
Efficient State-based CRDTs by Delta-Mutation

Paulo Sérgio ALMEIDA,
Ali SHOKER,
Carlos BAQUERO.

HASLab, INESC Tec &
Universidade do Minho
Braga, Portugal

{psa,shokerali,cbm}@di.uminho.pt

Abstract

CRDTs are distributed data types that make eventual consistency of a distributed object possible and non ad-hoc. Specifically, state-based CRDTs achieve this by sharing local state changes through shipping the entire state, that is then merged to other replicas with an idempotent, associative, and commutative join operation, ensuring convergence. This imposes a large communication overhead as the state size becomes larger. We introduce *Delta State Conflict-Free Replicated Datatypes* (δ -CRDT), which make use of δ -mutators, defined in such a way to return a *delta-state*, typically, with a much smaller size than the full state. Delta-states are joined to the local state as well as to the remote states (after being shipped). This can achieve the best of both worlds: small messages with an incremental nature, as in operation-based CRDTs, disseminated over unreliable communication channels, as in traditional state-based CRDTs. We introduce the δ -CRDT framework, and we explain it through establishing a correspondence to current state-based CRDTs. In addition, we present two anti-entropy algorithms: a basic one that provides eventual convergence, and another one that ensures both convergence and causal consistency. We also introduce two δ -CRDT specifications of well-known replicated datatypes.

May 16, 2014

Efficient State-based CRDTs by Delta-Mutation

Paulo Sérgio Almeida, Ali Shoker, and Carlos Baquero

HASLab/INESC TEC and Universidade do Minho, Portugal

1 Introduction

Eventual consistency (EC) is a relaxed consistency model that is often adopted by large-scale distributed systems [1,2,3,4] where availability must be maintained, despite outages and partitioning, whereas delayed consistency is acceptable. The limitations resulting from the CAP theorem [5] suggest trading strong consistency for high availability; a *like/unlike* action in social networks is a concrete example. A typical approach in EC systems is to allow replicas of a distributed object to temporarily diverge, provided that they can eventually be reconciled into a common consistent state. To avoid application-specific reconciliation methods, costly and error-prone, *Conflict-Free Replicated Data Types* (CRDTs) [6,7] were introduced, allowing the design of self-contained distributed data types that are always available and eventually converge when all operations are reflected at all replicas. Though CRDTs have been successfully deployed in practice [1], a lot of work is still required to improve their design and performance.

CRDTs support two complementary designs: *state-based* which disseminate object states and *operation-based* which disseminate operations [6,7]. In a state-based design [8,7] an operation is only executed on the local replica state. A replica periodically propagates its local changes to other replicas through shipping its entire state. A received state is incorporated with the local state via a *merge* function that deterministically reconciles both states. To maintain convergence, *merge* is defined as a *join*: a least upper bound over a join-semilattice [8,7].

A major drawback in current state-based CRDTs is the communication overhead of shipping the entire state, which can get very large in size. For instance, the state size of a *counter* CRDT (a vector of integer counters, one per replica) increases with the number of replicas; whereas in a *grow-only Set*, the state size depends on the set size, that grows as more operations are invoked. This communication overhead limits the use of state-based CRDTs to data-types with small state size (e.g., counters are reasonable while sets are not); recently there has been demand for CRDTs with large state sizes (e.g., in RIAK DT Maps [9] that can compose multiple CRDTs).

In this paper, we rethink the way state-based CRDTs should be designed, having in mind the problematic shipping of the entire state. Our aim is to ship a representation of the effect of recent update operations on the state, rather than the whole state, while preserving the idempotent nature of *join*; thus, allowing unreliable communication, on the contrary to operation-based CRDTs that demand exactly-once delivery. To achieve this, we introduce *Delta State-based CRDTs* (δ -CRDT): a state is a join-semilattice that results from the join of multiple fine-grained states, i.e., *deltas*, generated by what we call δ -mutators; these are new versions of the datatype mutators that return the effect of these mutators on the state. In this way, deltas can be retained in a buffer to be shipped individually (or joined in groups) instead of shipping the entire object. The changes to the local state are then incorporated at other replicas by joining the shipped deltas with their own states.

A key point in our approach is a simple equation relating the novel δ -mutators with the original CRDT mutators. The challenge when designing a new δ -CRDT that corresponds to an existing CRDT is to derive δ -mutators that obey this equation. In this paper, we prove that eventual consistency is guaranteed in δ -CRDT as long as all deltas produced by δ -mutators are delivered

and joined at other replicas, and we present a corresponding simple anti-entropy algorithm. We then focus on causal consistency, introducing the concept of *delta-interval* and the *causal delta-merging condition*. Based on these, we then present an anti-entropy algorithm for δ -CRDT, where sending and then joining delta-intervals into another replica state produces the same effect as if the entire state had been shipped and joined. We illustrate our approach by introducing two δ -CRDT specifications of well-known datatypes: Counters, and Optimized Add-Wins Observed-Remove Sets.

1.1 System Model

Consider a distributed system with nodes containing local memory, with no shared memory between them. Any node can send messages to any other node. The network is asynchronous, there being no global clock, no bound on the time it takes for a message to arrive, nor bounds on relative processing speeds. The network is unreliable: messages can be lost, duplicated or reordered (but are not corrupted). Some messages will, however, eventually get through: if a node sends infinitely many messages to another node, infinitely many of these will be delivered. In particular, this means that there can be arbitrarily long partitions, but these will eventually heal. Nodes have access to durable storage; nodes can crash but eventually will recover with the content of the durable storage as at the time of the crash. Durable state is written atomically at each state transition. Each node has access to its globally unique identifier in a set \mathbb{I} .

2 A Background of State-based CRDTs

Conflict-Free Replicated Data Types [6,7] (CRDTs) are distributed datatypes that allow different replicas of a distributed CRDT instance to diverge and ensures that, eventually, all replicas converge to the same state. State-based CRDTs achieve this through propagating updates of the local state by disseminating the entire state across replicas. The received states are then merged to remote states, leading to convergence.

A state-based CRDT consists of a triple (S, M, Q) , where S is a join-semilattice [10], Q is a set of query functions (which return some result without modifying the state), and M is a set of mutators that perform updates; a mutator $m \in M$ takes a state $X \in S$ as input and returns a new state $X' = m(X)$. A join-semilattice is a set with a *partial order* \sqsubseteq and a binary *join* operation \sqcup that returns the *least upper bound* (LUB) of two elements in S ; a *join* is designed to be commutative, associative, and idempotent. Mutators are defined in such a way to be *inflations*, i.e., for any mutator m and state X , the following holds:

$$X \sqsubseteq m(X)$$

In this way, for each replica there is a monotonic sequence of states, defined under the lattice partial order, where each subsequent state subsumes the previous state when joined elsewhere.

Both query and mutator operations are always available since they are performed using the local state without requiring inter-replica communication; however, as mutators are concurrently applied at distinct replicas, replica states will likely diverge. Eventual convergence is then obtained using an *anti-entropy* protocol that periodically ships the entire local state to other replicas. Each replica merges the received state with its local state using the *join* operation in S . Given the mathematical properties of *join*, if mutators stop being issued, all replicas eventually converge to the same state. i.e. the least upper-bound of all states involved. State-based CRDTs are interesting as they demand little guarantees from the dissemination layer, working under message loss, duplication, reordering, and temporary network partitioning, without impacting availability and eventual convergence.

Fig. 1 represents a state-based increment-only counter. The CRDT state Σ is a map from replica identifiers to positive integers. Initially, σ_i^0 is an empty map (with the assumption that unmapped keys implicitly map to zero, and only non zero mappings are stored). Only one mutator, i.e., `inc`, is defined which increments the value corresponding to the local replica i (returning the updated map). The query operation `value` returns the counter value by adding the integers in the map entries. The join of two states is the point-wise maximum of the maps.

The main weakness of state-based CRDTs is the cost of dissemination of updates, as the full state is sent. In this simple example of counters, even though increments only update the value corresponding to the local replica i , the whole map will always be sent in messages though the other map values remained intact (since no messages have been received and merged).

It would be interesting to only ship the recent modification incurred on the state. This is, however, not possible with the current model of state-based CRDTs as mutators always return a full state. Approaches which simply ship operations (e.g., an “increment n ” message), like in operation-based CRDTs, require reliable communication (e.g., because increment is not idempotent). In contrast, our approach allows producing and encoding recent mutations in an incremental way, while keeping the advantages of the state-based approach, namely the idempotent, associative, and commutative properties of join.

3 Delta-state CRDTs

We introduce *Delta-State Conflict-Free Replicated Data Types*, or δ -CRDT for short, as a new kind of state-based CRDTs, in which *delta-mutators* are defined to return a *delta-state*: a value in the same join-semilattice which represents the updates induced by the mutator on the current state.

Definition 1 (Delta-mutator). A *delta-mutator* m^δ is a function, corresponding to an update operation, which takes a state X in a join-semilattice S as parameter and returns a *delta-mutation* $m^\delta(X)$, also in S .

Definition 2 (Delta-group). A *delta-group* is inductively defined as either a *delta-mutation* or a *join of several delta-groups*.

Definition 3 (δ -CRDT). A δ -CRDT consists of a triple (S, M^δ, Q) , where S is a join-semilattice, M^δ is a set of *delta-mutators*, and Q a set of *query functions*, where the state transition at each replica is given by either joining the current state $X \in S$ with a *delta-mutation*:

$$X' = X \sqcup m^\delta(X),$$

or joining the current state with some received *delta-group* D :

$$X' = X \sqcup D.$$

In a δ -CRDT, the effect of applying a mutation, represented by a *delta-mutation* $\delta = m^\delta(X)$, is decoupled from the resulting state $X' = X \sqcup \delta$, which allows shipping this δ rather than the entire

$$\begin{aligned} \Sigma &= \mathbb{I} \leftrightarrow \mathbb{N} \\ \sigma_i^0 &= \{\} \\ \text{inc}_i(m) &= m\{i \mapsto m(i) + 1\} \\ \text{value}_i(m) &= \sum_{i \in \mathbb{I}} m(i) \\ m \sqcup m' &= \{(i, \max(m(i), m'(i))) \mid i \in \mathbb{I}\} \end{aligned}$$

Fig. 1: State-based Counter CRDT; replica i .

resulting state X' . All state transitions in a δ -CRDT, even upon applying mutations locally, are the result of some join with the current state. Unlike standard CRDT mutators, delta-mutators do not need to be inflations in order to inflate a state; this is however ensured by joining their output, i.e., deltas, into the current state.

In principle, a delta could be shipped immediately to remote replicas once applied locally. For efficiency reasons, multiple deltas returned by applying several delta-mutators can be joined locally into a delta-group and retained in a buffer. The delta-group can then be shipped to remote replicas to be joined with their local states. Received delta-groups can optionally be joined into their buffered delta-group, allowing transitive propagation of deltas. A full state can be seen as a special (extreme) case of a delta-group.

If the causal order of operations is not important and the attended aim is merely eventual convergence of states, then delta-groups can be shipped using an unreliable dissemination layer that may drop, reorder, or duplicate messages. Delta-groups can always be re-transmitted and re-joined, possibly out of order, or can simply be subsumed by a less frequent sending of the full state, e.g. for performance reasons or when doing state transfers to new members. In Section 4, we address state convergence when causal consistency is not required, and we address the latter in Section 5.

3.1 Delta-state decomposition of standard CRDTs

A δ -CRDT (S, M^δ, Q) is a *delta-state decomposition* of a state-based CRDT (S, M, Q) , if for every mutator $m \in M$, we have a corresponding mutator $m^\delta \in M^\delta$ such that, for every state $X \in S$:

$$m(X) = X \sqcup m^\delta(X)$$

This equation states that applying a delta-mutator and joining into the current state should produce the same state transition as applying the corresponding mutator of the standard CRDT.

Given an existing state-based CRDT (which is always a trivial decomposition of itself, i.e., $m(X) = X \sqcup m(X)$, as mutators are inflations), it will be useful to find a non-trivial decomposition such that delta-states returned by delta-mutators in M^δ are smaller than the resulting state:

$$\text{size}(m^\delta(X)) \ll \text{size}(m(X))$$

3.2 Example: δ -CRDT Counter

Fig. 2 depicts a δ -CRDT specification of a counter datatype that is a delta-state decomposition of the state-based counter in Fig. 1. The state, join and value query operation remain as before. Only the mutator inc^δ is newly defined, which increments the map entry corresponding to the local replica and only returns that entry, instead of the full map as inc in the state-based CRDT counter does. This maintains the original semantics of the counter while allowing the smaller deltas returned by the delta-mutator to be sent, instead of the full map. As before, the received payload (whether one or more deltas) might not include entries for all keys in \mathbb{I} , which

$$\begin{aligned} \Sigma &= \mathbb{I} \hookrightarrow \mathbb{N} \\ \sigma_i^0 &= \{\} \\ \text{inc}_i^\delta(m) &= \{i \mapsto m(i) + 1\} \\ \text{value}_i(m) &= \sum_{i \in \mathbb{I}} m(i) \\ m \sqcup m' &= \{(i, \max(m(i), m'(i))) \mid i \in \mathbb{I}\} \end{aligned}$$

Fig. 2: A δ -CRDT counter; replica i .

are assumed to have zero values. The decomposition is easy to understand in this example since the equation $\text{inc}_i(X) = X \sqcup \text{inc}_i^\delta(X)$ holds as $m\{i \mapsto m(i) + 1\} = m \sqcup \{i \mapsto m(i) + 1\}$. In other words, the single value for key i in the delta, corresponding to the local replica identifier, will overwrite the corresponding one in m since the former maps to a higher value (i.e., using max). Here it can be noticed that a delta *is* just a state, that can be joined possibly several times without requiring exactly-once delivery, and without being a representation of the “increment” operation (as in operation-based CRDTs), which is itself non-idempotent.

4 State Convergence

In the δ -CRDT execution model, and regardless of the anti-entropy algorithm used, a replica state always evolves by joining the current state with some *delta*: either the result of a delta-mutation, or some arbitrary delta-group (which itself can be expressed as a join of delta-mutations). Therefore, all states can be expressed as joins of delta-mutations, which makes state convergence in δ -CRDT easy to achieve: it is enough that all delta-mutations generated in the system reach every replica, as expressed by the following proposition.

Proposition 1. (*δ -CRDT convergence*) *Consider a set of replicas of a δ -CRDT object, replica i evolving along a sequence of states $X_i^0 = \perp, X_i^1, \dots$, each replica performing delta-mutations of the form $m_{i,k}^\delta(X_i^k)$ at some subset of its sequence of states, and evolving by joining the current state either with self-generated deltas or with delta-groups received from others. If each delta-mutation $m_{i,k}^\delta(X_i^k)$ produced at each replica is joined (directly or as part of a delta-group) at least once with every other replica, all replica states become equal.*

Proof. Trivial, given the associativity, commutativity, and idempotence of the join operation in any join-semilattice.

This opens up the possibility of having anti-entropy algorithms that are only devoted to enforce convergence, without necessarily providing causal consistency (enforced in standard CRDTs); thus, making a trade-off between performance and consistency guarantees. For instance, in a counter (e.g., for the number of *likes* on a social network), it may not be critical to have causal consistency, but merely not to lose increments and achieve convergence.

4.1 Basic Anti-Entropy Algorithm

A basic anti-entropy algorithm that ensures eventual convergence in δ -CRDT is presented in Algorithm 1. For the node corresponding to replica i , the durable state, which persists after a crash, is simply the δ -CRDT state X_i . The volatile state D stores a delta-group that is used to accumulate deltas before eventually sending it to other replicas. Without loss of generality, we assume that the join-semilattice has a bottom \perp , which is the initial value for both X_i and D_i .

When an operation is performed, the corresponding delta-mutator m^δ is applied to the current state of X_i , generating a delta d . This delta is joined both with X_i to produce a new state, and with D . In the same spirit of standard state based CRDTs, a node sends its messages in a periodic fashion, where the message payload is either the delta-group D_i or the full state X_i ; this decision is made by the function choose_i which returns one of them. To keep the algorithm simple, a node simply broadcasts its messages without distinguishing between neighbors. After each send, the delta-group is reset to \perp .

Once a message is received, the payload d is joined into the current δ -CRDT state. The basic algorithm operates in two modes: (1) a *transitive* mode (when t_i is true) in which m is also joined

<pre> 1 inputs: 2 $n_i \in \mathcal{P}(\mathbb{I})$, set of neighbors 3 $t_i \in \mathbb{B}$, true for transitive mode 4 $\text{choose}_i \in S \times S \rightarrow S$, ship state or delta 5 durable state: 6 $X_i \in S$, CRDT state; initially $X_i = \perp$ 7 volatile state: 8 $D_i \in S$, join of deltas; initially $D_i = \perp$ 9 on operation$_i(m^\delta)$ 10 $d = m^\delta(X_i)$ 11 $X'_i = X_i \sqcup d$ 12 $D'_i = D_i \sqcup d$ </pre>	<pre> 13 periodically 14 $m = \text{choose}_i(X_i, D_i)$ 15 for $j \in n_i$ do 16 $\text{send}_{i,j}(m)$ 17 $D'_i = \perp$ 18 on receive$_{j,i}(d)$ 19 $X'_i = X_i \sqcup d$ 20 if t_i then 21 $D'_i = D_i \sqcup d$ 22 else 23 $D'_i = D_i$ </pre>
--	--

Algorithm 1: Basic anti-entropy algorithm for δ -CRDT.

into D , allowing transitive propagation of delta-mutations; meaning that, deltas received at node i from some node j can later be sent to some other node k ; (2) a *direct* mode where a delta-group is exclusively the join of local delta-mutations (j must send its deltas directly to k). The decisions of whether to send a delta-group versus the full state (typically less periodically), and whether to use the transitive or direct mode are out of the scope of this paper. In general, decisions can be made considering many criteria like delta-groups size, state size, message loss distribution assumptions, and network topology.

5 Causal Consistency

For some CRDTs with commutative operations, like the counter in Fig. 2, eventual convergence of states may be enough, and thus any anti-entropy algorithm that satisfies the condition in Proposition 1, like Algorithm 1, can be used. In general, users of CRDTs likely want more than eventual convergence, namely causal consistency. An often used example is: some user removes the boss from the set allowed to read some folder and then adds some picture to the same folder. It may be undesirable that the picture is considered to be in one set, while the boss is still considered to be in the other. When using an anti-entropy mechanism which disseminates deltas with no order guarantees (like Algorithm 1) the execution is, in general, not causally consistent. In the above example, the boss could happen to see the picture, even if a single δ -CRDT containing the two sets were used.

Traditional state-based CRDTs converge using joins of the full state, which implicitly ensures per-object causal consistency [11]: each state of some replica of an object reflects the causal past of operations on the object (either applied locally, or applied at other replicas and transitively joined).

Therefore, it is desirable to have δ -CRDTs offer the same causal-consistency guarantees that standard state-based CRDTs offer. This raises the question about how can delta propagation and merging of δ -CRDT be constrained (and expressed in an anti-entropy algorithm) in such a manner to give the same results as if a standard state-based CRDT was used. Towards this objective, it is useful to define a particular kind of delta-group, which we call a *delta-interval*:

Definition 4 (Delta-interval). *Given a replica i evolving along 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} :*

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

The use of delta-intervals in anti-entropy algorithms will be a key ingredient towards achieving causal consistency. We now define a restricted kind of anti-entropy algorithms for δ -CRDTs.

Definition 5 (Delta-interval-based anti-entropy algorithm). *An anti-entropy algorithm for δ -CRDTs is delta-interval-based, if all deltas sent to other replicas are delta-intervals.*

The other ingredient towards achieving causal consistency is given by the following condition:

Definition 6 (Causal delta-merging condition). *A delta-interval based anti-entropy algorithm is said to satisfy the causal delta-merging condition if the algorithm only joins $\Delta_j^{a,b}$ from replica j into replica i states X_i that satisfy:*

$$X_i \supseteq X_j^a.$$

This means that a delta-interval is only joined into states that at least reflect (i.e., subsume) the state into which the first delta in the interval was previously joined. The causal delta-merging condition is important since any delta-interval based anti-entropy algorithm of a δ -CRDT that satisfies it, can be used to obtain the same outcome of standard CRDTs; this is more formally stated in the following proposition.

Proposition 2. (CRDT and δ -CRDT correspondence) *Let (S, M, Q) be a standard state-based CRDT and (S, M^δ, Q) a corresponding delta-state decomposition. Any δ -CRDT state reachable by an execution E^δ over (S, M^δ, Q) , by a delta-interval based anti-entropy algorithm A^δ satisfying the causal delta-merging condition, is equal to a state resulting from an execution E over (S, M, Q) , having the corresponding data-type operations, by an anti-entropy algorithm A for state-based CRDTs.*

Proof. See appendix.

Corollary 1. (δ -CRDT causal consistency) *Any δ -CRDT in which states are propagated and joined using a delta-interval-based anti-entropy algorithm satisfying the causal delta-merging condition ensures causal consistency.*

Proof. From Proposition 2 and causal consistency of state-based CRDTs.

5.1 Anti-Entropy Algorithm for Causal Consistency

Algorithm 2 is a delta-interval based anti-entropy algorithm which enforces the causal delta-merging condition. It can be used whenever the causal consistency guarantees of standard state-based CRDTs are needed. For presentation purposes, it excludes some optimizations that are important in practice, but easy to derive. The algorithm distinguishes neighbor nodes, and only sends them delta-intervals that are joined at the receiving node, obeying the delta-merging condition.

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 $l = \min(\text{dom}(M))$ and $u = \max(\text{dom}(M))$. 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 acknowledgements map A that stores, for each neighbor j , the index b such that $\Delta_i^{a,b}$ is the last delta-interval acknowledged by j (after it 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 acknowledgement message (ack, b); if an ack from j was successfully received by node i , it updates the entry of j in the acknowledgement map, using the max function. This handles possible old duplicates and messages arriving out of order.

```

1 inputs:
2    $n_i \in \mathcal{P}(\mathbb{I})$ , 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 $j,i$ (delta,  $d, n$ )
10   $X'_i = X_i \sqcup d$ 
11   $D'_i = D_i \{c_i \mapsto d\}$ 
12   $c'_i = c_i + 1$ 
13  send $i,j$ (ack,  $n$ )
14 on receive $j,i$ (ack,  $n$ )
15   $A'_i = A_i \{j \mapsto \max(A_i(j), n)\}$ 
16 on operation $i$ ( $m^\delta$ )
17   $d = m^\delta(X_i)$ 
18   $X'_i = X_i \sqcup d$ 
19   $D'_i = D_i \{c_i \mapsto d\}$ 
20   $c'_i = c_i + 1$ 
21 periodically // ship delta-interval or state
22   $j = \text{random}(n_i)$ 
23  if  $D_i = \{\} \vee \min(\text{dom}(D_i)) > A_i(j)$  then
24     $d = X_i$ 
25  else
26     $d = \bigsqcup \{D_i(l) \mid A_i(j) \leq l < c_i\}$ 
27  send $i,j$ (delta,  $d, c_i$ )
28 periodically // garbage collect deltas
29   $l = \min\{n \mid (\cdot, n) \in A_i\}$ 
30   $D'_i = \{(n, d) \in D_i \mid n \geq l\}$ 

```

Algorithm 2: An anti-entropy algorithm that ensures causal consistency of δ -CRDT.

Like the δ -CRDT state, 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, make the node 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 delta that j acked 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); 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.

Proposition 3. *Algorithm 2 produces the same reachable states as a standard algorithm over a CRDT for which the δ -CRDT is a decomposition.*

Proof. See appendix.

6 δ -CRDTs for Add-Wins OR-Sets

An Add-wins Observed-Remove Set is a well-known CRDT datatype that offers the same sequential semantics of a sequential set and adopts a specific resolution semantics for operations that concurrently add and remove the same element. Add-wins means that an add prevails over a concurrent remove. Remove operations, however, only affect elements added by causally preceding adds.

Fig. 3a depicts a simple, but inefficient, δ -CRDT implementation of a state-based add-wins OR-Set. The state Σ consists of a set of tagged elements and a set of tags, acting as tombstones. Globally unique tags of the form $\mathbb{I} \times \mathbb{N}$ are used and ensured by pairing a replica identifier in \mathbb{I} with a monotonically increasing natural counter. Once an element $e \in E$ is added to the set, the delta-mutator add^δ creates a globally unique tag by incrementing the highest tag present in its local state and that was created by replica i itself (max returns 0 if no tag is present). This tag is paired with value e and stored as a new unique triple. Since removes should only tombstone elements that are added before the remove operation, the delta-mutator rmv^δ retains in the tombstone set all tags associated to element e , being removed from the local state. Function elements only returns the elements that are added but not tombstoned yet. Join \sqcup simply unions the respective sets that are, therefore, both grow-only.

$\Sigma = \mathcal{P}(\mathbb{I} \times \mathbb{N} \times E) \times \mathcal{P}(\mathbb{I} \times \mathbb{N})$ $\sigma_i^0 = (\{\}, \{\})$ $\text{add}_i^\delta(e, (s, t)) = (\{(i, \max(\{n \mid (i, n, -) \in s\}) + 1, e)\}, \{\})$ $\text{rmv}_i^\delta(e, (s, t)) = (\{\}, \{(j, n) \mid (j, n, e) \in s\})$ $\text{elements}_i((s, t)) = \{e \mid (j, n, e) \in s \wedge (j, n) \notin t\}$ $(s, t) \sqcup (s', t') = (s \cup s', t \cup t')$	$\Sigma = \mathcal{P}(\mathbb{I} \times \mathbb{N} \times E) \times \mathcal{P}(\mathbb{I} \times \mathbb{N})$ $\sigma_i^0 = (\{\}, \{\})$ $\text{add}_i^\delta(e, (s, c)) = (\{(i, n, e)\}, \{(i, n)\})$ <p style="text-align: center; margin-left: 100px;">where $n = 1 + \max(\{k \mid (i, k) \in c\})$</p> $\text{rmv}_i^\delta(e, (s, c)) = (\{\}, \{(j, n) \mid (j, n, e) \in s\})$ $\text{elements}_i((s, c)) = \{e \mid (j, n, e) \in s\}$ $(s, c) \sqcup (s', c') = ((s \cap s') \cup \{(i, n, e) \in s \mid (i, n) \notin c'\} \cup \{(i, n, e) \in s' \mid (i, n) \notin c\}, c \cup c')$
(a) With Tombstones	(b) Without Tombstones (optimized)

Fig. 3: Add-wins observed-remove δ -CRDT set, replica i .

A more efficient design is presented in Fig. 3b, which offers the same semantics and have a similar state structure; however, it uses a different join-semilattice, allowing the set of tagged elements to shrink as elements are removed. Now, `elements` returns all elements in the tagged set s . Instead of the tombstone set, a *causal context set* is used. Adding an element creates a unique tag by resorting to the causal context c instead of s , that can now shrink; the new triple is added to s as before, but now the new tag is also added to causal context c . The delta-mutator rmv^δ is the same as before, collecting all tags associated to the element being removed. The desired semantics are achieved by the novel join operation \sqcup . To join two states, their causal contexts are simply unioned; whereas, the new tagged element set only preserves: (1) the triples present in both sets (therefore, not removed in either), and also (2) any triple present in one of the sets and whose tag is not present in the causal context of the other state. The other elements are simply discarded.

Causal Context Compression. For presentation simplicity, this optimized version of the add-wins OR-Set has a grow-only causal context that collects all the unique tags (even from elements added but no longer present). In practice the causal context can be efficiently compressed without any loss of information. When using an anti-entropy algorithm that provides causal consistency, e.g., Algorithm 2, then for each replica state $X_i = (s_i, c_i)$ and replica identifier $j \in \mathbb{I}$, we have a contiguous sequence:

$$1 \leq n \leq \max(\{k \mid (j, k) \in c_i\}) \Rightarrow (j, n) \in c_i.$$

Thus, the causal context can always be encoded as a compact version vector [12] $\mathbb{I} \leftrightarrow \mathbb{N}$ that keeps the maximum sequence number for each replica. Even under non-causal anti-entropy, compression is still possible by keeping a version vector that encodes the offset of the contiguous sequence of tags from each replica, together with a set for the non-contiguous tags. As anti-entropy proceeds, each tag is eventually encoded in the vector, and thus the set remains typically small. Compression is less likely for the causal context of delta-groups in transit or buffered to be sent, but those contexts are only transient and smaller than those in the actual replica states. Moreover, the same techniques that encode contiguous sequences of tags can also be used for transient context compression [13].

7 Related Work

Eventually convergent data types. The design of replicated systems that are always available and eventually converge can be traced back to historical designs in [14,15], among others. More recently, replicated data types that always eventually converge, both by reliably broadcasting operations

(called operation-based) or gossiping and merging states (called state-based), have been formalized as CRDTs [16,8,6,7]. These are also closely related to Bloom^L [17] and Cloud Types [18]. State semi-lattices were used for deterministic parallel programming in LVars [19], where variables progress in the lattice order by joining other values, and are only accessible by special threshold reads.

Message size. A key feature of δ -CRDT is message size reduction and coalescing, using small-sized deltas. A different sort of state-based deltas was introduced for Computational CRDTs [20]; however, it requires a special merge for deltas (in addition to the standard join) that does not guarantee idempotence. Operation-based CRDTs [6,7,21] also support small message sizes, and in particular, *pure* flavours [21] that restrict messages to the operation name, and possible arguments. Though pure operation-based CRDTs allow for compact states and are very fast at the source (since operations are broadcast without consulting the local state), the model requires more systems guarantees than δ -CRDT do, e.g., exactly-once reliable delivery and membership information, and impose more complex integration of new replicas.

Encoding causal histories. State-based CRDT are always designed to be causally consistent [8,7]. Optimized implementations of sets, maps, and multi-value registers can build on this assumption to keep the meta-data small [11]. In δ -CRDT, however, deltas and delta-groups are normally not causally consistent, and thus the design of *join*, the meta-data state, as well as the anti-entropy algorithm used must ensure this. Without causal consistency, the causal context in δ -CRDT can not always be summarized with version vectors, and consequently, techniques that allow for gaps are often used. A well known mechanism that allows for encoding of gaps is found in Concise Version Vectors [22]. Interval Version Vectors [13], later on, introduced an encoding that optimizes sequences and allows gaps, while preserving efficiency when gaps are absent.

8 Conclusion

CRDTs allow flexible, while principled, design of distributed protocols that trade strict consistency for improved availability, faster response time, and support for disconnected operation. These benefits are harvested whenever a given application can model all, or a part, of its behavior using CRDTs. In particular, state-based CRDTs allow idempotent gossiping of states and require very basic guarantees from the network, as they cope with message loss, re-ordering, and duplication.

The price that comes with state-based CRDTs is that states have a potential to get very large. This problem is made worse if multiple objects are composed into a single CRDT object to benefit from intra-replica consistency and atomicity. In this paper, we addressed these limitations by introducing the new concept of δ -CRDT. By devising *delta-mutators* over state-based datatypes, it is now possible to detach the changes that an operation induces on the state. This brings a significant performance gain as it allows only shipping small states, i.e., *deltas*, instead of the entire state. The significant property in δ -CRDT is that it preserves the crucial properties (idempotence, associativity and commutativity) of standard state-based CRDT joins.

We have shown how δ -CRDT can achieve convergence possibly with, or without, causal consistency; and we presented an anti-entropy algorithm for each case. In particular, the causally consistent algorithm allows replacing classical state-based CRDTs by more efficient ones, while preserving their properties. As a first application of our approach, we designed a novel δ -CRDT specification for a well-known and widely used datatype: an optimized observed-remove set.

References

- [1] Cribbs, S., Brown, R.: Data structures in Riak. In: Riak Conference (RICON), San Francisco, CA, USA (oct 2012)
- [2] Bailis, P., Ghodsi, A.: Eventual consistency today: Limitations, extensions, and beyond. *Queue* **11**(3) (March 2013) 20:20–20:32
- [3] Terry, D.B., Theimer, M.M., Petersen, K., Demers, A.J., Spreitzer, M.J., Hauser, C.H.: Managing update conflicts in Bayou, a weakly connected replicated storage system. In: *Symp. on Op. Sys. Principles (SOSP)*, Copper Mountain, CO, USA, ACM SIGOPS, ACM Press (December 1995) 172–182
- [4] DeCandia, G., Hastorun, D., Jampani, M., Kakulapati, G., Lakshman, A., Pilchin, A., Sivasubramanian, S., Voshall, P., Vogels, W.: Dynamo: Amazon’s highly available key-value store. In: *Symp. on Op. Sys. Principles (SOSP)*. Volume 41 of *Operating Systems Review.*, Stevenson, Washington, USA, Assoc. for Computing Machinery (October 2007) 205–220
- [5] Gilbert, S., Lynch, N.: Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *SIGACT News* **33**(2) (2002) 51–59
- [6] Shapiro, M., Pregoça, N., Baquero, C., Zawirski, M.: A comprehensive study of Convergent and Commutative Replicated Data Types. *Rapp. Rech. 7506*, Institut National de la Recherche en Informatique et Automatique (INRIA), Rocquencourt, France (January 2011)
- [7] Shapiro, M., Pregoça, N., Baquero, C., Zawirski, M.: Conflict-free replicated data types. In Défago, X., Petit, F., Villain, V., eds.: *Int. Symp. on Stabilization, Safety, and Security of Distributed Systems (SSS)*. Volume 6976 of *Lecture Notes in Comp. Sc.*, Grenoble, France, Springer-Verlag (October 2011) 386–400
- [8] Baquero, C., Moura, F.: Using structural characteristics for autonomous operation. *Operating Systems Review* **33**(4) (1999) 90–96
- [9] Brown, R., Cribbs, S., Meiklejohn, C., Elliott, S.: Riak dt map: A composable, convergent replicated dictionary. In: *Proceedings of the First Workshop on Principles and Practice of Eventual Consistency. PaPEC ’14*, New York, NY, USA, ACM (2014) 1:1–1:1
- [10] Davey, B.A., Priestley, H.A.: *Introduction to Lattices and Order* (2. ed.). Cambridge University Press (2002)
- [11] Burckhardt, S., Gotsman, A., Yang, H., Zawirski, M.: Replicated data types: specification, verification, optimality. In Jagannathan, S., Sewell, P., eds.: *POPL*, ACM (2014) 271–284
- [12] Parker, D.S., Popek, G.J., Rudisin, G., Stoughton, A., Walker, B.J., Walton, E., Chow, J.M., Edwards, D., Kiser, S., Kline, C.: Detection of mutual inconsistency in distributed systems. *IEEE Trans. Softw. Eng.* **9**(3) (May 1983) 240–247
- [13] Mukund, M., R., G.S., Suresh, S.P.: Optimized or-sets without ordering constraints. In: *Proceedings of the International Conference on Distributed Computing and Networking*, New York, NY, USA, ACM (2014) 227–241
- [14] Wu, G.T.J., Bernstein, A.J.: Efficient solutions to the replicated log and dictionary problems. In: *Symp. on Principles of Dist. Comp. (PODC)*, Vancouver, BC, Canada (August 1984) 233–242
- [15] Johnson, P.R., Thomas, R.H.: The maintenance of duplicate databases. *Internet Request for Comments RFC 677*, Information Sciences Institute (January 1976)
- [16] Letia, M., Pregoça, N., Shapiro, M.: CRDTs: Consistency without concurrency control. *Rapp. Rech. RR-6956*, Institut National de la Recherche en Informatique et Automatique (INRIA), Rocquencourt, France (June 2009)

- [17] Conway, N., Marczak, W.R., Alvaro, P., Hellerstein, J.M., Maier, D.: Logic and lattices for distributed programming. In: Proceedings of the Third ACM Symposium on Cloud Computing, ACM (2012) 1
- [18] Burckhardt, S., Fähndrich, M., Leijen, D., Wood, B.P.: Cloud types for eventual consistency. In: ECOOP 2012–Object-Oriented Programming. Springer (2012) 283–307
- [19] Kuper, L., Newton, R.R.: Lvars: lattice-based data structures for deterministic parallelism. In: Proceedings of the 2nd ACM SIGPLAN workshop on Functional high-performance computing, ACM (2013) 71–84
- [20] Navalho, D., Duarte, S., Preguiça, N., Shapiro, M.: Incremental stream processing using computational conflict-free replicated data types. In: Proceedings of the 3rd International Workshop on Cloud Data and Platforms, ACM (2013) 31–36
- [21] Baquero, C., Almeida, P.S., Shoker, A.: Making operation-based CRDTs operation-based. In: to appear in Proceedings of Distributed Applications and Interoperable Systems: 14th IFIP WG 6.1 International Conference, Springer (2014)
- [22] Malkhi, D., Terry, D.: Concise version vectors in winfs. Distributed Computing **20**(3) (2007) 209–219

A Proof of Proposition 1

Proof. By simulation, establishing a correspondence between an execution E^δ , and execution E of a standard CRDT of which (S, M^δ, Q) is a decomposition, as follows: 1) the state (X_i, D_i, \dots) of each node in E^δ containing CRDT state X_i , information about delta-intervals D_i and possibly other information, corresponds to only X_i component (in the same join-semilattice); 2) for each action which is a delta-mutation m^δ in E^δ , E executes the corresponding mutation m , satisfying $m(X) = X \sqcup m^\delta(X)$; 3) whenever E^δ contains a send action of a delta-interval $\Delta_i^{a,b}$, execution E contains a send action containing the full state X_i^b ; 4) whenever E^δ performs a join into some X_i of a delta-interval $\Delta_j^{a,b}$, execution E delivers and joins the corresponding message containing the full CRDT state X_j^b . By induction on the length of the trace, assume that for each replica i , each node state X_i in E is equal to the corresponding component in the node state in E^δ , up to the last action in the global trace. A send action does not change replica state, preserving the correspondence. Replica states only change either by performing data-type update operations or upon message delivery by merging deltas/states respectively. If the next action is an update operation, the correspondence is preserved due to the delta-state decomposition property $m(X) = X \sqcup m^\delta(X)$. If the next action is a message delivery at replica i , with a merging of delta-interval/state from other replica j , because algorithm A^δ satisfies the causal merging-condition, it only joins into state X_i^k a delta-interval $\Delta_j^{a,b}$ if $X_i^k \sqsupseteq X_j^a$. In this case, the outcome will be:

$$\begin{aligned}
X_i^{k+1} &= X_i^k \sqcup \Delta_j^{a,b} \\
&= X_i^k \sqcup \bigsqcup \{d_j^k \mid a \leq k < b\} \\
&= X_i^k \sqcup X_j^a \sqcup \bigsqcup \{d_j^k \mid a \leq k < b\} \\
&= X_i^k \sqcup X_j^a \sqcup d_j^a \sqcup d_j^{a+1} \sqcup \dots \sqcup d_j^{b-1} \\
&= X_i^k \sqcup X_j^{a+1} \sqcup d_j^{a+1} \sqcup \dots \sqcup d_j^{b-1} \\
&= \dots \\
&= X_i^k \sqcup X_j^{b-1} \sqcup d_j^{b-1} \\
&= X_i^k \sqcup X_j^b
\end{aligned}$$

The resulting state X_i^{k+1} in E^δ will be, therefore, the same as the corresponding one in E where the full CRDT state from j has been joined, preserving the correspondence between E^δ and E .

B Proof of Proposition 3

Proof. From Proposition 1, it is enough to prove that the algorithm satisfies the causal delta-merging condition. The algorithm explicitly keeps deltas d_i^k tagged with increasing sequence numbers (even after a crash), according with the definition; node j only sends to i a delta-interval $\Delta_j^{a,b}$ if i has acked a ; this ack is sent only if i has already joined some delta-interval (possibly a full state) $\Delta_j^{k,a}$. Either $k = 0$ or, by the same reasoning, this $\Delta_j^{k,a}$ could only have been joined at i if some other interval $\Delta_j^{l,k}$ had already been joined into i . This reasoning can be recursed until a delta-interval starting from zero is reached. Therefore, $X_i \sqsupseteq \bigsqcup \{d_j^k \mid 0 \leq k < a\} = \Delta_j^{0,a} = X_j^a$.