i, 1)\})$, merges and broadcasts it. Replica j delivers the delta, and changes its state to $(\{(i, 1) \mapsto v\}, \{(i, 1)\})$. Now, replica i executes $\text{WRITE}_i(u)$ concurrently to replica j executing $\text{WRITE}_j(w)$. The delta generated by $\text{WRITE}_i(u)$ is $(\{(i, 2) \mapsto u\}, \{(i, 1), (i, 2)\})$, and the delta generated by $\text{WRITE}_j(w)$ is $(\{(j, 1) \mapsto w\}, \{(i, 1), (j, 1)\})$. After all deltas have been merged, both replicas are in the state:

$$(\{(i, 2) \mapsto u, (j, 1) \mapsto w\}, \{(i, 1), (i, 2), (j, 1)\}),$$

so the MVReg holds both values u and w .

Algorithm 1 MVReg algorithm for node i .

```

1: Global state  $(m, c)$  ▷ CAUSAL(DOTFUN)
2: procedure WRITE $_i(v)$ 
3:    $d \leftarrow \text{next}_i(c)$  ▷ New dot
4:    $c' \leftarrow \text{dom } m$  ▷ Currently stored dots
5:   return  $(\{d \mapsto v\}, c' \cup \{d\})$  ▷ New  $X^{\delta}$ 
6: procedure READ $_i$ 
7:   return  $\text{ran } m$ 

```

4 OBSERVED-REMOVE SEMANTICS

Observed-remove semantics are a commonly used in production systems (e.g., Riak [29]). Under such semantics, an operation can overwrite (“undo”) preceding operations executed on the same element. The MVReg is a simple example of such an object, where a WRITE overwrites any preceding WRITE operations albeit not concurrent ones: If two WRITE operations are concurrent, then both values are reflected in the state. We now present two more elaborate examples, and formalize their OR semantics.

4.1 δ -Based ORMap

We first formalize the *Observed Remove Map (ORMap)* [18]. It is a mapping of some set of keys K , each to an arbitrary δ -based CRDT, possibly an ORMap itself. Thus, the ORMAP is *composable*. The map has no schema, in that each key can be mapped to a different type of CRDT. Its API consists of the following methods:

- APPLY** (k, o_i^{δ}) – given a key $k \in K$ mapped to some δ -based CRDT v of type V , and a method o_i^{δ} from the API of V , The method applies o_i^{δ} to the CRDT by invoking o_i^{δ} on v .
- DELETE** (k) – deletes key k from the map.
- GET** (k) – returns the embedded value of the CRDT mapped to by k .

For example, consider a map from keys to MVRegs. Writing value v to key k_1 is done by executing $\text{APPLY}(k_1, \text{WRITE}(v))$ to the map, and $\text{GET}(k_1)$ reads the same MVReg.

Recall that the deltas generated by API methods and the replica states are both members of the same join semilattice. We next define the semantics of the δ -based ORMap. Given execution σ , define:

$$\begin{aligned} \mathcal{X}_{\text{apply}}(\sigma) &\triangleq \{X^{\delta} \mid \exists \text{APPLY}_i(k, o_i^{\delta}) \in \text{Ops}(\sigma) \text{ that returns } X^{\delta} \\ &\wedge \nexists \text{DELETE}_j(k) \in \text{Ops}(\sigma) : \text{APPLY}_i(k, o_i^{\delta}) <_{\sigma} \text{DELETE}_j(k)\} \\ X(\sigma) &\triangleq \begin{cases} \bigsqcup_{X^{\delta} \in \mathcal{X}_{\text{apply}}(\sigma)} X^{\delta} & \mathcal{X}_{\text{apply}}(\sigma) \neq \emptyset \\ \perp & \text{otherwise} \end{cases} \end{aligned}$$

Intuitively, $\mathcal{X}_{\text{apply}}$ is the set of deltas generated by APPLY operations whose keys were not subsequently deleted, and X is the join of these deltas.

Furthermore, assume the existence of an abstract function $\text{VALUE}(X, k)$, which retrieves the value of key k from any element X in the join semilattice, similarly to executing $\text{GET}(k)$ on a replica with state X . Then a $\text{GET}(k)$ returns $\text{VALUE}(X(\sigma), k)$.

Under this definition, in the presence of a concurrent APPLY and DELETE of the same key, the DELETE removes all APPLY operations that causally precede it and leaves the key in the map with the

concurrently written value. A DELETE can thus be seen as an “undo” of all operations leading up to it.

The semantics of the ORMap were previously defined less formally in [18]. Here we formally define the semantics of the return value of a GET for the first time.

Algorithm 2 presents the ORMap pseudo-code, as proposed in [18]. The state is an element in $\text{CAUSAL}(\text{DOTMAP})$. The causal context is shared by all keys and corresponding nested CRDTs. The APPLY method fetches the value at key k from m and pairs it with the shared causal context c to obtain the value from the embedded type. It then invokes the operation over the pair. From the resulting pair it extracts the value to create a new mapping for that key, which it pairs with the resulting causal context. As the DELETE method uses the dots method from within a DOTMAP, the dots returned are the dots of all embedded values, therefore deleting a key recursively removes the corresponding embedded values.

Algorithm 2 ORMap algorithm for node i .

```

1: Global state  $(m, c)$  ▷ CAUSAL(DOTMAP)
2: procedure APPLY $_i(k, o_i^\delta)$ 
3:    $(v, c') \leftarrow o_i^\delta(m[k], c)$  ▷ Apply the mutator.
4:   return  $(\{k \mapsto v\}, c')$ 
5: procedure DELETE $_i(k)$ 
6:    $c' \leftarrow \text{dots}(m[k])$  ▷ Get all embedded dots.
7:   return  $(\perp, c')$ 

```

4.2 δ -Based ORArray

We now introduce a new OR data type, the *Observed Remove Array* (ORArray). To this end, we represent an array as a set of value-position pairs, whose values are general CRDTs. Using an integer index as the position would require updating multiple element positions following an insertion or deletion. Therefore, we adopt the notion of *stable identifiers* [27], which allow insertion and deletion without updating other elements, for example, the set of real numbers. An insertion of an element between an element at position p_1 and one at p_2 , inserts an element at position $(p_1 + p_2)/2$. The array is the sequence of values sorted by the corresponding position in ascending order. For example, the array $[a, b, c]$ can be represented by the set $\{(a, 1.0), (b, 7.5), (c, 42.7)\}$.

As the ORArray allows concurrent operations, the position may be ambiguous, similarly to the value of an MVReg under concurrent writes. To this end, we extend the position to be a set of possible positions. For example, if element b is concurrently moved before a and after c , its position may be $\{0.3, 54\}$. To sort the array, replicas deterministically choose an element from the position set (e.g., the maximum) to use for ordering the values. To keep track of values as they change positions under concurrent operations, we identify each array element by a unique identifier. For example, the unique identifier can be the dot with which the element was created (as the creation point is unique). The representation of the ORArray is therefore a map from unique identifiers to value-positions pairs. For example, the aforementioned array can be represented as the map: $\{(i, 1) \mapsto (a, \{1.0\}), (j, 1) \mapsto (b, \{0.3, 54\}), (i, 2) \mapsto (c, \{42.7\})\}$, where $(i, 1)$ is the dot marking the creation even of element a , $(j, 1)$

is the creation event of b , and $(i, 2)$ is the creation event of c . When an element is inserted for the first time, the API receives the unique identifier that identifies the inserted element until its deletion.

The value of an element can be any δ -based CRDT, possibly an ORArray itself. Thus, the ORArray is also composable. The array has no schema; each value can be a different type of CRDT.

Its API consists of the following methods:

APPLY (uid, o_i^δ, p) – given a unique identifier uid mapped to some pair consisting of a δ -based CRDT v of type V and a set of positions, a method o_i^δ from the API of V , and a target position p , The method applies o_i^δ to the CRDT by invoking o_i^δ on v , and overwriting the set of positions with p .

MOVE (uid, p) – given a unique identifier uid and a position p , overwrites the position of the element to be p .

DELETE (uid) – deletes any element mapped to uid from the map.

GET (uid) – given a uid mapped to a value v of some CRDT of type V and a set of positions P , returns the value of v based on V paired with the set of positions P .

While the API uses a unique id to refer to an element in an array and uses stable position identifiers to sort it, traditional use of an array consists of using integer indices, i.e., the first element of an array arr is $arr[0]$. To this end, we also provide a wrapper API and implemented functions converting an integer index to a unique id and to a stable position identifier. Converting an index idx to a unique identifier uid is done by sorting the array using the position identifiers, and returning the uid of the element at index idx . Converting an index to a position identifier is done in a similar fashion. The wrapper API consists of the following methods, mapped to the methods above:

INSERT (idx, o_i^δ) – given an index idx and a method o_i^δ from the API of some CRDT of type V , The method assigns a unique id uid , assigns a stable position identifier p such that the new element in the sorted array appears at index idx , and invokes **APPLY** (uid, o_i^δ, p) .

UPDATE (idx, o_i^δ) – given an index idx and a method o_i^δ of some CRDT type V , The method finds the uid corresponding to the element at index idx , finds the position p , and invokes **APPLY** (uid, o_i^δ, p) .

MOVE (old_idx, new_idx) – given two indexes, finds the element uid corresponding to the element at index old_idx , calculates the stable position identifier p such that the element in the sorted array will be at index new_idx , and invokes **MOVE** (uid, p) .

DELETE (idx) – given an index idx , finds the element uid corresponding to the element at index idx , and invokes **DELETE** (uid) .

GET (idx) – given an index idx , finds the element uid corresponding to the element at index idx , and invokes **GET** (uid) .

We define the semantics of the ORArray in two steps: First the value and next the position. We define the semantics of the value

in a similar fashion to Section 4.1. Given execution an σ , define:

$$\begin{aligned} \mathcal{X}_{\text{apply}}(\sigma) &\triangleq \{X^\delta \mid \exists \text{APPLY}_i(\text{uid}, o_i^\delta, p) \in \text{Ops}(\sigma) \text{ that returns } X^\delta \\ &\wedge \nexists \text{DELETE}_j(\text{uid}) \in \text{Ops}(\sigma) : \text{APPLY}_i(\text{uid}, o_i^\delta, p) <_\sigma \text{DELETE}_j(\text{uid})\} \\ X(\sigma) &\triangleq \begin{cases} \bigsqcup_{X^\delta \in \mathcal{X}_{\text{apply}}(\sigma)} X^\delta & \mathcal{X}_{\text{apply}}(\sigma) \neq \emptyset \\ \perp & \text{otherwise} \end{cases} \end{aligned}$$

Furthermore, assume the existence of an abstract function which gets the value of element uid from any element X in the join semilattice, similarly to executing $\text{GET}(\text{uid})$ on a replica with state X , denoted $\text{VALUE}(X, \text{uid})$. Then $\text{fst}(\text{GET}(\text{uid}))$ returns $\text{VALUE}(X, \text{uid})$.

To define the positional semantics, we first identify the subset of $\text{Ops}(\sigma)$ consisting of all operations that alter element uid (i.e., $\text{APPLY}(\text{uid}, o_i^\delta, p)$, $\text{MOVE}(\text{uid}, p)$, and $\text{DELETE}(\text{uid})$), and do not precede another operation that alters element uid in σ . We denote this subset as $\text{maximal}(\sigma, \text{uid})$.

Next, given an execution σ , define:

$$\begin{aligned} P_{\text{apply}}(\text{uid}, \sigma) &\triangleq \{p \mid \exists \text{APPLY}(\text{uid}, o_i^\delta, p) \in \text{maximal}(\sigma, \text{uid})\} \\ P_{\text{move}}(\text{uid}, \sigma) &\triangleq \{p \mid \exists \text{MOVE}(\text{uid}, p) \in \text{maximal}(\sigma, \text{uid})\} \end{aligned}$$

At the end of σ , if $\text{fst}(\text{GET}(\text{uid})) \neq \perp$ then $\text{scnd}(\text{GET}(\text{uid}))$ returns a non-empty set of positions $P(\text{uid})$ such that: $P_{\text{apply}}(\text{uid}, \sigma) \subseteq P(\text{uid}) \subseteq (P_{\text{apply}}(\text{uid}, \sigma) \cup P_{\text{move}}(\text{uid}, \sigma))$. Note that if $P_{\text{apply}}(\text{uid}, \sigma)$ is empty, and the value is not \perp , then the position has at least one element from $P_{\text{move}}(\sigma)$.

Under these semantics, the position set of an element contains all positions set by APPLY operations, as such operations update the value and therefore it is important they not be overwritten. It may or may not contain positions set by MOVE operations that occur concurrently with other updates, as these do not impact the underlying data.

5 ORARRAY ALGORITHM

In Section 5.1 we present a new dot store, and in Section 5.2 we present our algorithm implementing the δ -based ORArray.

5.1 The COMPDOTFUN

To construct the ORArray we first propose a new dot store. We combine the DOTMAP and DOTFUN to get a dot store which maps dots to dot stores. We assume the dot stores in the range of the DOTFUN are all join-semilattices. We will show that the resulting mapping is also a join-semilattice. Thus, this construction yields a composable DOTFUN . We denote this dot store as COMPDOTFUN . The dot store is presented in Equation 6. Note that the dots method returns the dots in the domain, as well as a union on recursive calls of dots on all dot stores in the range.

$$\begin{aligned} \text{COMPDOTFUN}(V : \text{DotStore}) : \mathcal{I} \times \mathbb{N} &\mapsto V \\ \text{dots}(m) = \text{dom } m \cup \bigcup_{v \in \text{ran } m} &\text{dots}(v) \end{aligned} \quad (6)$$

We define its corresponding join-semilattice in Equation 7. The join operation keeps keys that have not been deleted (as in the DOTFUN), or the values themselves haven't been deleted (as in the

DOTMAP).

$$\begin{aligned} (m, c) \sqcup (m', c') &= \\ &(\{d \mapsto v(d) \mid d \in \text{dom } m \cap \text{dom } m' \wedge v(d) \neq \perp\} \\ &\cup \{(d, v) \in m \mid d \notin c'\} \cup \{(d, v) \in m' \mid d \notin c\}, c \cup c') \end{aligned} \quad (7)$$

where $v(d) = \text{fst}((m(d), c) \sqcup (m'(d), c'))$

In Appendix A, we show and prove the claim that given three elements in the $\text{Causal}(\text{COMPDOTFUN})$, the resulting object is depends on set union and the join of the inner dot store of the three elements. As set union is commutative, associative, and idempotent, and as we assume that the inner dot store is a join semilattice we arrive at the following corollary:

COROLLARY 5.1. *The $\text{Causal}(\text{COMPDOTFUN})$ is a join-semilattice.*

5.2 The ORArray

We now propose an algorithm implementing the δ -based ORArray. The ORArray is instantiated with some set of stable position identifiers \mathcal{P} . Algorithm 3 presents our proposal. We store the position in the following dot store:

$$\text{Pos} \triangleq \text{COMPDOTFUN}(\text{DOTFUN}(\mathcal{P}))$$

The state and deltas are generated by the algorithm are from the following join-semilattice:

$$\text{CAUSAL}(\text{DOTMAP}(\mathcal{I} \times \mathbb{N}, (\text{DotStore}, \text{Pos}))).$$

In other words, the ORArray is stored as a mapping from dots to pairs of value-positions, where the value is any CRDT, and the set of positions is stored using a DOTFUN nested in a COMPDOTFUN . This position construction is what allows for concurrent updates. This construction stores the set of positions as a forest, where each tree is of depth 3. For example, an element whose set of positions is $\{0.3, 54\}$ may have the position stored in the following forest:

$$\{d_1 \mapsto \{d_2 \mapsto 0.3\}, d_3 \mapsto \{d_4 \mapsto 54\}\},$$

where d_1, d_2, d_3 , and d_4 are dots. The APPLY function adds a new tree and removes all the old ones by removing their roots. In the example above, the roots d_1 and d_3 are removed and a new tree is inserted. As the root is stored in a COMPDOTFUN , once a root is removed it never reappears, contrary to keys in a DOTMAP which may remain undeleted if there is a concurrent update. The MOVE function removes the child of all roots and adds a single child to every root with the new position. Consider the example above where the position is set to 42. The dots d_2 and d_4 , and the forest is altered to

$$\{d_1 \mapsto \{d_5 \mapsto 42\}, d_3 \mapsto \{d_5 \mapsto 42\}\}.$$

The DELETE function removes all roots, thus removing the position. A MOVE may be overwritten by a concurrent APPLY or DELETE , as the APPLY and DELETE remove the roots. But an APPLY always adds a new root, thus an element cannot have a non-empty value and be without any possible position.

As an APPLY inserts a new tree into the forest, concurrent APPLY operations create multiple trees. For example, if one apply sets the position to $\{d_1 \mapsto \{d_2 \mapsto 0.3\}\}$ and another one sets the position to $\{d_3 \mapsto \{d_4 \mapsto 54\}\}$, after they are merged the resulting forest is the one appearing above.

Algorithm 3 ORArray algorithm for node i .

```
1: Global state  $(m, c)$  ▷ CAUSAL(DOTMAP( $\mathcal{I} \times$   
   $\mathbb{N}$ ,  $(\text{DotStore}, \text{COMPDOTFUN}(\text{DOTFUN}(\mathcal{P}))))$ )  
2: procedure APPLY $_i(\text{uid}, o_i^\delta, p)$   
3:    $d \leftarrow \text{next}_i(c)$   
4:    $(v, c') \leftarrow o_i^\delta(\text{fst}(m[\text{uid}]), c \cup \{d\})$   
5:    $\text{roots} \leftarrow \{d' \mid d' \in \text{dom } \text{scnd}(m[d])\}$  ▷ Get all roots  
6:    $m' \leftarrow \{\text{uid} \mapsto (v, \{d \mapsto \{d \mapsto p\})\})\}$   
7:    $c' \leftarrow c' \cup \{d\} \cup \text{roots}$   
8:   return  $(m', c')$   
9: procedure MOVE $_i(\text{uid}, p)$   
10:   $d \leftarrow \text{next}_i(c)$   
11:   $\text{roots} \leftarrow \{d' \mid d' \in \text{dom } \text{scnd}(m[d])\}$  ▷ Get all roots  
12:   $\text{children} \leftarrow \bigcup_{\text{child} \in \text{dom}} (\text{ran } \text{scnd}(m[\text{uid}])) \text{ dots}(\text{child})$   
13:   $\text{position} \leftarrow \{r \mapsto \{d \mapsto p\} \mid r \in \text{roots}\}$   
14:   $m' \leftarrow \{\text{uid} \mapsto (\{\}, \text{position})\}$   
15:   $c' \leftarrow \{d\} \cup \text{children}$   
16:  return  $(m', c')$   
17: procedure DELETE $_i(\text{uid})$   
18:   $c' \leftarrow \text{dots}(m[\text{uid}])$   
19:  return  $(\{\}, c')$ 
```

We now prove that Algorithm 3 implements an ORArray. We begin with the values. Note that the implementations of the APPLY $_i$ and the DELETE $_i$ operations in Algorithm 3 are, disregarding Line 3 and Line 5, identical to those in Algorithm 2. The difference in the APPLY $_i$ operation is due to updating the position identifier. Furthermore, the semantics as defined in Section 4 are identical. Therefore, we begin with the following observation:

OBSERVATION 1. *The value of an element in an object implemented by Algorithm 3 adheres to OR semantics.*

We next show that the set of positions of an element P satisfies the semantics described in Section 4.2. We show each containment separately.

LEMMA 5.2. *Consider an execution σ of Algorithm 3. Let P be the set of possible positions of some element at the end of σ . Then $P_{\text{apply}}(\text{uid}, \sigma) \subseteq P$.*

PROOF. If $p \in P_{\text{apply}}(\text{uid}, \sigma)$, then there exists some APPLY $(\text{uid}, o_i^\delta, p)$ in $\text{maximal}(\sigma, \text{uid})$. When APPLY $(\text{uid}, o_i^\delta, p)$ executes, it execute Line 3, generating a new dot d . It then adds $\{d \mapsto \{d \mapsto p\}\}$ as the position in the delta it generates X^δ . As the operation is a maximal operation, no other operation has seen dot d . In particular, the local causal context does not contain d . Therefore, the position survives the join, and $\{p \mapsto \{p \mapsto d\}\} \in \text{scnd}(\text{GET}(\text{uid}))$ for a GET that was to execute at this point. Therefore, $p \in P$.

This is true for all $p \in P_{\text{apply}}(\text{uid}, \sigma)$, and therefore:

$$P_{\text{apply}}(\text{uid}, \sigma) \subseteq P. \quad \square$$

We now show the right hand side.

LEMMA 5.3. *Consider an execution σ of Algorithm 3. Let P be the set of possible positions of some element at the end of σ . Then $P \subseteq P_{\text{apply}}(\text{uid}, \sigma) \cup P_{\text{move}}(\text{uid}, \sigma)$.*

PROOF. We prove the lemma by showing that if $p \notin P_{\text{apply}}(\text{uid}, \sigma) \cup P_{\text{move}}(\text{uid}, \sigma)$, then $p \notin P$ for any p . Consider some $p \notin P_{\text{apply}}(\text{uid}, \sigma) \cup P_{\text{move}}(\text{uid}, \sigma)$. Trivially, if there exists no operation in $\text{Ops}(\sigma)$ that sets the position to p , then $p \notin P$.

Consider the case that there exists either a APPLY $(\text{uid}, o_i^\delta, p) \in \text{Ops}(\sigma)$ or MOVE $(\text{uid}, o_i^\delta, p) \in \text{Ops}(\sigma)$. Denote the operation as o_p . As $p \notin P_{\text{apply}}(\text{uid}, \sigma) \cup P_{\text{move}}(\text{uid}, \sigma)$, then $o_p \notin \text{maximal}(\sigma, \text{uid})$. Therefore, there exists some other operation $o \neq \text{GET} \in \text{Ops}(\sigma)$ such that $o_p <_\sigma o$.

When o_p returns a delta X^δ , then the position in X^δ is of the form $\{d_1 \mapsto \{d^1 \mapsto p\}, d_2 \mapsto \{d^2 \mapsto p\}, \dots, d_k \mapsto \{d^k \mapsto p\}\}$ for some integer k . We now show that no matter what type of operation o is, p is removed after o . If o is an APPLY, then on Line 5 dots $\{d_1, d_2, \dots, d_k\}$ are added to the causal context and not added to the map. Denote the delta returned by o as $X^{\delta'}$. Then when X^δ is joined with $X^{\delta'}$, these dots do not survive that join, as they are keys in a COMPDOTFUN. If o is a DELETE, then dots $\{d_1, d_2, \dots, d_k\}$ are added to the causal context and the result is identical. If o is a MOVE, then on Line 12 dots $\{d^1, d^2, \dots, d^k\}$ are added to the delta's causal context. As this join is in a DOTFUN, the dots don't survive the join.

In either one of the cases, the position does not survive the execution of o , and therefore $p \notin P$. \square

LEMMA 5.4. *Consider an execution σ of Algorithm 3. Let P be the set of possible positions of some element at the end of σ , and denote the value of the element v . If $v \neq \perp$ then $P \neq \emptyset$*

PROOF. Consider some execution σ of Algorithm 3, and consider some element uid , such that its corresponding value v is not \perp . If $P_{\text{apply}}(\text{uid}, \sigma) \neq \emptyset$ then the lemma is proven, as $P_{\text{apply}}(\text{uid}, \sigma) \subseteq P$. Assume, therefore, that $P_{\text{apply}}(\text{uid}, \sigma) = \emptyset$.

As $v \neq \perp$, there exists some APPLY operation that does not precede a DELETE in σ . Consider the maximal such APPLY operation, i.e., an APPLY that does not precede a DELETE or another APPLY operation in σ . Let d be the new dot generated on Line 3 during the execution of the operation. As $P_{\text{apply}}(\text{uid}, \sigma) = \emptyset$, then there exists a MOVE operation such that the APPLY precedes the MOVE in σ . Consider the maximal such MOVE in σ . During the execution of the operation, the MOVE adds $\{d \mapsto \{d' \mapsto p\}\}$ to the position map. As d wasn't yet deleted (as it is only removed due a subsequent APPLY or DELETE), and d' is a new dot generated by the MOVE, the local state satisfies: $\{d \mapsto \{d' \mapsto p\}\} \in \text{scnd}(m[\text{uid}])$. Therefore $p \in P$, and $P \neq \emptyset$. \square

6 DSON

We now describe the construction of DSON. A JSON document [12] is defined recursively as follows:

```
REGISTER = String|Number|null  
MAP = {String : VALUE, ..., String : VALUE}  
ARRAY = [VALUE, VALUE, ..., VALUE]  
VALUE = MAP|ARRAY|REGISTER
```

A JSON document is a VALUE itself. The identifier of a VALUE in a JSON document is the path followed from the root to get to the

VALUE, starting from doc . For example, given the JSON document:

$$\{k_1 : v_1, k_2 : [v_2, v_3], k_3 : \{k_4 : v_4\}\},$$

the identifier of v_1 is $doc.k_1$, the identifier of v_3 is $doc.k_2[1]$, and the identifier of the value v_4 is $doc.k_3.k_4$.

The JSON document is one of three objects: REGISTER, MAP, and ARRAY. We define each object's API methods:

- REGISTER:** (1) WRITE(x) – sets register value to x ;
(2) READ() – returns the register value.
- MAP:** (1) APPLY(k, o_i^δ) – given a key k and function o_i^δ , applies o_i^δ to key k ;
(2) GET(k) – returns the VALUE of key k ;
(3) REMOVE(k) – removes key k from the map.
- ARRAY:** (1) APPLY(uid, o_i^δ, p) – given a unique identifier uid , function o_i^δ and position p , applies o_i^δ to element uid and sets its position to p ;
(2) MOVE(uid, o_i^δ, p) – given a unique identifier uid and position p , set the position of element uid to p ;
(3) GET(uid) – returns the VALUE of element uid ;
(4) REMOVE(uid) – removes the element uid .

An *Observed-Remove* JSON (OR-JSON) CRDT is one in which the register, map, and array all adhere to OR semantics, as defined in Section 4. We achieve this by implementing the register using the MVReg, the map using the ORMap, and the array using the ORArray. The replicas communicate with each other via a causal anti-entropy algorithm, and each replica is sequential.

Our solution provides the following guarantees:

- *Read-your-writes* [43] – if a read method is executed following a write method in the same replica, then the read sees the write;
- *Causal Consistency* [30] – if a read method sees a write w_1 , and w_2 causally precedes w_1 , then the read sees w_2 ;
- *Strong Eventual Consistency* – all nodes eventually converge on the same state, such that identical reads on different nodes return identical values.

Read-your-writes is achieved because every replica is sequential, and it merges the new delta with its local state one an operation returns, before any subsequent reads. Eventual consistency is achieved by (1) using some anti-entropy algorithm that ensures that the latest states or updates eventually reach all replicas, e.g., as defined in [18]; and (2) having the states and the deltas be elements of a join-semilattice. Causal consistency is achieved by using a causal anti entropy algorithm, e.g., one that only delivers a delta once all deltas causally preceding it have been delivered.

As previously noted, an important property for our JSON CRDT is bounding the stored metadata. In our case, the state is the pair (m, c) , where m stores the state of the underlying document and c is the causal context. We now analyze the metadata overhead.

We first consider the steady state, i.e., after all conflicts have been resolved and the maximal updates are not concurrent to any others. In this case, as explained in Section 3.3, the causal context can be efficiently stored in a map from node identifiers to integers, therefore its size is $O(n \log n)$, where n is the number of replicas. The ORMap does not store any dots, therefore comes with no metadata overhead. The MVReg stores a dot per value, therefore has a metadata size of $O(1)$ per value. An element in the ORArray has

1 dot as the uid , and 2 dots for the position, therefore has a size of $O(1)$ per array element. Therefore, the metadata stored for a document with D registers and array elements is $O(D + n \log n)$.

We now analyze the metadata in the presence of k concurrent operations. The worst case is for $k/2$ operations to be APPLY methods to the same array element (creating $k/2$ new roots), followed by all replicas receiving all deltas, and then $k/2$ operations to be MOVE methods to that same element (creating $k/2$ children for each root). In this case, the metadata for the position has $O(k^2)$ dots. As we use a causal anti-entropy algorithm, causal context remains $O(n \log n)$. Therefore, in this case the metadata stored for a document of size D is bounded by $O(k^2 D + n \log n)$.

7 IMPLEMENTATION

We implemented DSON in around 2500 lines of JavaScript code[8]. We expose 5 methods from the external API of Automerge [2]:

- init()** initializes a new empty document.
from(obj) creates a new document and populates it with the contents of a given JSON object obj .
change(doc, f): applies function f to document doc locally as well as locally storing the delta to be broadcast.
applyChanges(doc, delta): merges delta $delta$ with the local document doc .
getChanges(doc): coalesces all deltas generated for a document since the previous getChanges operation for sending to another replica.

In a nutshell, $change(doc, f)$ returns the delta X^δ that represents the alteration of document doc by function f , and $applyChanges(doc, delta)$ computes $X \sqcup X^\delta$ where X is the state of document doc , as defined in previous sections. Separate libraries implement disk persistence and network communication, e.g., encryption, authentication, access control, and implementation of an anti-entropy algorithm.

To represent a JSON object, the top level of DSON is a map. Recall that ORMaps and ORArrays recursively apply changes via a delta mutator. These recursive functions can be complex to write and reason about. To make the interaction with DSON more accessible, we follow a path similar to Automerge [2], utilizing JavaScript proxies [11] to allow a user to interact with DSON using the standard JavaScript programming model. To that end, we implemented a mechanism that transforms these operations into delta mutators and applies them. For example, figure 2 depicts the update $doc.k_1[0] = \{k_2 : 1\}$, we first create a delta mutator representing the left hand side to identify the path in the document to which this mutator is to be applied. Then the value is created by creating an MVReg for the literal and a new ORMap is created into which the MVReg is inserted under the key k_2 .

We also implemented the wrapper described in Section 5.2, presented in Algorithm 4. This wrapper enables programmatic access to arrays using JavaScript indices (integers), while internally array positions are encoded using stable identifiers as defined in Logoot [44] and LSEQ [33]. Note that the code altering the state is decoupled from the state itself, so that any choice of stable identifiers can be plugged in. The function $GETUIDFROMIDX(idx, (m, c))$ ($GETPOSFROMIDX(idx, (m, c))$) is implemented by sorting the array and returning the unique identifier (position) of the element at the index idx .

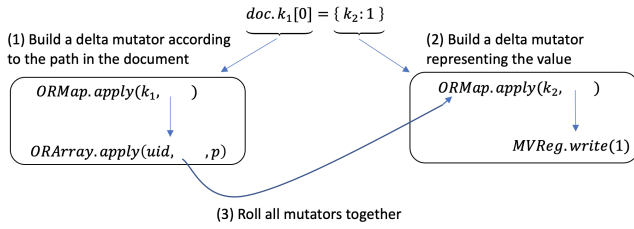


Figure 2: Transform a JavaScript operation $doc.k_1[0] = \{k_2: 1\}$ to a delta mutation in the JSON CRDT, where $p = \text{GETPosFROMIDX}(0)$ and $uid = \text{GETUIDFROMIDX}(0)$.

Concurrent writes support both *value* and *type* conflicts. A value conflict is where the same variable is concurrently set to different values, e.g., replica 1 sets $user.last_access = 14:49$ and replica 2 sets $user.last_access = 14:52$. This is handled by the MVReg.

Type conflicts arise because JSON objects are weakly typed. The type of a variable can change over time, and variables can be set concurrently to different types, for example, replica 1 may set $a = 1$, while replica 2 may concurrently set $a = \{\}$. The resulting variable is not defined by a single type. The challenge is that our δ -based CRDTs do not support joins across dot stores, e.g., a join between a *DotFun* and a *DotMap* is undefined. To circumvent this, we nest each variable in an ORMap containing three keys, one corresponding to a map, one to an array and one to a value, and arbitrarily choose the following order $Map > Array > MVReg$. A type change is implemented by adding an element to this triple, and removes the other keys. The CRDT eventually converges to a single type, although values that were written with previous types may remain until they are overwritten, which can be arbitrarily long.

8 EVALUATION

We compare DSON to two open-source libraries implementing JSON CRDTs: Automerge [2] and Yjs [15]. Both of these libraries optimize for collaborative text editing, but nevertheless have full support for JSON CRDTs. We note that while these libraries optimize for different behavior, they are the only CRDT libraries known to us providing full JSON datatype. Riak [14] does not provide an array data type, therefore we do not compare to it. The goal of the evaluation section is to see how the metadata size varies with the number of operations. Unlike Automerge and Yjs, our code is a research prototype and has not been optimized for performance, so we emphasize overall trends. We ran all tests on a 4-core Intel i7-8565U laptop. All documents are initialized empty, with no warmup done. Measuring document size was done using the byte size of the latest encoding provided by each library. All benchmarks are available together with the DSON source code [8].

We evaluate the map and array separately, and first evaluate the map. In our first benchmark, we update the same key in a map consecutively with different values. The document size remains constant (it remains a single key value pair) because we do not add keys. The results are presented in Figure 3a. DSON and Yjs both maintain constant metadata size, whereas the metadata in Automerge grows linearly with the number of operations.

Algorithm 4 Wrapper for ORArray.

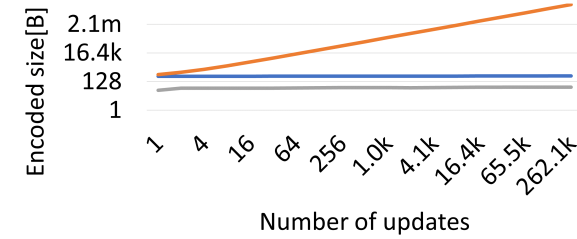
IDX represents a JavaScript index (integer), UID represents a dot, and POS represents a stable identifier.

- 1: Shared state (m, c) with Algorithm 3
 - 2: Function $\text{GETPosFROMIDX}(idx, (m, c))$ \triangleright Converts index idx to stable position identifier $p \in P$.
 - 3: Function $\text{GETUIDFROMIDX}(idx, (m, c))$ \triangleright Returns the uid of the element at index idx .
 - 4: **procedure** $\text{INSERT}(idx, o_i^\delta)$
 - 5: $p \leftarrow \text{GETPosFROMIDX}(idx, (m, c))$
 - 6: $uid \leftarrow \text{next}_i(c)$
 - 7: **return** $\text{APPLY}_i(uid, p, o_i^\delta)$
 - 8: **procedure** $\text{UPDATE}(idx, o_i^\delta)$
 - 9: $d \leftarrow \text{next}_i(c)$
 - 10: $p \leftarrow \text{GETPosFROMIDX}(idx, (m, c))$
 - 11: $uid \leftarrow \text{GETUIDFROMIDX}(idx)$
 - 12: **return** $\text{APPLY}_i(uid, p, o_i^\delta)$
 - 13: **procedure** $\text{MOVE}(old_idx, new_idx)$
 - 14: $uid \leftarrow \text{GETUIDFROMIDX}(old_idx, (m, c))$
 - 15: $p \leftarrow \text{GETPosFROMIDX}(new_idx, (m, c))$
 - 16: **return** $\text{MOVE}_i(uid, p)$
 - 17: **procedure** $\text{DELETE}(idx)$
 - 18: $uid \leftarrow \text{GETUIDFROMIDX}(idx, state)$
 - 19: **return** $\text{DELETE}_i(uid)$
-

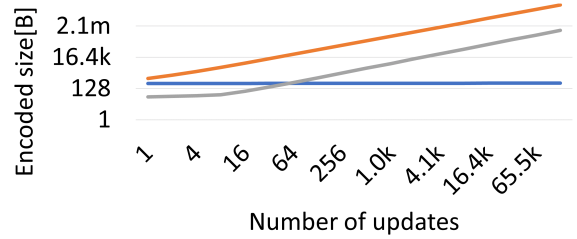
Next, we continuously insert and subsequently delete a key in a map. In this case the document is empty after every insert-delete pair. Figure 3b presents the results. In this workload, DSON maintains constant metadata size, whereas the metadata in both Automerge and Yjs is unbounded. The Yjs behavior is a result of its garbage collector - it creates a local object for every key which it can only clean up after its parent object has been deleted, to avoid sync conflicts [16]. Yjs attempts to garbage collect more efficiently for certain patterns but a complete and efficient solution is inherently difficult. Our delta based approach avoids the need for garbage collection altogether.

We now evaluate the array. We first measure an update only workload, i.e., the same index is consecutively updated. The results are shown in Figure 3c. As expected, the DSON's metadata size remains constant. This is also the case when continuously inserting and removing array elements, as shown in Figure 3d. In the latter case Yjs also achieves constant metadata, because YATA [34] handles arrays efficiently. Note, however, that Yjs is optimized for character insertions and deletions. When inserting and removing a complex data type, the optimizations break down and the metadata grows. Figure 3e presents the metadata growth as a function of consecutively inserting and removing a map in an array, and Figure 3f presents the metadata growth as a function of inserting and removing an array in an array. While DSON keeps constant metadata, that of Yjs grows.

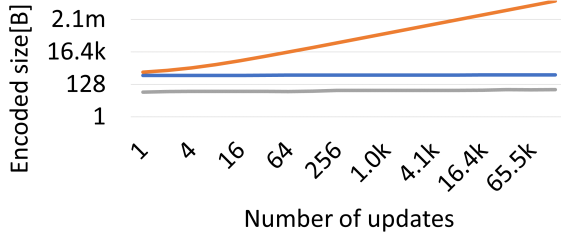
In Section 6 we outlined our worst case metadata growth as $O(k^2)$ for k concurrent operations. This stems from how we store the position of an array element. To achieve this worst case growth, the k operations need to concurrently execute an APPLY to the same element, then deliver each other's messages, then execute MOVE



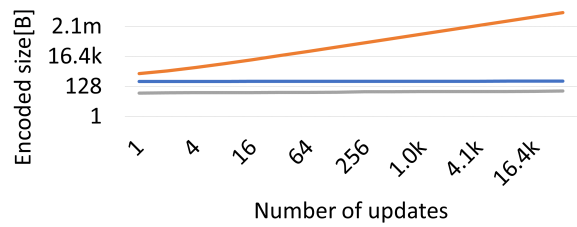
(a) Map size when updating the same key in a map.



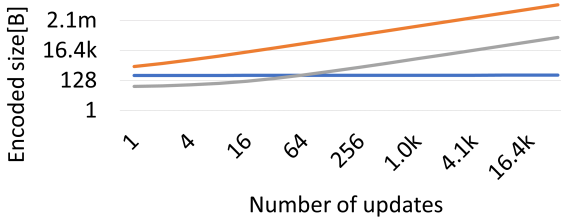
(b) Map size when inserting and removing a map key.



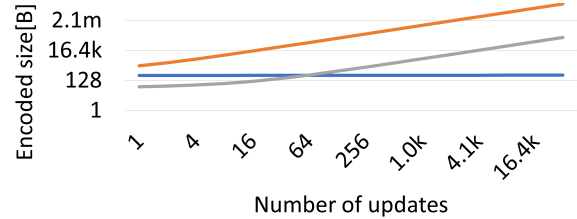
(c) Array size when consecutively updating the same index in an array.



(d) Array size when consecutively inserting and removing a character in an array.



(e) Array size when consecutively inserting and removing a map in an array.



(f) Array size when consecutively inserting and removing an array in an array.

—DSON —Automerge —Yjs

Figure 3: Micro benchmarks measuring map and array sizes.

operations to that same element concurrently, and finally deliver each other’s messages. In this execution, the position of the element has k roots, each with k children, resulting in k^2 dots, giving $O(k^2)$ metadata growth.

To show the worst-case behavior, all replicas update every array element, then all replicas deliver each other’s messages, and then all replicas sort the array. The results are presented in Figure 4. In this case, the empirical analysis matches the theoretical expectation. We argue that this scenario is rare, and under better conditions the metadata grows significantly slower than in the worst case

To show this, we implement the SORT operation. Sorting an array is a common and useful operation, and is efficiently done with our MOVE operation. As explained in Section 5, the array values are not copied at all, rather only their position changes. Furthermore, sorting can be concurrently done by different replicas, which sort the array into similar states.

We use a randomized micro-benchmark with 5 replicas to capture as many scenarios as possible. The benchmark is defined by

two parameters, an update probability p_u and a sync probability p_s . Each run executes 200 steps, where in every step a randomly chosen replica performs an UPDATE with probability p_u , and a SORT with probability $1 - p_u$. After every step, the replicas sync with probability p_s . Each run is repeated 10 times, where every run measures the maximum document size within the run. The result of the micro-benchmark is the average of the maximums, with error bars depicting the standard deviation. The results are presented in Figure 5. The worst case memory size is taken from Figure 4, and the best case is measured when the array is created, before any concurrent operations occur, for the appropriate number of replicas. The analysis shows that the document size in these tests remains far from the worst-case scenario.

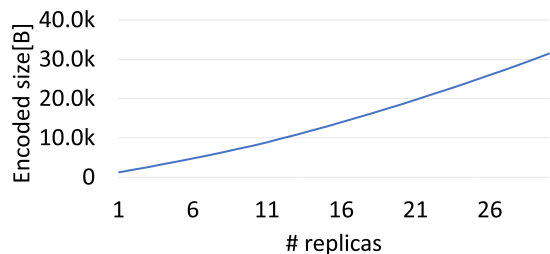


Figure 4: Document size where all replicas update all elements of an array, then sync, and then sort the array, as a function of the number of replicas.

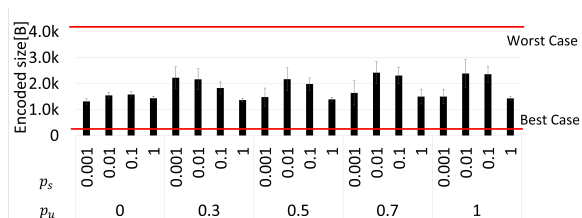


Figure 5: Document size for 5 replicas, where each replica either performs an update with probability p_u or sorts with probability $1 - p_u$. All replicas synchronize with probability p_s .

9 FUTURE WORK

9.1 Performance Optimizations

Enes et al. have shown that practical anti-entropy algorithms are wasteful, and in the end δ -based CRDTs perform no better than the state based solution [23]. They identify the concept of *join decomposition* which can be used to obtain optimal deltas, and show how this concept is applied to existing δ -based CRDTs. Future work is needed to construct the join decomposition for the δ -based CRDTs mentioned in this paper.

9.2 Ad-hoc Clients

We have shown that the size of the causal context is $O(n \log n)$ where n is the number of replicas. For a system with a bounded number of long-lived replicas this overhead remains reasonably small. For example, when deploying distributed document stores as a global cloud service, one could provide strong consistency within a region, e.g., as is done for IBM Cloudant on Transaction Engine [10], and use DSON to coordinate updates globally. In this case, each region would be considered a single replica, and the number of replicas is bounded and fairly static. Note that this requires clients to be connected to some replica. For systems with an unbounded number of short-lived replicas this overhead may become significant, since any replica that connects and executes operations leaves a permanent footprint.

We now differentiate between long-lived replicas, which typically represent servers connecting the network, and short-lived replicas which represent ad-hoc, possibly disconnected, clients. For simplicity, we refer to the latter as *clients*. We propose a solution

which supports an unbounded number of clients while maintaining low metadata overheads.

To this end, we note that the purpose of the dot is to uniquely identify an operation. Moreover, dots needn't dictate a causal order – up until now given two dots (c, i) and $(c, i + k)$ for some $k \geq 0$ we could say that the operation that generated (c, i) causally precedes the operation that generated $(c, i + k)$. However, this is not an inherent requirement, as the causal context is used to maintain the causal order.

We propose using a mixture of local and global dots. A client c is connected to a single replica r , and runs the same algorithm as for mutating the JSON CRDT. When the client generates a dot, it considers itself a replica and does so normally, generating a dot (c, i) for some integer i . When the client propagates its delta, it passes along a list of its locally generated dots contained in the delta $(c, i), (c, i + 1), \dots, (c, i + k)$. When the server receives the delta and list, for every client dot in the list it generates its dot, mapping $(c, i + j)$ to $(r, l + j)$ for all $1 \leq j \leq k$, where (r, l) would be the next dot generated right before the delta is received.

After generating the mapping, the replica scans the delta and alters all client dots to the corresponding replica dots. It then sends the mapping back to client c , and propagates the transformed delta to all other clients and replicas. This method allows the clients to operate as if they were replicas, while maintaining a small causal context. One drawback is if the mapping sent back from replica r to client c is lost or delayed for long enough, the client may assume the replica crashed. In such a case it may send its delta to a different replica r' , which will reapply the changes. In such a case we may have a double commit of an update done by the client.

9.3 Efficient Replica Removal

The CRDT approach allows replicas to be added to the network dynamically at no additional cost, as this case is equivalent to the replica joining from the start without sending messages. Efficiently removing replicas from the network is more challenging. The changes made by a replica are retained in the causal context, which leads to an increase in metadata over time. Simply removing the relevant dots from the causal context may harm causality, since it may lead to unwanted re-application of updates. A safe method of cleaning up metadata is left for further work.

10 CONCLUSIONS

We presented DSON, a space efficient δ -based CRDT approach for distributed JSON document stores, enabling high availability at a global scale. We formally defined DSON's semantics, and proved its correctness and convergence, providing consistency guarantees which relieve application developers of the burden of reasoning about complex distributed application logic. The amount of metadata DSON stores is not dependent on the number of document updates, unlike previous approaches. We evaluated DSON empirically in this respect and showed that the metadata size stays constant as expected while varying the number of document updates, unlike the case for Automerge and Yjs which are optimized for collaborative document editing. DSON is a sound theoretical basis for efficient and robust highly available distributed document stores, which we believe can lead to efficient systems in practice.

REFERENCES

- [1] 2021. Akka. <https://akka.io/>.
- [2] 2021. Automerge. <https://github.com/automerge/automerge>.
- [3] 2021. Couchbase. <https://www.couchbase.com/>.
- [4] 2021. Couchbase Availability. <https://docs.couchbase.com/server/5.5/understanding-couchbase/clusters-and-availability/replication-architecture.html>.
- [5] 2021. CouchDB. <https://couchdb.apache.org/>.
- [6] 2021. DB-Engines Ranking of Document Stores. <https://db-engines.com/en/ranking/document+store>.
- [7] 2021. DBMS popularity broken down by database model. https://db-engines.com/en/ranking_categories.
- [8] 2021. DSON - JSON CRDT Using Delta-Mutations. <https://github.com/crdt-ibm-research/json-delta-crdt>.
- [9] 2021. DynamoDB. <https://aws.amazon.com/dynamodb/>.
- [10] 2021. IBM Cloudant on Transaction Engine Documentation. <https://cloud.ibm.com/docs/Cloudant?topic=Cloudant-overview-te>.
- [11] 2021. JavaScript Proxy. https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Global_Objects/Proxy.
- [12] 2021. JSON. <https://www.json.org/json-en.html>.
- [13] 2021. Peer-base/JS-delta-crdts: Delta State-based CRDTs in Javascript. <https://github.com/peer-base/js-delta-crdts>
- [14] 2021. Riak. <https://riak.com/index.html>.
- [15] 2021. Yjs. <https://github.com/yjs/yjs>.
- [16] 2021. Yjs Forum: Map Metadata Overhead. <https://discuss.yjs.dev/t/map-metadata-overhead/492>.
- [17] Paulo Sérgio Almeida, Ali Shoker, and Carlos Baquero. 2015. Efficient state-based crdts by delta-mutation. In *International Conference on Networked Systems*. Springer, 62–76.
- [18] Paulo Sérgio Almeida, Ali Shoker, and Carlos Baquero. 2018. Delta state replicated data types. *J. Parallel and Distrib. Comput.* 111 (2018), 162–173.
- [19] J Chris Anderson, Jan Lehnardt, and Noah Slater. 2010. *CouchDB: The Definitive Guide, chapter 17: Conflict Management*. " O'Reilly Media, Inc."
- [20] Carlos Baquero, Paulo Sérgio Almeida, and Ali Shoker. 2014. Making operation-based CRDTs operation-based. In *IFIP International Conference on Distributed Applications and Interoperable Systems*. Springer, 126–140.
- [21] Juan Benet. 2014. IpfS-content addressed, versioned, p2p file system. *arXiv preprint arXiv:1407.3561* (2014).
- [22] Ivan Chajda, Radomír Halaš, and Jan Kühr. 2007. *Semilattice structures*. Vol. 30. Heldermann Lemgo.
- [23] Vitor Enes, Paulo Sérgio Almeida, Carlos Baquero, and João Leitão. 2019. Efficient synchronization of state-based CRDTs. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. IEEE, 148–159.
- [24] Pascal Grosch, Roman Krafft, Marcel Wölki, and Annette Bieniusa. 2020. AutoCouch: A JSON CRDT Framework. In *7th Workshop on Principles and Practice of Consistency for Distributed Data (PaPoC 2020)*. ACM, Article 6. <https://doi.org/10.1145/3380787.3393679>
- [25] Pat Helland. 2012. Idempotence Is Not a Medical Condition: An essential property for reliable systems. *Queue* 10, 4 (2012), 30–46.
- [26] Andrew Jeffery, Heidi Howard, and Richard Mortier. 2021. Rearchitecting Kubernetes for the Edge. In *Proceedings of the 4th International Workshop on Edge Systems, Analytics and Networking*. 7–12.
- [27] Martin Kleppmann. 2020. Moving elements in list CRDTs. In *Proceedings of the 7th Workshop on Principles and Practice of Consistency for Distributed Data*. 1–6.
- [28] Martin Kleppmann and Alastair R Beresford. 2017. A Conflict-Free Replicated JSON Datatype. *IEEE Transactions on Parallel and Distributed Systems* 28, 10 (April 2017), 2733–2746. <https://doi.org/10.1109/TPDS.2017.2697382> arXiv:1608.03960
- [29] Rusty Klophaus. 2010. Riak core: Building distributed applications without shared state. In *ACM SIGPLAN Commercial Users of Functional Programming*. 1–1.
- [30] Leslie Lamport. 1978. Time, clocks, and the ordering of events in a distributed system. In *Concurrency: the Works of Leslie Lamport*. 179–196.
- [31] Lars Larsson, Harald Gustafsson, Cristian Klein, and Erik Elmroth. 2020. Decentralized Kubernetes Federation Control Plane. In *2020 IEEE/ACM 13th International Conference on Utility and Cloud Computing (UCC)*. IEEE, 354–359.
- [32] M Mihai Letia, N Pregaica, and M Shapiro. 2009. CRDTs: Consistency without concurrency control. *RR-6956, INRIA* (2009).
- [33] Brice Nédelec, Pascal Molli, Achour Mostefaoui, and Emmanuel Desmontils. 2013. LSEQ: an adaptive structure for sequences in distributed collaborative editing. In *Proceedings of the 2013 ACM symposium on Document engineering*. 37–46.
- [34] Petru Nicolaescu, Kevin Jahns, Michael Derntl, and Ralf Klamma. 2016. Near Real-Time Peer-to-Peer Shared Editing on Extensible Data Types. In *19th International Conference on Supporting Group Work (GROUP 2016)*. ACM, 39–49. <https://doi.org/10.1145/2957276.2957310>
- [35] Gérard Oster, Pascal Urso, Pascal Molli, and Abdessamad Imine. 2005. *Real time group editors without operational transformation*. Research Report RR-5580. INRIA. <https://hal.inria.fr/inria-00071240/document>
- [36] Gérard Oster, Pascal Urso, Pascal Molli, and Abdessamad Imine. 2006. Data consistency for P2P collaborative editing. In *Proceedings of the 2006 20th anniversary conference on Computer supported cooperative work*. 259–268.
- [37] Nuno Pregaica, Joan Manuel Marques, Marc Shapiro, and Mihai Letia. 2009. A commutative replicated data type for cooperative editing. In *2009 29th IEEE International Conference on Distributed Computing Systems*. IEEE, 395–403.
- [38] Nuno Pregaica, Joan Manuel Marques, Marc Shapiro, and Mihai Letia. 2009. A Commutative Replicated Data Type for Cooperative Editing. In *29th IEEE International Conference on Distributed Computing Systems (ICDCS 2009)*. IEEE, 395–403. <https://doi.org/10.1109/ICDCS.2009.20>
- [39] Arik Rinberg, Tomer Solomon, Guy Khazma, Gal Lushi, Roei Shlomo, and Paula Ta-Shma. 2021. Array CRDTs Using Delta-Mutations. In *Proceedings of the 8th Workshop on Principles and Practice of Consistency for Distributed Data*. 1–3.
- [40] Hyun-Gul Roh, Myeongjae Jeon, Jin-Soo Kim, and Joonwon Lee. 2011. Replicated abstract data types: Building blocks for collaborative applications. *J. Parallel and Distrib. Comput.* 71, 3 (2011), 354–368.
- [41] Hyun-Gul Roh, Myeongjae Jeon, Jin-Soo Kim, and Joonwon Lee. 2011. Replicated Abstract Data Types: Building Blocks for Collaborative Applications. *J. Parallel and Distrib. Comput.* 71, 3 (March 2011), 354–368. <https://doi.org/10.1016/j.jpdc.2010.12.006>
- [42] Marc Shapiro, Nuno Pregaica, Carlos Baquero, and Marek Zawirski. 2011. Conflict-free replicated data types. In *Symposium on Self-Stabilizing Systems*. Springer, 386–400.
- [43] Douglas B Terry, Alan J Demers, Karin Petersen, Mike J Spreitzer, Marvin M Theimer, and Brent B Welch. 1994. Session guarantees for weakly consistent replicated data. In *Proceedings of 3rd International Conference on Parallel and Distributed Information Systems*. IEEE, 140–149.
- [44] Stéphane Weiss, Pascal Urso, and Pascal Molli. 2009. Logoot: A scalable optimistic replication algorithm for collaborative editing on p2p networks. In *2009 29th IEEE International Conference on Distributed Computing Systems*. IEEE, 404–412.

A PROOF OF HELPER CLAIM FOR SECTION 5.1

CLAIM 1. Let (m_1, c_1) , (m_2, c_2) , and (m_3, c_3) be elements from the Causal \langle COMPDOTFUN \rangle set. Then:

$$\begin{aligned}
& ((m_1, c_1) \sqcup (m_2, c_2)) \sqcup (m_3, c_3) = \\
& (\{k \mapsto v(k) \mid k \in \text{dom } m_1 \cap \text{dom } m_2 \cap \text{dom } m_3\} \\
& \cup \{k \mapsto \text{fst}((m_1(k), c_1) \sqcup (m_3(k), c_3)) \mid \\
& \quad k \in \text{dom } m_1 \cap \text{dom } m_3 \wedge k \notin c_2\} \\
& \cup \{k \mapsto \text{fst}((m_2(k), c_2) \sqcup (m_3(k), c_3)) \mid \\
& \quad k \in \text{dom } m_2 \cap \text{dom } m_3 \wedge k \notin c_1\} \\
& \cup \{k \mapsto \text{fst}((m_1(k), c_1) \sqcup (m_2(k), c_2)) \mid \\
& \quad k \in \text{dom } m_1 \cap \text{dom } m_2 \wedge k \notin c_3\} \\
& \cup \{k \mapsto m_3(k) \mid k \in \text{dom } m_3 \wedge k \notin c_1 \cup c_2\} \\
& \cup \{k \mapsto m_2(k) \mid k \in \text{dom } m_2 \wedge k \notin c_1 \cup c_3\} \\
& \cup \{k \mapsto m_1(k) \mid k \in \text{dom } m_1 \wedge k \notin c_2 \cup c_3\}, c_1 \cup c_2 \cup c_3) \\
& \textbf{where } v(k) = \text{fst}(((m_1(k), c_1) \sqcup (m_2(k), c_2)) \sqcup (m_3(k), c_3)).
\end{aligned}$$

PROOF. Consider (m_1, c_1) , (m_2, c_2) and (m_3, c_3) .

$$\begin{aligned}
(m', c') & \triangleq (m_1, c_1) \sqcup (m_2, c_2) \\
\text{dom } m' & = (\text{dom } m_1 \cap \text{dom } m_2) \cup \text{dom } m_1 \setminus c_2 \cup \text{dom } m_2 \setminus c_1
\end{aligned}$$

Consider $(m', c') \sqcup (m_3, c_3)$. We consider each set separately:

$$\begin{aligned}
& \{k \mapsto \text{fst}((m'(k), c') \sqcup (m_3(k), c_3)) \mid k \in \text{dom } m' \cap \text{dom } m_3\} = \\
& \{k \mapsto v(k) \mid k \in \text{dom } m_1 \cap \text{dom } m_2 \cap \text{dom } m_3\} \\
& \cup \{k \mapsto \text{fst}((m_1(k), c_1) \sqcup (m_3(k), c_3)) \mid k \in \text{dom } m_1 \setminus c_2 \cap \text{dom } m_3\} \\
& \cup \{k \mapsto \text{fst}((m_2(k), c_2) \sqcup (m_3(k), c_3)) \mid k \in \text{dom } m_2 \setminus c_1 \cap \text{dom } m_3\} = \\
& \{k \mapsto v(k) \mid k \in \text{dom } m_1 \cap \text{dom } m_2 \cap \text{dom } m_3\} \\
& \cup \{k \mapsto \text{fst}((m_1(k), c_1) \sqcup (m_3(k), c_3)) \mid k \in \text{dom } m_1 \cap \text{dom } m_3 \wedge k \notin c_2\} \\
& \cup \{k \mapsto \text{fst}((m_2(k), c_2) \sqcup (m_3(k), c_3)) \mid k \in \text{dom } m_2 \cap \text{dom } m_3 \wedge k \notin c_1\} \\
& \{(d, v) \in m' \mid k \notin c_3\} = \\
& \{k \mapsto m'(k) \mid k \in \text{dom } m' \wedge k \notin c_3\} = \\
& \{k \mapsto m'(k) \mid k \in \text{dom } m' \wedge k \notin c_3\} = \\
& \{k \mapsto \text{fst}((m_1(k), c_1) \sqcup (m_2(k), c_2)) \mid k \in \text{dom } m_1 \cap \text{dom } m_2 \wedge k \notin c_3\} \\
& \cup \{k \mapsto m_1(k) \mid k \in \text{dom } m_1 \wedge k \notin c_2 \wedge k \notin c_3\} \\
& \cup \{k \mapsto m_2(k) \mid k \in \text{dom } m_2 \wedge k \notin c_1 \wedge k \notin c_3\} = \\
& \{k \mapsto \text{fst}((m_1(k), c_1) \sqcup (m_2(k), c_2)) \mid k \in \text{dom } m_1 \cap \text{dom } m_2 \wedge k \notin c_3\} \\
& \cup \{k \mapsto m_1(k) \mid k \in \text{dom } m_1 \wedge k \notin c_2 \cup c_3\} \\
& \cup \{k \mapsto m_2(k) \mid k \in \text{dom } m_2 \wedge k \notin c_1 \cup c_3\} \\
& \{(d, v) \in m_3 \mid k \notin c'\} = \{k \mapsto m_3(k) \mid k \in \text{dom } m_3 \wedge k \notin c_1 \cup c_2\}
\end{aligned}$$

Finally:

$$c' \cup c_3 = c_1 \cup c_2 \cup c_3$$

□

B ANTI-ENTROPY ALGORITHM FOR CAUSAL CONSISTENCY

An anti-entropy algorithm as described by Alemida et al. in [17]. The algorithm ensures eventual convergence and the same causal-consistency guarantees that standard state-based CRDTs offer.

Definition B.1 (Delta-interval). Given a replica i progressing along the states X_i^0, X_i^1, \dots , by joining delta d_i^k (either local delta-mutation or received delta-group) into X_i^k to obtain X_i^{k+1} , a delta-interval $\Delta_i^{a,b}$ is a delta-group resulting from joining deltas d_i^a, \dots, d_i^{b-1} :

$$\Delta_i^{a,b} = \bigsqcup \{d_i^k \mid a \leq k < b\}$$

Each node i keeps a contiguous sequence of deltas d_i^l, \dots, d_i^u in a map D from integers to deltas, with $l = \min(\text{dom}(D))$ and $u = \max(\text{dom}(D))$. The sequence numbers of deltas are obtained from the counter c_i that is incremented when a delta (whether a delta-mutation or delta-interval received) is joined with the current state. Each node i keeps an acknowledgments map A that stores, for each neighbor j , the largest index b for all delta-intervals $\Delta_i^{a,b}$ acknowledged by j (after j receives $\Delta_i^{a,b}$ from i and joins it into X_j).

Node i sends a delta-interval $d = \Delta_i^{a,b}$ with a *(delta, d, b)* message; the receiving node j , after joining $\Delta_i^{a,b}$ into its replica state, replies with an acknowledgment message *(ack, b)*; if an ack from j was successfully received by node i , it updates the entry of j in the acknowledgment map, using the max function. This handles possible old duplicates and messages arriving out of order.

The counter c_i is also kept in a durable storage. This is essential to avoid conflicts after potential crash and recovery incidents. Otherwise, there would be the danger of receiving some delayed ack, for a delta-interval sent before crashing, causing the node to skip sending some deltas generated after recovery, thus violating the delta-merging condition.

The algorithm for node i periodically picks a random neighbor j . In principle, i sends the join of all deltas starting from the latest delta acked by j and forward. Exceptionally, i sends the entire state in two cases: (1) if the sequence of deltas D_i is empty, or (2) if j is expecting from i a delta that was already removed from D_i (e.g., after a crash and recovery, when both deltas and the ack map, being volatile state, are lost); i tracks this in $A_i(j)$. To garbage collect old deltas, the algorithm periodically removes the deltas that have been acked by *all* neighbors.

<pre> 1 inputs: 2 $n_i \in P(\mathbb{I})$, set of neighbors 3 durable state: 4 $X_i \in S$, CRDT state; initially $X_i = \perp$ 5 $c_i \in \mathbb{N}$, sequence number; initially $c_i = 0$ 6 volatile state: 7 $D_i \in \mathbb{N} \leftrightarrow S$, sequence of deltas; initially $D_i = \{\}$ 8 $A_i \in \mathbb{I} \leftrightarrow \mathbb{N}$, acknowledges map; initially $A_i = \{\}$ 9 on receive$_{j,i}$(<i>delta, d, n</i>) 10 if $d \not\sqsubseteq X_i$ then 11 $X'_i = X_i \sqcup d$ 12 $D'_i = D_i \{c_i \mapsto d\}$ 13 $c'_i = c_i + 1$ 14 send$_{i,j}$(<i>ack, n</i>) 15 on receive$_{j,i}$(<i>ack, n</i>) 16 $A'_i = A_i \{j \mapsto \max(A_i(j), n)\}$ </pre>	<pre> 17 on operation$_i$(m^δ) 18 $d = m^\delta(X_i)$ 19 $X'_i = X_i \sqcup d$ 20 $D'_i = D_i \{c_i \mapsto d\}$ 21 $c'_i = c_i + 1$ 22 periodically // ship delta-interval or state 23 $j = \text{random}(n_i)$ 24 if $D_i = \{\} \vee \min(\text{dom}(D_i)) > A_i(j)$ then 25 $d = X_i$ 26 else 27 $d = \bigsqcup \{D_i(l) \mid A_i(j) \leq l < c_i\}$ 28 if $A_i(j) < c_i$ then 29 send$_{i,j}$(<i>delta, d, c_i</i>) 30 periodically // garbage collect deltas 31 $l = \min\{n \mid (_, n) \in A_i\}$ 32 $D'_i = \{(n, d) \in D_i \mid n \geq l\}$ </pre>
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Algorithm 1: Anti-entropy algorithm ensuring causal consistency of δ -CRDT.