

Automatically Verifying Replication-Aware Linearizability

VIMALA SOUNDARAPANDIAN, IIT Madras, India

KARTIK NAGAR, IIT Madras, India

ASEEM RASTOGI, Microsoft Research, India

KC SIVARAMAKRISHNAN, IIT Madras, India and Tarides, India

Data replication is crucial for enabling fault tolerance and uniform low latency in modern decentralized applications. Replicated Data Types (RDTs) have emerged as a principled approach for developing replicated implementations of basic data structures such as counter, flag, set, map, etc. While the correctness of RDTs is generally specified using the notion of strong eventual consistency—which guarantees that replicas that have received the same set of updates would converge to the same state—a more expressive specification which relates the converged state to updates received at a replica would be more beneficial to RDT users. Replication-aware linearizability is one such specification, which requires all replicas to always be in a state which can be obtained by linearizing the updates received at the replica. In this work, we develop a novel fully automated technique for verifying replication-aware linearizability for Mergeable Replicated Data Types (MRDTs). We identify novel algebraic properties for MRDT operations and the merge function which are sufficient for proving an implementation to be linearizable and which go beyond the standard notions of commutativity, associativity, and idempotence. We also develop a novel inductive technique called bottom-up linearization to automatically verify the required algebraic properties. Our technique can be used to verify both MRDTs and state-based CRDTs. We have successfully applied our approach to a number of complex MRDT and CRDT implementations including a novel JSON MRDT.

CCS Concepts: • **Software and its engineering** → **Formal software verification**; • **Computing methodologies** → **Distributed programming languages**.

Additional Key Words and Phrases: MRDTs, Eventual consistency, Automated verification, Replication-aware linearizability

ACM Reference Format:

Vimala Soundarapandian, Kartik Nagar, Aseem Rastogi, and KC Sivaramakrishnan. 2025. Automatically Verifying Replication-Aware Linearizability. *Proc. ACM Program. Lang.* 9, OOPSLA1, Article 111 (April 2025), 27 pages. <https://doi.org/10.1145/3720452>

1 Introduction

Modern decentralized applications often employ data replication across geographically distributed locations to enhance fault tolerance, minimize data access latency, and improve scalability. This practice is crucial for mitigating the impact of network failures and reducing data transmission delays to end users. However, these systems encounter the challenge of concurrent conflicting data updates across different replicas.

Recently, Mergeable Replicated Data Types (MRDTs) [11, 12, 23] have emerged as a systematic approach to the problem of ensuring that replicas remain eventually consistent despite concurrent conflicting updates. MRDTs draw inspiration from the Git version control system, where each

Authors' Contact Information: Vimala Soundarapandian, IIT Madras, Chennai, India, cs19d750@cse.iitm.ac.in; Kartik Nagar, IIT Madras, Chennai, India, nagark@cse.iitm.ac.in; Aseem Rastogi, Microsoft Research, Bangalore, India, aseemr@microsoft.com; KC Sivaramakrishnan, IIT Madras, Chennai, India and Tarides, Chennai, India, kcsrkr@cse.iitm.ac.in.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2025 Copyright held by the owner/author(s).

ACM 2475-1421/2025/4-ART111

<https://doi.org/10.1145/3720452>

update creates a new version, and any two versions can be merged explicitly through a user-defined merge function. `merge` is a ternary function that takes as input the two versions to be merged and their Lowest Common Ancestor (LCA), i.e., the most recent version from which the two versions diverged. As opposed to Conflict-Free Replicated Data Types (CRDTs)[22], which may have to carry around causal context metadata to ensure consistency, MRDTs can rely on the underlying system model to provide the causal context through the LCA. This results in implementations that are comparatively simpler and also more efficient. For example, if we consider state-based CRDTs, which are the closest analogue to the MRDT model, then any counter CRDT implementation would require $O(n)$ space, where n is the number of replicas (a lower bound proved by [4]), whereas a counter MRDT implementation only requires $O(1)$ space. The states maintained by CRDT implementations need to form a join semi-lattice, with all CRDT operations restricted to being monotonic functions and merge restricted to the lattice join. While these restrictions simplify the task of reasoning about correctness [5, 13, 18], crafting correct and efficient CRDT implementations itself becomes much harder.

MRDTs do not require any of the above restrictions, which helps in developing implementations with better space and time complexity. However, reasoning about correctness now becomes harder. Indeed, the MRDT system model allows arbitrary replicas to merge their states at arbitrary points in time, and this can result in subtle bugs requiring a very specific orchestration of merge actions. As part of this work, we discovered such subtle bugs in MRDT implementations claimed to be verified by previous works [23] (more details can be found in §5.2). The MRDT state as well as the implementation of data type operations and the merge function have to be cleverly designed to ensure strong eventual consistency. That is, despite concurrent conflicting updates and arbitrary ordering of merges, all replicas will eventually converge to the same state. Further, we would also like to show that an MRDT satisfies the functional behavior of the data type, along with the user-defined conflict resolution policy for concurrent conflicting updates (e.g., for a set data type, an *add-wins* policy that favors the add operation over a concurrent remove of the same element at different replicas). There have been a few works [11, 12, 23] that have looked at the problem of specifying and verifying MRDTs. However, they either restrict the system model by disallowing concurrent merges [12], focus only on convergence as the correctness specification [11, 12], or do not support automated verification [23].

In this work, we couch the correctness of MRDTs using the notion of *Replication-Aware Linearizability* (RA-linearizability) [28], which says that the state at any replica must be obtained by linearizing (i.e., constructing a sequence of) update operations that have been applied at the replica. As a first contribution, we adapt RA-linearizability to the MRDT system model (§3), and develop a simple specification framework for MRDTs based on conflict resolution policy for concurrent update operations. We show that an MRDT implementation can be linearized only under certain technical constraints on the conflict resolution policy and if the merge operation satisfies a weaker notion of commutativity called *conditional commutativity*. By ensuring that the linearization order obeys the conflict resolution policy for concurrent update operations and it remains the same across all replicas, we guarantee both strong eventual consistency and adherence to the user-provided specification.

Next, we propose a sound but not complete technique for proving RA-linearizability for MRDT implementations. The main challenge lies in showing that the merge function generates a state which is a linearization of its inputs. We develop a technique called *bottom-up linearization*, which relies on certain simple algebraic properties of the merge function to prove that it generates the correct linearization. We then design an induction scheme to *automatically* verify the required algebraic properties of merge for an arbitrary MRDT implementation. Our main insight here is to leverage the fact that the merge inputs are themselves linearizations, and hence, we can use

induction over their operation sequences. We extract a set of verification conditions (VCs) that are amenable to automated reasoning, and prove that if an MRDT implementation satisfies the VCs, it is linearizable (§4). While our development is focussed on MRDTs, our technique can be directly applied on state-based CRDTs. State-based CRDTs also have a merge-based system model which is slightly simpler than MRDTs as the merge function does not require any LCA.

Finally, we develop a framework in the F^{*} [27] programming language that allows implementing MRDTs and automatically mechanically proving the VCs required by our technique. The framework provides several advantages over previous works. First, we require the programmer to specify only the MRDT operations, the merge function, and the conflict resolution policy, in contrast to the earlier work that also requires proof constructs such as abstract simulation relations [23]. Second, the VCs are simple enough that in *all* the case studies we have done, including data types such as counter, set, map, boolean flag, and list, they are automatically discharged by F^{*}. Finally, we extract the verified implementations to OCaml using the F^{*} extraction pipeline and run them (§5). We have also implemented and verified a few state-based CRDTs using our framework. In the next section, we present the main ideas of our work informally through a series of examples.

2 Overview

2.1 System Model

The MRDT system model resembles a distributed version control system, such as Git [6], with replication centred around versioned states in branches and explicit merges. A replicated data store handles multiple objects independently [9, 19]; in our presentation, we focus on modeling a store with a single object. The state of the object is replicated across multiple replicas $r_1, r_2, \dots \in \mathcal{R}$ in the store. Clients interact with the store by performing query or update operations on one of the replicas, with update operations modifying its state. These replicas operate concurrently, allowing independent modifications without synchronization. They periodically (and non-deterministically) exchange updates with each other through a process called *merge*. Due to concurrent operations happening at multiple replicas, conflicts may arise, which must be resolved by the merge operation in an appropriate and consistent manner. An object has a type $\tau \in \text{Type}$, whose type signature $\langle O_\tau, Q_\tau, \text{Val}_\tau \rangle$ contains the set of supported update operations O_τ , query operations Q_τ and their return values Val_τ .

Definition 2.1. An MRDT implementation for a data type τ is a tuple $\mathcal{D}_\tau = \langle \Sigma, \sigma_0, \text{do}, \text{merge}, \text{query}, \text{rc} \rangle$, where:

- Σ is the set of states, $\sigma_0 \in \Sigma$ is the initial state.
- $\text{do} : \Sigma \times \mathcal{T} \times \mathcal{R} \times O_\tau \rightarrow \Sigma$ implements all update operations in O_τ , where \mathcal{T} is the set of timestamps.
- $\text{merge} : \Sigma \times \Sigma \times \Sigma \rightarrow \Sigma$ is a three-way merge function.
- $\text{query} : \Sigma \times Q_\tau \rightarrow \text{Val}_\tau$ implements all query operations in Q_τ , returning a value in Val_τ .
- $\text{rc} \subseteq O_\tau \times O_\tau$ is the conflict resolution policy to be followed for concurrent update operations.

An MRDT \mathcal{D}_τ provides implementations of *do*, *merge* and *query* which will be invoked by the data store appropriately. A client request to perform an update operation $o \in O_\tau$ at a replica r triggers the call $\text{do}(\sigma, t, r, o)$. This takes as input the current state $\sigma \in \Sigma$ of r , a unique timestamp $t \in \mathcal{T}$ and produces an updated state which is then installed at r . The data store ensures that timestamps are unique across all operations (which can be achieved through e.g. Lamport timestamps [14]).

Replicas can also receive states from other replicas, which are merged with the receiver's state using *merge*. The merge function is called with the current states of both the sender and receiver replicas and their lowest common ancestor (LCA), which represents the most recent common state

from which the two replicas diverged. Clients can query the state of the MRDT using the query method. This takes a MRDT state $\sigma \in \Sigma$ and a query operation as input and produces a return value. Note that a query operation cannot change the state at a replica.

While merging, it may happen that conflicting update operations may have been performed on the two states, in which case, the implementation also provides a conflict resolution policy rc . The merge function should make sure that this policy is followed while computing the merged state. To illustrate, we now present a couple of MRDT implementations: an increment-only counter and an observed-remove set.

The counter MRDT implementation is given in Fig. 1. The state space of the counter MRDT is simply the set of natural numbers, and it allows clients to perform only one update operation (inc) which increments the value of the counter. For merging two counter states σ_1 and σ_2 , whose lowest common ancestor is σ_\top , intuitively, we want to find the total number of increment operations across σ_1 and σ_2 . Since σ_\top already accounts for the effect of the common increments in σ_1 and σ_2 , we need to count the newer increments and then add them to σ_\top . This is achieved by adding $\sigma_1 - \sigma_\top$ and $\sigma_2 - \sigma_\top$ to σ_\top , which simplifies to the merge definition in Fig. 1. For example, suppose we have replicas r_1 and r_2 whose initial state was $\sigma_\top = 2$. Now, if there are 2 inc operations at r_1 and 3 inc operation at r_2 , their states will be $\sigma_1 = 4$ and $\sigma_2 = 5$. On merging r_2 at r_1 , $merge(\sigma_\top, \sigma_1, \sigma_2)$ will return 7, which reflects the total number of increments. The query method simply returns the current state of the counter. Finally, the increment operation commutes with itself, so there is no need to define a conflict resolution policy.

```

1:  $\Sigma = \mathbb{N}$ 
2:  $O = \{inc\}$ 
3:  $Q = \{rd\}$ 
4:  $\sigma_0 = 0$ 
5:  $do(\sigma, \_, \_, inc) = \sigma + 1$ 
6:  $merge(\sigma_\top, \sigma_1, \sigma_2) = \sigma_1 + \sigma_2 - \sigma_\top$ 
7:  $query(\sigma, rd) = \sigma$ 
8:  $rc = \emptyset$ 

```

Fig. 1. Counter MRDT implementation

```

1:  $\Sigma = \mathcal{P}(\mathbb{E} \times \mathcal{T})$ 
2:  $O = \{add_a, rem_a \mid a \in \mathbb{E}\}$ 
3:  $Q = \{rd\}$ 
4:  $\sigma_0 = \{\}$ 
5:  $do(\sigma, t, \_, add_a) = \sigma \cup \{(a, t)\}$ 
6:  $do(\sigma, \_, \_, rem_a) = \sigma \setminus \{(a, i) \mid (a, i) \in \sigma\}$ 
7:  $merge(\sigma_\top, \sigma_1, \sigma_2) =$ 
    $(\sigma_\top \cap \sigma_1 \cap \sigma_2) \cup (\sigma_1 \setminus \sigma_\top) \cup (\sigma_2 \setminus \sigma_\top)$ 
8:  $query(\sigma, rd) = \{a \mid (a, \_) \in \sigma\}$ 
9:  $rc = \{(rem_a, add_a) \mid a \in \mathbb{E}\}$ 

```

Fig. 2. OR-set MRDT implementation

rem_a method removes all entries in the set corresponding to the element a . An element a is considered to be present in the set if there is some (a, t) in the state.

The merge method takes as input the LCA set σ_\top and the two sets σ_1 and σ_2 to be merged, retains elements of σ_\top that were not removed in both sets, and includes the newly added elements from both sets. Since σ_\top is the most recent state from which the two sets diverged, the intersection of all three sets is the set of elements that were not removed from σ_\top in either branch, while the difference of either set with the σ_\top corresponds to the newly added elements. The query operation rd returns all the elements in the set. The conflict resolution relation rc orders rem_a before add_a of the same element in order to achieve the add-wins semantics. Note that all other pairs of operations ($add__$ and $add__$, $rem__$ and $rem__$, and add_x and rem_y with $x \neq y$) commute with each other, hence rc does not specify their ordering. We now consider whether the merge operation adheres to the conflict resolution policy.

An observed-remove set (OR-set) [22] is an implementation of a set data type that employs an add-wins conflict-resolution strategy, prioritizing addition in cases of concurrent addition and removal of the same element. Fig. 2 shows the OR-set MRDT implementation. This implementation is quite similar to the operation-based (op-based) CRDT implementation of OR-set [21]. The state of the OR-set is a set of element-timestamp pairs, with the initial state being an empty set. Clients can perform two operations for every element $a \in \mathbb{E}$: add_a and rem_a . The add_a method adds the element a along with the (unique) timestamp at which the operation was performed. The

2.2 RA-linearizability for MRDTs

We would like to verify that an MRDT implementation is correct, in the sense that in every execution, (a) replicas which have observed the same set of update operations converge to the same state, and (b) this state reflects the semantics of the implemented data type and the conflict resolution policy. Note that an update operation o is considered to be visible to a replica r either if o is directly applied by a client at r , or indirectly through merge with another replica r' on which o was visible. To specify MRDT correctness, we propose to use the notion of RA-linearizability [28]: the state at any replica during any execution must be achievable by applying a sequence (or linearization) of the update operations visible to the replica. Further, this linearization should obey the user-specified conflict resolution policy for concurrent operations, and the local replica order for non-concurrent operations.

Our definition of RA-linearizability allows viewing the state of an MRDT replica as a sequence of update operations applied on the initial state, thus abstracting over the merge function and how it handles concurrent operations. Consequently, any formal reasoning (e.g. assertion checking, functional correctness, equivalence checking etc.) can now essentially forget about the presence of merges, and only focus on update operations, with the additional guarantee that operations would have been correctly linearized, taking into account the conflict resolution policy and local replica ordering.

Proving RA-linearizability for MRDTs is straightforward when there is only a single replica on which all operations are performed, since there is no interleaving among operations on a single replica. Complexity arises when update operations happen concurrently across replicas, which are then merged. For a merge operation, we need to show that the output can be obtained by applying a linearization of update operations witnessed by both replicas being merged. However, the states being merged would have been obtained after an arbitrary number of update operations or even other merges. Further, the MRDT framework maintains only the states, but not the update operations leading to those states, thus requiring the verification technique to somehow infer the update operations leading to a state, and then show that merge constructs the correct linearization.

We break down this difficult problem gradually with a series of observations. We will start with an intuitively correct approach, show how it could be broken through examples, and gradually refine it to make it work. As a starting point, we first observe that we can leverage the following algebraic properties of the MRDT update operations and the merge function: (i) commutativity of merge and update operations, (ii) commutativity of merge, (iii) idempotence of merge, and (iv) commutativity of update operations. To motivate this, we first introduce some terminology. An event $e = \langle t, r, o \rangle$ is generated for every update operation instance, where t is the event's timestamp and r is the replica on which the update operation o is applied. Applying an event e on a replica with state σ changes the replica state to $e(\sigma) = \text{do}(\sigma, t, r, o)$ using the implementation of the operation o . Given a sequence of events $\pi = e_1 e_2 \dots e_n$, we use the notation $\pi(\sigma)$ to denote $e_n(\dots(e_2(e_1(\sigma))))$. Now, the properties described above can be formally defined as follows (forall $\sigma_\top, \sigma_1, \sigma_2, e, e'$):

- (P1) $\text{merge}(\sigma_\top, e(\sigma_1), \sigma_2) = e(\text{merge}(\sigma_\top, \sigma_1, \sigma_2))$
- (P2) $\text{merge}(\sigma_\top, \sigma_1, \sigma_2) = \text{merge}(\sigma_\top, \sigma_2, \sigma_1)$
- (P3) $\text{merge}(\sigma_\top, \sigma_\top, \sigma_\top) = \sigma_\top$
- (P4) $e(e'(\sigma)) = e'(e(\sigma))$

As per our proposed definition of RA-linearizability, we need to show that there exists a linearization of events visible at the replica such that the state of the replica can be obtained by applying this linearization. As mentioned earlier, an event can become visible at a replica either by a direct client application, or by merging with another replica. To illustrate this, consider the scenario

shown in Fig. 3 where two replicas with states σ_1 and σ_2 are being merged. These states were obtained by applying a sequence of events π_1 and π_2 respectively on the LCA state σ_\top . We call the events in π_1 and π_2 as local to their respective replicas. Now, when the two states are merged to create a new state σ_m we would need to show that the state $\sigma_m (= \text{merge}(\sigma_\top, \sigma_1, \sigma_2))$ can be obtained by linearizing all the events in π_1 and π_2 , and applying this linearization on the state σ_\top .

To show that the merge function constructs a linearization, we can take advantage of properties (P1)-(P4). In particular, commutativity of merge and update operation application (P1) allows us to move an event from the second argument of merge to outside, and we can then repeatedly apply this property to peel off all the events in π_1 . More formally, by performing induction on the sequence π_1 and using (P1), we can show that $\text{merge}(\sigma_\top, \pi_1(\sigma_\top), \sigma_2) = \pi_1(\text{merge}(\sigma_\top, \sigma_\top, \sigma_2))$. We can then use commutativity of merge (P2) to swap the last two arguments of merge, and then apply (P1) again to peel off all the events in π_2 , thus establishing that $\text{merge}(\sigma_\top, \sigma_\top, \pi_2(\sigma_\top)) = \pi_2(\text{merge}(\sigma_\top, \sigma_\top, \sigma_\top))$. Finally, using merge idempotence (P3), and combining all the previous results, we can infer that $\text{merge}(\sigma_\top, \sigma_1, \sigma_2) = \pi_2(\pi_1(\sigma_\top))$. Commutativity of update operations (P4) ensures that all linearizations of events in π_1 and π_2 lead to the same state, thus ratifying the specific linearization order $\pi_1\pi_2$ that we constructed using properties P1-P3. We call this process as bottom-up linearization, since we built the sequence from end through property (P1), linearizing one event at a time.

It is also easy to see that the counter MRDT implementation in Fig. 1 satisfies (P1)-(P4). In particular, commutativity of integer addition and subtraction essentially gives us (P1)-(P4) for free. While this strategy works for the counter MRDT, commutativity of all update operations is in general a very strong requirement, and would fail for other datatypes. For example, the OR-set MRDT of Fig. 2 does not satisfy (P4), as the add_a and rem_a operations do not commute.

In the presence of non-commutative update operations, the property (P1) now needs to be altered, as we need to consider the conflict resolution policy to decide the replica from which an event needs to be peeled off. To illustrate this, consider an OR-set execution depicted in Fig. 4. We show the version graph of the execution, where each oval represents a version. The state of the version is depicted inside the oval. The versions v_1 and v_2 are obtained by applying rem_a and add_a operations to the version v_\top on two different replicas (r_1 and r_2). Each edge is labeled with the event corresponding to the application of an operation. Let $\sigma_\top = \{\}$ denote the state of the LCA v_\top . The versions v_1 and v_2 are then merged at r_2 which gives rise to a new version v_m with state $\text{merge}(\sigma_\top, e_1(\sigma_\top), e_2(\sigma_\top))$. Now, since e_1 and e_2 do not commute, the conflict resolution policy of OR-set places e_1 (i.e. the remove operation) before e_2 (i.e. the add operation). Hence, we want the merged version to follow the linearization order $e_2(e_1(\sigma_\top))$. This requires us to first peel off the event e_2 from the third argument of merge. To achieve this, we can alter the property (P1) by making it aware of the conflict resolution policy as follows:

$$(P1') \quad (e_1, e_2) \in \text{rc} \implies \text{merge}(\sigma_\top, e_1(\sigma_1), e_2(\sigma_2)) = e_2(\text{merge}(\sigma_\top, e_1(\sigma_1), \sigma_2))^1$$

¹Note that we are abusing the rc notation slightly, since rc is a relation over operations O , but we are considering it over operation instances (i.e. events)

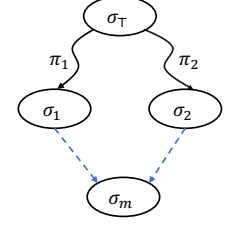


Fig. 3. Linearizing a merge operation

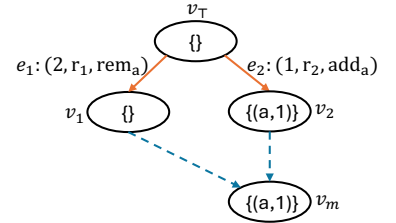


Fig. 4. OR-set execution

Property (P1') would then allow us to establish the required linearization order. Property (P4) also needs to be altered due to the presence of non-commutative update operations. We modify (P4) to enforce commutativity for non-rc related events, which gives us flexibility to include such events in any order while constructing the linearization sequence:

$$(P4') \quad (e_1, e_2) \notin rc \wedge (e_2, e_1) \notin rc \implies e_1(e_2(\sigma)) = e_2(e_1(\sigma))$$

However, we now face another major challenge: proving (P1') for the OR-set MRDT. For the counter MRDT, the operations and merge function used integer addition and subtraction, which commute with each other. But for the OR-set, add_a uses set union, while merge uses set difference and intersection, which do not commute in general. Hence, (P1') does not hold for arbitrary $\sigma_\top, \sigma_1, \sigma_2$.

To illustrate this concretely, consider the same execution of Fig. 4, except assume that the state σ_\top of the LCA v_\top is $\{(a, 1)\}$. Let us try to establish (P1') for the merge of versions v_1 and v_2 . First, note that as per the OR-set rc, the antecedent of (P1') is satisfied, as $(e_1, e_2) \in rc$. Now, the RHS in the consequent must contain the tuple $(a, 1)$, since the event e_2 adds $(a, 1)$ to the result of the merge. Does the LHS also contain $(a, 1)$? Expanding the definition of merge in the LHS, $(a, 1)$ will not be present in $(\sigma_\top \cap e_1(\sigma_\top) \cap e_2(\sigma_\top))$ (because $(a, 1) \notin e_1(\sigma_\top)$, as e_1 removes a). Similarly, since $(a, 1)$ is in σ_\top , it will not be present in $e_2(\sigma_\top) \setminus \sigma_\top$. It will not be in $e_1(\sigma_\top) \setminus \sigma_\top$, as e_1 removes a . To conclude, $(a, 1)$ will not be present in the LHS, thus invalidating the consequent of (P1').

However, we note that this particular execution is actually spurious, because the tuple $(a, 1)$ in the LCA could only have been added by another add_a operation whose timestamp is the same as e_2 . But this is not possible as the data store ensures that timestamps are unique across all events. In the general case, we would not be able to show (P1') for OR-set because the tuple (a, t) being added by the add_a operation (event e_2) could also be present in the LCA state. However, this situation cannot occur.

Thus, it is possible to show (P1') for all *feasible* states $\sigma_\top, \sigma_1, \sigma_2$ that may occur during an actual execution. In the case of OR-set, there are two arguments which are required to infer this: (i) timestamps are unique across all events and (ii) if a tuple (a, t) is present in the state σ , then there must have been an add_a operation with timestamp t in the history of events leading to σ . While the first argument is a property of the data store, the second argument is an invariant linking a state with the history of events leading to that state. Such arguments are in general hard to infer, and would also change across different MRDTs. We now present our second major observation which allows us to automatically verify (P1') for feasible states without requiring invariants like argument (ii) linking MRDT states and events.

2.3 Verification using Induction on Event Sequences

In order to show property (P1') for an MRDT implementation, we need to consider the feasible states which would be given as input to the merge function during an actual execution. We observe that we can leverage the RA-linearizability of the MRDT implementation, and hence characterize these feasible states by sequences of MRDT update operations (more precisely, events corresponding to update operation instances). We can now use induction over these sequences to establish property (P1'). Note that the input states to merge may themselves have been obtained through prior merges, but we can inductively assume that these prior merges resulted in correct linearizations. Since merge takes as input three states $(\sigma_\top, \sigma_1, \sigma_2)$, we need to consider three sequences that led to these states and induct on all the three separately.

Concretely, let π_\top be a sequence of events which when applied on the initial MRDT state σ_0 results in the state σ_\top . Since the LCA state always contains events which are common to the states σ_1 and σ_2 , π_\top will be the common prefix of the sequences leading to both σ_1 and σ_2 . We consider

the sequences π_1 and π_2 that consist of the local events which when applied on σ_τ led to σ_1 and σ_2 respectively. Fig. 5 depicts the situation. Notice that the last two events on each replica before the merge are fixed to be e_1 and e_2 , which would be related by the rc relation, as per the requirement of property (P1').

$$\text{merge}(\sigma_0, e_1(\sigma_0), e_2(\sigma_0)) = e_2(\text{merge}(\sigma_0, e_1(\sigma_0), \sigma_0)) \quad (1)$$

$$\begin{aligned} \text{merge}(\sigma_\tau, e_1(\sigma_\tau), e_2(\sigma_\tau)) &= e_2(\text{merge}(\sigma_\tau, e_1(\sigma_\tau), \sigma_\tau)) \\ \implies \text{merge}(e(\sigma_\tau), e_1(e(\sigma_\tau)), e_2(e(\sigma_\tau))) &= e_2(\text{merge}(e(\sigma_\tau), e_1(e(\sigma_\tau)), e(\sigma_\tau))) \end{aligned} \quad (2)$$

We first induct on the sequence π_τ which leads to the state σ_τ . For this, we assume that $\pi_1 = \pi_2 = \epsilon$, and hence $\sigma_\tau = \sigma_1 = \sigma_2 = \pi_\tau(\sigma_0)$. We also assume the antecedent of property (P1'), i.e. $(e_1, e_2) \in \text{rc}$, and hence our goal is to show its consequent. For the OR-set, e_1 will be a rem_a event, while e_2 will be an add_a event (say with timestamp t).

Eqn. (1) is the base-case of the induction (where $\pi_\tau = \epsilon$), and this can be now directly discharged since σ_0 is an empty set, and hence clearly won't contain (a, t) . Eqn. (2) is the inductive case, which assumes that (P1') is true for some LCA state σ_τ , and tries to prove the property when one more update operation (signified by the event e) is applied on the LCA (and also on both σ_1 and σ_2 , since LCA

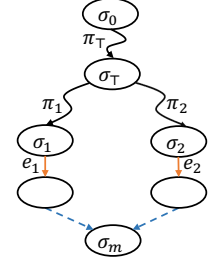


Fig. 5. Induction on event sequences

operations are common to both states to be merged). This can also be automatically discharged with the property that events e, e_1, e_2 have different timestamps. Intuitively, the inductive hypothesis establishes that $(a, t) \notin \sigma_\tau$, and since the timestamp of event e is different from e_1 and e_2 , it cannot add (a, t) to the LCA, thus preserving the property that $(a, t) \notin e(\sigma_\tau)$, thereby implying the consequent. This completes the proof for property (P1') for any arbitrary LCA state σ_τ that may be feasible in an actual execution. A similar inductive strategy is used for proving property (P1') for feasible states σ_1 and σ_2 (more details in §4).

2.4 Intermediate Merges

In our linearization strategy for merges (given by properties (P1'-P4')), we first considered the local update operations of each branch, linearized them according to the conflict-resolution policy, and then applied this sequence on the LCA. This effectively orders the update operations that led to the LCA before the update operations local to each branch.

However, in a Git-based execution model, due to a phenomenon known as intermediate merges, it may happen that update operations of the LCA may need to be linearized after update operations local to a branch. To illustrate this, consider an execution of the OR-set MRDT as shown in Fig. 6. There are 3 operations and 2 merges being performed in this execution, with the events e_1, e_3 at replica r_1 and event e_2 at replica r_2 .

Instead of merging with the latest version v_3 at replica r_1 , replica r_2 first merges with an intermediate version v_1 to generate the version v_4 . Next, this version v_4 is merged with the latest version v_3 of replica r_1 . However, note that for this merge, the LCA will be version v_1 . This is because the set of events associated with version v_3 is $\{e_1, e_3\}$, while for version v_4 , it is $\{e_1, e_2\}$. Hence, the set of common

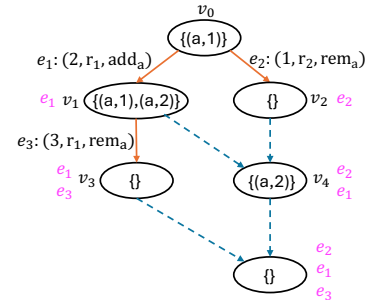


Fig. 6. Intermediate merge

events among both versions would be $\{e_1\}$, which corresponds to the version v_1 . Indeed, in the version graph, both v_1 and v_0 are ancestors of v_3 and v_4 , but v_1 is the lowest common ancestor².

In Fig. 6, we have also provided the linearization of events associated with each version. Notice that for version v_4 , which is obtained through a merge of v_1 and v_2 , the conflict resolution policy of the OR-set linearizes e_2 before e_1 . Now, for the merge of v_3 and v_4 , we have a situation where a local event (e_2 in v_4) needs to be linearized before an event of the LCA (e_1 in v_1). This does not fit our linearization strategy. Let us see why. If we were to try to apply (P1'), it would linearize e_1 after e_3 , since these are the last operations in the two states to be merged and the conflict resolution policy orders $\text{add}_a(e_1)$ after $\text{rem}_a(e_3)$. However, in the execution, e_1 and e_3 are causally related, i.e. e_1 occurs before e_3 on the same replica, and hence they should be linearized in that order. Intuitively, property (P1') does not work because it does not consider the possibility that the last event in one replica could be visible to the last event in another replica, and hence the linearization must obey the visibility relation.

In order to handle this situation, we consider another algebraic property (P1-1), which explicitly forces visibility relation among the last events by making one of them part of the LCA:

$$(P1-1) \text{ merge}(e_1(\sigma_0), e_3(\sigma_1), e_1(\sigma_2)) = e_3(\text{merge}(e_1(\sigma_0), \sigma_1, e_1(\sigma_2)))$$

Note that events in the LCA are visible to events on both replicas being merged. Hence, by having the same event e_1 in both the first and third argument to merge in the LHS, e_3 would have to be linearized after e_1 to respect the visibility order, thus over-riding the rc ordering among them. Property (P1-1) can be directly applied to the execution in Fig. 6 for the merge of v_3 and v_4 (with σ_0 as the state of v_0 , σ_1 as the state of v_1 and σ_2 as the state of v_2), constructing the correct linearization.

We will revisit the example in Fig. 6 and properties (P1') and (P1-1) in a more formal setting in §4, renaming them as BOTTOMUP-2-OP and BOTTOMUP-1-OP. We will also identify the conditions under which these properties can guarantee the existence of a correct linearization.

3 Problem Definition

In this section, we formally define the semantics of the replicated data store on top of which the MRDT implementations operate (§3.1), the notion of RA-linearizability for MRDTs (§3.2), and the process of bottom-up linearization (§3.3).

3.1 Semantics of the Replicated Data Store

The semantics of the replicated store defines all possible executions of an MRDT implementation. Formally, the semantics are parameterized by an MRDT implementation $\mathcal{D} = \langle \Sigma, \sigma_0, \text{do}, \text{merge}, \text{query}, \text{rc} \rangle$ of type $\tau = \langle \mathcal{O}_\tau, \mathcal{Q}_\tau, \text{Val}_\tau \rangle$ and are represented by a labeled transition system $\mathcal{S}_{\mathcal{D}} = (\Phi, \rightarrow)$. Each configuration in Φ maintains a set of versions, where each version is created either by applying an MRDT operation to an existing version, or by merging two versions. Each replica is associated with a head version, which is the most recent version seen at the replica. Formally, each configuration C in Φ is a tuple $\langle N, H, L, G, \text{vis} \rangle$, where:

- $N : \text{Version} \rightarrow \Sigma$ is a partial function that maps versions to their states (Version is the set of all possible versions).
- $H : \mathcal{R} \rightarrow \text{Version}$ is also a partial function that maps replicas to their head versions. A replica is considered active if it is in the domain of H of the configuration.
- $L : \text{Version} \rightarrow \mathbb{P}(\mathcal{E})$ maps a version to the set of events that led to this version. Each event $e \in \mathcal{E}$ is an update operation instance, uniquely identified by a timestamp value (we define $\mathcal{E} = \mathcal{T} \times \mathcal{R} \times \mathcal{O}$).

²in §3, we will formally prove that the LCA of two versions according to the version graph contains the intersection of events in both versions.

[CREATEBRANCH]

$$\frac{N' = N[v \mapsto N(H(r))] \quad \begin{array}{c} r \in \text{dom}(H) \quad r' \notin \text{dom}(H) \quad v \notin \text{dom}(N) \\ H' = H[r \mapsto v] \quad L' = L[v \mapsto L(H(r))] \quad G' = (\text{dom}(N) \cup \{v\}, E \cup \{(H(r), v)\}) \end{array}}{(N, H, L, G, \text{vis}) \xrightarrow{\text{createBranch}(r', r)} (N', H', L', G', \text{vis})}$$

[APPLY]

$$\frac{\begin{array}{c} o \in O_\tau \quad \forall e' \in \bigcup \text{range}(L). \text{time}(e') \neq t \quad r \in \text{dom}(H) \quad v \notin \text{dom}(N) \quad N' = N[v \mapsto \text{do}(N(H(r)), e)] \\ H' = H[r \mapsto v] \quad L' = L[v \mapsto L(H(r)) \cup \{e\}] \quad G' = (\text{dom}(N'), E \cup \{(H(r), v)\}) \quad \text{vis}' = \text{vis} \cup (L(H(r)) \times \{e\}) \end{array}}{(N, H, L, G, \text{vis}) \xrightarrow{\text{apply}(t, r, o)} (N', H', L', G', \text{vis}')} e = (t, r, o)$$

[MERGE]

$$\frac{\begin{array}{c} r_1, r_2 \in \text{dom}(H) \quad v \notin \text{dom}(N) \quad v_\top = \text{LCA}(H(r_1), H(r_2)) \quad N' = N[v \mapsto \text{merge}(N(v_\top), N(H(r_1)), N(H(r_2)))] \\ H' = H[r_1 \mapsto v] \quad L' = L[v \mapsto L(H(r_1)) \cup L(H(r_2))] \quad G' = (\text{dom}(N'), E \cup \{(H(r_1), v), (H(r_2), v)\}) \end{array}}{(N, H, L, G, \text{vis}) \xrightarrow{\text{merge}(r_1, r_2)} (N', H', L', G', \text{vis})}$$

[QUERY]

$$\frac{r \in \text{dom}(H) \quad q \in Q_\tau \quad a = \text{query}(N(H(r)), q)}{(N, H, L, G, \text{vis}) \xrightarrow{\text{query}(r, q, a)} (N, H, L, G, \text{vis})}$$

Fig. 7. Semantics of the replicated datastore

- $G = (\text{dom}(N), E)$ is the version graph, whose vertices represent the versions in the configuration (i.e. those in the domain of N) and whose edges represent a relationship between different versions (we explain the different types of edges below).
- $\text{vis} \subseteq \mathcal{E} \times \mathcal{E}$ is a partial order over events.

Figure 7 gives a formal description of the transition rules. CREATEBRANCH forks a new replica r' from an existing replica r , installing a new version v at r' with the same state as the head version $H(r)$ of r , and adding an edge $(H(r), v')$ in the version graph. APPLY applies an update operation o on some replica r , generating a new event e with a timestamp different than all events generated so far. $\bigcup \text{range}(L)$ denotes the set of events witnessed across all versions. A new version v is also created whose state is obtained by applying o on the current state of the replica r . The version graph is updated by adding the edge $(H(r), v)$. The vis relation as well as the function L , which tracks events applied at each version, are also updated. In particular, each event e' already applied at r , i.e. $e' \in L(H(r))$, is made visible to e : $(e', e) \in \text{vis}$, while $L'(v)$ is obtained by adding e to $L(H(r))$.

MERGE takes two replicas r_1 and r_2 , applies the merge function on the states of their head versions to generate a new version v , which is installed as the new head version at r_1 . Edges are added in the version graph from the previous head versions of r_1 and r_2 to v . $L(v)$ is obtained by taking a union of $L(r_1)$ and $L(r_2)$, and there is no change in the visibility relation. QUERY takes a replica r and a query operation q and applies q to the state at the head version of r , returning an output value a . Note that the QUERY transition does not modify the configuration and the return value of the query is stored as part of the transition label. While our operational semantics is based on and inspired by previous works [11, 23], we note that it is more general and precisely captures the MRDT system model as opposed to previous works. In particular, Kaki et al. [11] place significant restrictions on the MERGE transition, disallowing arbitrary replicas to be merged to ensure that there is a total order on the merge transitions. While the semantics in the work by Soundarapandian et al. [23] does allow arbitrary merges, it is more abstract and high-level, and does not even keep track of versions and the version graph.

Notation: We now introduce some notation that will be used throughout the paper. Given a configuration C , we use $X(C)$ to project the component X of C . For a relation R , we use $x \xrightarrow{R} y$ to signify that $(x, y) \in R$. We use $R|_S$ to indicate the relation as given by R but restricted to elements of the set S . Let R^* denote the reflexive-transitive closure of R , and let R^+ denote the transitive closure of R . For an event e , we use the projection functions op , time , rep to obtain the update operation, timestamp and replica resp. For a sequence of events π , $\pi|_S(\sigma)$ denotes application of the sub-sequence of π restricted to events in S . For a configuration C , we use $e_1 \parallel_C e_2$ to denote that e_1 and e_2 are concurrent, that is $\neg(e_1 \xrightarrow{\text{vis}(C)} e_2 \vee e_2 \xrightarrow{\text{vis}(C)} e_1)$. Given a total order over a set of events \mathcal{E} , represented by a sequence π , and $\text{lo} \subseteq \mathcal{E} \times \mathcal{E}$, we say that π extends lo if $\text{lo} \subseteq \pi$. The relation rc orders update operations, but for convenience we sometime use it for ordering events, with the intention that it is actually being applied to the underlying update operations. We use $e_1 \neq e_2$ to indicate that $\text{time}(e_1) \neq \text{time}(e_2)$.

We define the initial configuration of $\mathcal{S}_{\mathcal{D}}$ as $C_0 = \langle N_0, H_0, L_0, G_0, \emptyset \rangle$, which consists of only one replica r_0 . Here, $H_0 = [r_0 \mapsto v_0]$, $N_0 = [v_0 \mapsto \sigma_0]$, where σ_0 is the initial state as given by \mathcal{D}_τ , while v_0 denotes the initial version and $L_0 = [v_0 \mapsto \emptyset]$. The graph $G_0 = (\{v_0\}, \emptyset)$ is the initial version graph. An execution of $\mathcal{S}_{\mathcal{D}}$ is defined to be a finite sequence of transitions, $C_0 \xrightarrow{t_1} C_1 \xrightarrow{t_2} C_2 \dots \xrightarrow{t_n} C_n$. Note that the label of a transition corresponds to its type. Let $\llbracket \mathcal{S}_{\mathcal{D}} \rrbracket$ denote the set of all possible executions of $\mathcal{S}_{\mathcal{D}}$.

Finally, as mentioned earlier, merge is a ternary function, taking as input the states of two versions to be merged, and the state of the lowest common ancestor (LCA) of the two versions. Version $v_1 \in V$ is defined to be a causal ancestor of version $v_2 \in V$ if and only if $(v_1, v_2) \in E^*$.

Definition 3.1 (LCA). Given a version graph $G = (V, E)$ and versions $v_1, v_2 \in V$, $v_\top \in V$ is defined to be the lowest common ancestor of v_1 and v_2 (denoted by $\text{LCA}(v_1, v_2)$) if (i) $(v_\top, v_1) \in E^*$ and $(v_\top, v_2) \in E^*$, (ii) $\forall v \in V. (v, v_1) \in E^* \wedge (v, v_2) \in E^* \implies (v, v_\top) \in E^*$.

Note that the version history graph at any point in any execution is guaranteed to be acyclic (i.e. a DAG), and hence the LCA (if it exists) is guaranteed to be unique. We now present an important property linking the LCA of two versions with events applied at each version.

LEMMA 3.2. *Given a configuration $C = \langle N, H, L, G, \text{vis} \rangle$ reachable in some execution $\tau \in \llbracket \mathcal{S}_{\mathcal{D}} \rrbracket$ and two versions $v_1, v_2 \in \text{dom}(N)$, if v_\top is the LCA of v_1 and v_2 in G , then $L(v_\top) = L(v_1) \cap L(v_2)$ ³.*

Thus, the events of the LCA are exactly those applied at both the versions. This intuitively corresponds to the fact that $\text{LCA}(v_1, v_2)$ is the most recent version from which the two versions v_1 and v_2 diverged. Note that it is possible that the LCA may not exist for two versions. Fig. 8 depicts the version graph of such an execution. Vertices with in-degree 1 (i.e. v_1, v_2, v_3, v_4) have been generated by applying a new update operation (with the orange edges labeled by the corresponding events e_1, e_2, e_3, e_4), while vertices with in-degree 2 have been obtained by merging two other versions (depicted by blue edges). The merge of v_1 and v_4 (leading to v_6) has a unique LCA v_0 , similarly, merge of v_2 and v_3 (leading to v_5) also has a unique LCA v_0 . However, if we now want to merge v_5 and v_6 , both v_1 and v_2 are ancestors, but there is no LCA. We note that this execution will actually be prohibited by the semantics of Kaki et al. [11], since the two merges leading to v_5 and v_6 are concurrent.

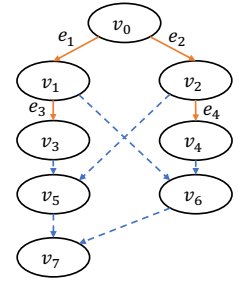


Fig. 8. Version Graph with no LCA for v_5 and v_6

³All proofs can be found in Appendix §A of the extended version [25] of the paper.

Notice that $L(v_5) = \{e_1, e_2, e_3\}$, while $L(v_6) = \{e_1, e_2, e_4\}$. Hence, by Lemma 3.2, $L(LCA(v_5, v_6)) = \{e_1, e_2\}$, but such a version is not generated during the execution. To resolve this issue, we introduce the notion of *potential* LCAs.

Definition 3.3 (Potential LCAs). Given a version graph $G = (V, E)$ and versions $v_1, v_2 \in V$, $v_\top \in V$ is defined to be a potential LCA of v_1 and v_2 if (i) $(v_\top, v_1) \in E^*$ and $(v_\top, v_2) \in E^*$, (ii) $\neg(\exists v. (v, v_1) \in E^* \wedge (v, v_2) \in E^* \wedge (v_\top, v) \in E^*)$.

For merging v_1 and v_2 , we first find all the potential LCAs, and recursively merge them to obtain the actual LCA state. For the execution in Fig. 8, the potential LCAs of v_5 and v_6 would be v_1 and v_2 (with $L(v_1) = \{e_1\}$ and $L(v_2) = \{e_2\}$); merging them would get us the actual LCA. In §A.1 of the extended version [25] of the paper, we prove that this recursive merge-based strategy is guaranteed to generate the actual LCA.

3.2 Replication-aware Linearizability for MRDTs

As mentioned in §2, our goal is to show that the state of every version v generated during an execution is a linearization of the events in $L(v)$. We use the notation lo to indicate the linearization relation, which is a binary relation over events. For an execution in \mathcal{SD} , we want lo between the events of the execution to satisfy certain desirable properties: (i) lo between two events should not change during an execution, (ii) lo should obey the conflict resolution policy for concurrent events and (iii) lo should obey the replica-local vis ordering for non-concurrent events. This would ensure that two versions which have observed the same set of events will have the same state (i.e. *strong eventual consistency*), and this state would also be a linearization of update operations of the data type satisfying the conflict resolution policy.

While the lo relation in classical linearizability literature is typically a total order, in our context, we take advantage of commutativity of update operations, and only define lo over non-commutative events. As we will see later, this flexibility allows us to have different sequences of events which extend the same lo relation between non-commutative events, and hence are guaranteed to lead to the same state. We use the notation $e \rightleftarrows e'$ to indicate that events e and e' commute with each other. Formally, this means that $\forall \sigma. e(e'(\sigma)) = e'(e(\sigma))$. Two update operations o, o' commute if $\forall e, e'. \text{op}(e) = o \wedge \text{op}(e') = o' \implies e \rightleftarrows e'$. As mentioned earlier, the rc relation is also only defined between non-commutative update operations.

LEMMA 3.4. *Given a set of events \mathcal{E} , if $\text{lo} \subseteq \mathcal{E} \times \mathcal{E}$ is defined over every pair of non-commutative events in \mathcal{E} , then for any two sequences π_1, π_2 which extend lo , for any state σ , $\pi_1(\sigma) = \pi_2(\sigma)$.*

Given a configuration $C = \langle N, H, L, G, \text{vis} \rangle$, let $\mathcal{E}_C = \bigcup \text{range}(L(C))$ denote the set of events witnessed across all versions in C . Then, our goal is to define an appropriate linearization relation $\text{lo}_C \subseteq \mathcal{E}_C \times \mathcal{E}_C$, which adheres to the rc relation for concurrent events, the vis relation for non-concurrent events, and for every version $v \in \text{dom}(N)$, $N(v)$ should be obtained by sequentializing the events in $L(v)$, with the sequence extending lo_C . Note that this requires lo^+ to be irreflexive⁴.

We now demonstrate that an lo relation with all the desirable properties may not exist for all executions. Suppose there are MRDT update operations o, o' such that $o \xrightarrow{\text{rc}} o'$. Fig. 9 contains a part of the version graph generated during some execution, containing two instances of both o and o' . We use $e_i : o_i$ to denote that event $\text{op}(e_i) = o_i$. Notice that e_1 and e_4 , e_2 and e_3 are concurrent, while e_1 and e_3 , e_2 and e_4 are non-concurrent. Applying the rc ordering on concurrent events, we would want $e_3 \xrightarrow{\text{lo}} e_2$ and $e_4 \xrightarrow{\text{lo}} e_1$, while applying vis ordering, we would want $e_1 \xrightarrow{\text{lo}} e_3$ and $e_2 \xrightarrow{\text{lo}} e_4$.

⁴ lo need not be transitive, as we only want to define lo between non-commutative events, and non-commutativity is not a transitive property

However, this results in a lo-cycle, thus making it impossible to construct a sequence of update operations for the merge of v_5 and v_6 , which adheres to the lo ordering.

Notice that the above execution only requires the rc relation to be non-empty (i.e. there should exist some $(o, o') \in rc$). If the rc relation is empty, then all update operations would commute with each other, and hence the lo relation would also be empty. If rc is non-empty, rc^+ should be irreflexive to ensure irreflexivity of lo^+ . Note that rc^+ being irreflexive means that for any MRDT update operation o , $(o, o) \notin rc$, and hence o must commute with itself, since rc relation is defined for all pairs of non-commutative update operations. Furthermore, Fig. 9 shows that even if rc^+ is irreflexive, it may still not be possible to construct an lo relation which can be extended to a total order and which adheres to the rc relation between all pairs of concurrent events. To ensure existence of an lo relation such that lo^+ is irreflexive when rc^+ is irreflexive, we define it as follows:

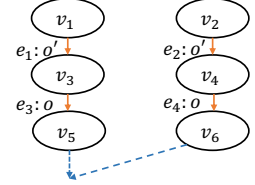


Fig. 9. Example demonstrating cycle in lo

Definition 3.5 (Linearization relation). Let C be a configuration reachable in some execution in $\llbracket \mathcal{SD} \rrbracket$. Let \mathcal{E}_C be the set of events in C . Then, lo_C is defined as:

$$\begin{aligned} \forall e_1, e_2 \in \mathcal{E}_C. e_1 \xrightarrow{lo_C} e_2 &\Leftrightarrow (e_1 \xrightarrow{vis(C)} e_2 \wedge \neg e_1 \rightharpoonup e_2) \\ &\vee (e_1 \parallel_C e_2 \wedge e_1 \xrightarrow{rc} e_2 \wedge \neg(\exists e_3 \in \mathcal{E}. e_2 \xrightarrow{vis(C)} e_3 \wedge \neg e_2 \rightharpoonup e_3)) \end{aligned}$$

lo_C follows the visibility relation only between non-commutative events. For concurrent non-commutative events e_1 and e_2 with $e_1 \xrightarrow{rc} e_2$, lo_C follows the rc relation only if there is no event e_3 such that e_2 is visible to e_3 and e_2 does not commute with e_3 . Applying this definition to the execution in Fig. 9, for the configuration obtained after merge, we would have neither $e_4 \xrightarrow{lo} e_1$, nor $e_3 \xrightarrow{lo} e_2$, thus avoiding the cycle in lo.

LEMMA 3.6. *For an MRDT \mathcal{D} such that rc^+ is irreflexive, for any configuration C reachable in \mathcal{SD} , lo_C^+ is irreflexive.*

Going forward, we will assume that rc^+ is irreflexive for any MRDT \mathcal{D} . We note that restricting lo to not always obey the rc relation by considering non-commutative update operations happening locally (and thus related by vis) is also sensible from a practical perspective. For example, in the case of OR-set, even though we have $rem_a \xrightarrow{rc} add_a$, if add_a is locally followed by another rem_a , it does not make sense to order a concurrent rem_a event before the add_a event. More generally, if an event e_2 is visible to another event e_3 with which it does not commute, then e_2 is effectively "overwritten" by e_3 , and hence there is no need to linearize a concurrent event e_1 before e_2 .

While lo_C is now guaranteed to be irreflexive for any configuration C , and hence can be extended to a sequence, it now no longer enforces an ordering among all non-commutative pairs of events. Thus, there could exist sequences π_1, π_2 extending an lo_C relation which may contain a pair of non-commutative events in different orders. For example, in Fig. 9, for the configuration C obtained after the merge, $lo_C = \{(e_1, e_3), (e_2, e_4)\}$, resulting in sequences $\pi_1 = e_1 e_2 e_3 e_4$ and $\pi_2 = e_1 e_3 e_2 e_4$ which both extend lo_C , but contain the non-commutative events e_2 and e_3 in different orders. Thus, Lemma 3.4 can no longer be applied, and it is not guaranteed that π_1 and π_2 would lead to the same state. Notice that in the sequences π_1 and π_2 above, even though e_2 and e_3 appear in different orders, e_4 always appears after both. Indeed, e_4 must appear after e_2 due to visibility relation, and since e_3 and e_4 commute with each other (since both correspond to the same operation o), it is

enough to consider sequences where e_4 appears after e_3 . Based on the above observation, we now introduce a notion called conditional commutativity to ensure that sequences such as π_1, π_2 would lead to the same state:

Definition 3.7 (Conditional Commutativity). Events e and e' are said to conditionally commute with respect to event e'' (denoted by $e \stackrel{e''}{\rightleftarrows} e'$) if $\forall \sigma \in \Sigma. \forall \pi \in \mathcal{E}^*. e''(\pi(e(e'(\sigma)))) = e''(\pi(e'(e(\sigma))))$.

Update operations o and o' conditionally commute w.r.t. update operation o'' if $\forall e, e', e''. op(e) = o \wedge op(e') = o' \wedge op(e'') = o'' \Rightarrow e \stackrel{e''}{\rightleftarrows} e'$. For example, for the OR-set MRDT of Fig. 2, $add_a \stackrel{rem_a}{\rightleftarrows} rem_a$. Even though add and $remove$ operations of the same element do not commute with each other, if there is guaranteed to be a future $remove$ operation, then they do commute. For the execution in Fig. 9, if e_2 and e_3 conditionally commute w.r.t. e_4 , then both the sequences π_1 and π_2 will lead to the same state. For non-commutative update operations that are not ordered by lo , we enforce their conditional commutativity through the following property:

$$COND-COMM(\mathcal{D}) \triangleq \forall o_1, o_2, o_3 \in O. (o_1 \xrightarrow{rc} o_2 \wedge \neg o_2 \rightleftarrows o_3) \Rightarrow o_1 \stackrel{o_3}{\rightleftarrows} o_2$$

$COND-COMM(\mathcal{D})$ is a property of an MRDT \mathcal{D} , enforcing conditional commutativity of update operations o_1 and o_2 w.r.t. o_3 if o_2 does not commute with o_3 . Connecting this with the definition of linearization relation, if there are events e_1, e_2, e_3 performing operations o_1, o_2, o_3 resp., and if $e_1 \xrightarrow{rc} e_2, e_2 \xrightarrow{vis} e_3$ and $\neg e_2 \rightleftarrows e_3$, then there will not be a linearization relation between e_1 and e_2 . However, $COND-COMM(\mathcal{D})$ would then ensure that the ordering of e_1 and e_2 will not matter, due to the presence of the event e_3 . We also formalize the requirement of an rc relation between all pairs of non-commutative update operations:

$$RC-NON-COMM(\mathcal{D}) \triangleq \forall o_1, o_2 \in O. \neg o_1 \rightleftarrows o_2 \Leftrightarrow o_1 \xrightarrow{rc} o_2 \vee o_2 \xrightarrow{rc} o_1$$

LEMMA 3.8. *For an MRDT \mathcal{D} which satisfies $RC-NON-COMM(\mathcal{D})$ and $COND-COMM(\mathcal{D})$, for any reachable configuration C in $\mathcal{S}_{\mathcal{D}}$, for any two sequences π_1, π_2 over \mathcal{E}_C which extend lo_C , for any state σ , $\pi_1(\sigma) = \pi_2(\sigma)$.*

Definition 3.9 (RA-linearizability of MRDT). Let \mathcal{D} be an MRDT which satisfies $RC-NON-COMM(\mathcal{D})$ and $COND-COMM(\mathcal{D})$. Then, a configuration $C = \langle N, H, L, G, vis \rangle$ of $\mathcal{S}_{\mathcal{D}}$ is RA-linearizable if, for every active replica $r \in range(H)$, there exists a sequence π consisting of all events in $L(H(r))$ such that $lo(C)_{|L(H(r))} \subseteq \pi$ and $N(H(r)) = \pi(\sigma_0)$. An execution $\tau \in \llbracket \mathcal{S}_{\mathcal{D}} \rrbracket$ is RA-linearizable if all of its configurations are RA-linearizable. Finally, \mathcal{D} is RA-linearizable if all of its executions are RA-linearizable.

For a configuration to be RA-linearizable, every active replica must have a state which can be obtained by applying a sequence of events witnessed at that replica, and that sequence must obey the linearization relation of the configuration. For an execution to be RA-linearizable, all of its configurations must be RA-linearizable. Lemma 3.6 ensures the existence of a sequence extending the linearization relation, while Lemma 3.8 ensures that two versions which have witnessed the same set of events will have the same state (i.e. strong eventual consistency). Further, we also show that if an MRDT is RA-linearizable, then for any query operation in any execution, the query result is derived from the state obtained by applying the update events seen at the corresponding replica right before the query:

LEMMA 3.10. *If MRDT \mathcal{D} is RA-linearizable, then for all executions $\tau \in \llbracket \mathcal{S}_{\mathcal{D}} \rrbracket$, for all transitions $C \xrightarrow{\text{query}(r,q,a)} C'$ in τ where $C = \langle N, H, L, G, \text{vis} \rangle$, there exists a sequence π consisting of all events in $L(H(r))$ such that $\text{lo}(C)_{|L(H(r))} \subseteq \pi$ and $a = \text{query}(\pi(\sigma_0), q)$.*

Compared to the definition of RA-linearizability in the work by Wang et. al. [28], there is one major difference: Wang et. al. also consider a sequential specification in the form of a set of valid sequences of data-type operations, and requires the linearization sequence to belong to the specification. Our definition simply requires the state of a replica to be a linearization of the update operations applied to the replica, without appealing to a separate sequential specification. Once this is done, we can separately show that a linearization of the MRDT operations obeys the sequential specification. For this, we can ignore the presence of the merge operation as well as the MRDT system model (which are taken care of by the RA-linearizability definition), thus boiling down to proving a specification over a sequential functional implementation, which is a well-studied problem.

3.3 Bottom-up Linearization

As demonstrated in §2, our approach to show RA-linearizability of an MRDT implementation is based on using algebraic properties of merge (specifically, commutativity of merge and update operation application) which allows us to show that the result of a merge operation is a linearization of the events in each of the versions being merged. We first describe a generic template for the algebraic properties which can be used to prove RA-linearizability:

$$\frac{\forall j. \pi_j \in \mathcal{E} \cup \{\epsilon\} \quad l, a, b \in \Sigma \quad \pi \in \{\pi_0, \pi_1, \pi_2\} \quad \forall j. \pi'_j = \pi_j - \pi}{\text{merge}(\pi_0(l), \pi_1(a), \pi_2(b)) = \pi(\text{merge}(\pi'_0(l), \pi'_1(a), \pi'_2(b)))} \quad [\text{BOTTOMUPTEMPLATE}]$$

The template for the algebraic property is given in the conclusion of the above rule, while the premises describe certain conditions. Each π_j for $j \in \{0, 1, 2\}$ is a sequence of 0 or 1 event (i.e. either ϵ or a single event e_j), while l, a, b are arbitrary states of the MRDT. Note that applying the ϵ event on a state leaves it unchanged (i.e. $\epsilon(s) = s$). Then, we can select one event π which has been applied to the arguments of merge on the LHS, and bring it outside, i.e. remove the event from each argument on which it was applied, and instead apply the event to the result of merge. Note that the notation $\pi'_j = \pi_j - \pi$ means that if $\pi = \pi_j$, then $\pi'_j = \epsilon$, else $\pi'_j = \pi_j - \pi$.

The rule (P1') given in §2.2 can be seen as an instantiation of the above template with $\pi_0 = \epsilon, \pi_1 = e_1, \pi_2 = e_2$ and $\pi = e_2$ where $e_1 \xrightarrow{\text{rc}} e_2$. Similarly, (P1-1) is another instantiation with $\pi_0 = \pi_2 = e_1, \pi_1 = e_3$ and $\pi = e_3$ where $e_3 \neq e_1$. Assuming that the input arguments to merge are obtained through sequences of events τ_0, τ_1, τ_2 , the template rule builds the linearization sequence $\tau = \tau' e$ where e is the last event in one of the τ_i s, and τ' is recursively generated by applying the rule on $\tau' = \tau - e$. We call this procedure as *bottom-up linearization*. The event e should be chosen in such a way that the sequence τ is an extension of the linearization relation (Def. 3.5).

However, bottom-up linearization might fail if the last event in the merge output is not the last event in any of the three arguments to merge. For example, consider the execution shown in Fig. 10, where there exists an rc-chain: $o_2 \xrightarrow{\text{rc}} o_3 \xrightarrow{\text{rc}} o_1$, and o_1 and o_2 are non-commutative. e_1 is visible to e_2 , while event e_3 is concurrent to e_1 and e_2 . Now, for the version obtained after merging v_3 and v_4 , the linearization relation would be $e_1 \xrightarrow{\text{lo}} e_2$ and $e_2 \xrightarrow{\text{lo}} e_3$.

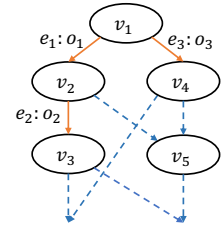


Fig. 10. Example demonstrating the failure of bottom-up linearization in the presence of an rc-chain

Notably, even though e_1 and e_3 are also concurrent, and rc orders o_3 before o_1 , this will not result in a linearization relation from e_3 to e_1 , due to the presence of a non-commutative update operation e_2 to which e_1 is visible. The bottom-up linearization for the merge of v_3 and v_4 , will result in the sequence $e_1e_2e_3$, which is an extension of the linearization order.

However, suppose we first merge versions v_2 and v_4 , to obtain the version v_5 , where the linearization relation is $e_3 \xrightarrow[\text{rc}]{\text{lo}} e_1$. Merging v_3 and v_5 (with LCA v_2) would have the same linearization relation as merging v_3 and v_4 . However, the sequences leading to v_3 and v_5 are e_1e_2 and e_3e_1 respectively, while the only sequence which extends the linearization relation for their merge is $e_1e_2e_3$. Bottom-up linearization will then be constrained to pick either e_1 or e_2 to appear at the end, but such a sequence will not extend the linearization relation resulting in the failure of bottom-up linearization. To avoid such cases, we place an additional constraint which prohibits the presence of an rc-chain:

$$\text{NO-RC-CHAIN}(\mathcal{D}) \triangleq \neg(\exists o_1, o_2, o_3 \in O. o_1 \xrightarrow{\text{rc}} o_2 \xrightarrow{\text{rc}} o_3)$$

If there is an rc-chain, executions such as Fig. 10 are possible, resulting in infeasibility of bottom-up linearization. However, we will show that if an MRDT satisfies $\text{NO-RC-CHAIN}(\mathcal{D})$, then we can use bottom-up linearization to prove that \mathcal{D} is linearizable. We note that NO-RC-CHAIN is a pragmatic restriction and consistent with standard conflict-resolution strategies such as add/remove-wins, enable/disable-wins, update/delete-wins, etc. which are typically used in MRDT implementations.

4 Verifying RA-linearizability of MRDTs

In this section, we present our verification strategy for proving RA-linearizability of MRDTs using bottom-up linearization. According to Def. 3.9, in order to prove that an MRDT \mathcal{D} is linearizable, we need to consider every configuration C reachable in any execution, and show that all replicas in C have states which can be obtained by linearizing the events applied to the replica, i.e. finding a sequence which obeys the linearization relation (Def. 3.5). We will assume that \mathcal{D} satisfies the three constraints (RC-NON-COMM, COND-COMM and NO-RC-CHAIN) necessary for an MRDT to be linearizable, and for bottom-up linearization to succeed.

Our overall proof strategy is to use induction on the length of the execution and to extract generic verification conditions (VCs) which help us to discharge the inductive case. These VCs would essentially be instantiations of the BOTTOMUPTEMPLATE rule, proving that the merge operation results in a linearization of the events of the two versions being merged. Proving these VCs for arbitrary MRDTs is not straightforward (as discussed in §2.3), and hence we propose another induction scheme over event sequences. We first discuss the instantiations of the BOTTOMUPTEMPLATE rule required for linearizing merges.

4.1 Linearizing Merge Operations

Consider an execution $\tau \in \llbracket \mathcal{S}_{\mathcal{D}} \rrbracket$ such that all configurations in τ are linearizable. Suppose τ ends in the configuration C . Now, we extend τ by one more transition, resulting in the new configuration C' ; we need to prove that C' is also linearizable. Let $C = \langle N, H, L, G, \text{vis} \rangle$, $C' = \langle N', H', L', G', \text{vis}' \rangle$. It is easy to see if that this transition is caused due to CREATEBRANCH or APPLY rules, then C' will be linearizable. For example, in the [APPLY] transition, where a new update operation o is applied on a replica r (generating a new event e), only the state at r changes, and this new state is obtained by directly applying e on the original state σ at r . Since σ was assumed to be linearizable, there exists a sequence π which extends $\text{lo}(C)_{|L(H(r))}$, with $\sigma = \pi(\sigma_0)$ (recall that $L(H(r))$ denotes the set of events applied at r). Then, the new state $e(\sigma)$ is clearly linearizable through the sequence πe which extends $\text{lo}(C')_{|L'(H'(r))}$.

We focus on the difficult case when there is a MERGE transition from C to C' which merges the replicas r_1 and r_2 . Let σ_1 and σ_2 be the states of the head versions v_1 and v_2 at r_1 and r_2 respectively. Let σ_\top be the state of the LCA version v_\top of v_1 and v_2 . Recall that $L(v_\top) = L(v_1) \cap L(v_2)$. The transition will install a new version with state $\sigma_m = \text{merge}(\sigma_\top, \sigma_1, \sigma_2)$ at the replica r_1 , leaving the other replicas unchanged. Also, $L'(v_m) = L(v_1) \cup L(v_2)$. We need to show that there exists a sequence π of events in $L'(v_m)$ such that π extends $\text{lo}(C')|_{L'(v_m)}$ and $\sigma_m = \pi(\sigma_0)$.

We first describe the structure of a sequence π which extends $\text{lo}(C')|_{L'(v_m)}$. For ease of readability, we use L_1 for $L(v_1)$, L_2 for $L(v_2)$ and L_\top for $L(v_\top)$, and lo_m for $\text{lo}(C')|_{L'(v_m)}$. We define the following sets of events:

$$\begin{aligned} L'_1 &= L_1 \setminus L_\top & L'_2 &= L_2 \setminus L_\top \\ L_1^b &= \{e \in L'_1 \mid \exists e_\top \in L_\top. (e \xrightarrow{\text{lo}_m} e_\top \vee \exists e' \in L'_1. e \xrightarrow{\text{lo}_m} e' \xrightarrow{\text{lo}_m} e_\top)\} \\ L_2^b &= \{e \in L'_2 \mid \exists e_\top \in L_\top. (e \xrightarrow{\text{lo}_m} e_\top \vee \exists e' \in L'_2. e \xrightarrow{\text{lo}_m} e' \xrightarrow{\text{lo}_m} e_\top)\} \\ L_\top^a &= \{e_\top \in L_\top \mid \exists e \in L_1^b \cup L_2^b. e \xrightarrow{\text{lo}_m} e_\top\} & L_1^a &= L'_1 \setminus L_1^b & L_2^a &= L'_2 \setminus L_2^b & L_\top^b &= L_\top \setminus L_\top^a \end{aligned}$$

L'_1 and L'_2 are the local events in each version. Note that any pair of events $e_1 \in L'_1, e_2 \in L'_2$ will necessarily be concurrent. This is because, in any reachable configuration, any version v is always **causally closed**, which means that if $e_1 \xrightarrow{\text{vis}} e_2$ and $e_2 \in L(v)$, then $e_1 \in L(v)$. Hence, for events $e_1 \in L'_1, e_2 \in L'_2$, if $e_1 \xrightarrow{\text{vis}} e_2$ then $e_1 \in L'_2$, which would make e_1 a non-local event (i.e. part of the LCA). Bottom-up linearization first linearizes the local events across the two versions using the rc relation for non-commutative events, and then linearizes events of the LCA. However, as demonstrated by the example in §2.4, local events may also need to be linearized before events of the LCA (due to possible intermediate merges), and these events are collected in the sets L_1^b and L_2^b . Specifically, $L_i^b (i = 1, 2)$ contains those local events e in L'_i which either occur lo_m before some event in the LCA, or which occur lo_m before another local event e' which occurs lo_m before an LCA event. The events of the LCA which need to be linearized after local events are collected in L_\top^a . Finally, L_1^a and L_2^a contain local events which do not occur lo_m before an LCA event.

Example 4.1. Consider the execution in Fig. 6, and the merge of versions v_3 and v_4 , for which the LCA is v_1 . For this merge, $L'_1 = \{e_3\}$, $L'_2 = \{e_2\}$, $L_1^b = \emptyset$, $L_2^b = \{e_2\}$, $L_\top^a = \{e_1\}$. For the merge of versions v_1 and v_2 (whose LCA is v_0), $L'_1 = \{e_1\}$, $L'_2 = \{e_2\}$, while L_1^b, L_2^b, L_\top^a will all be empty (since no local event comes lo -before an LCA event).

We now show that there exists a sequence π which extends lo_m and which has events in $S_1 = L_\top^b$ followed by $S_2 = L_\top^a \cup L_1^b \cup L_2^b$ followed by $S_3 = L_1^a \cup L_2^a$ (later, we will discuss the ordering of events inside each set S_i). To prove this, we will demonstrate that there is no lo_m from events in S_i to events in S_{i-1} . Based on the definitions of the S_i sets, we can deduce some obvious facts: (i) there cannot be events $e \in S_3, e' \in L_\top$ such that $e \xrightarrow{\text{lo}_m} e'$, because otherwise, such an event e would be in $L_1^b \cup L_2^b$ (and hence not in S_3), (ii) there cannot be events $e \in L_1^b \cup L_2^b, e' \in L_\top^b$ such that $e \xrightarrow{\text{lo}_m} e'$, because otherwise, such an event e' would be in L_\top^a . In addition, using NO-RC-CHAIN and RC-NON-COMM, we also prove the following:

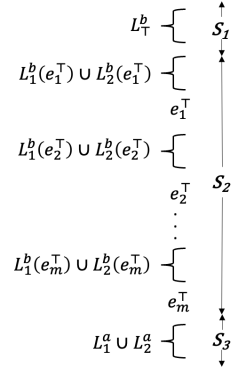


Fig. 11. Structure of sequence π extending lo_m

LEMMA 4.2. (1) For events $e \in L_1^a \cup L_2^a$, $e' \in L_1^b \cup L_2^b$, $\neg(e \xrightarrow{\text{lo}_m} e')$.

(2) For events $e \in L_\top^a$, $e' \in L_\top^b$, $\neg(e \xrightarrow{\text{lo}_m} e')$.

(1) from the above lemma ensures that there is no lo_m relation from S_3 to S_2 , while (2) ensures the same from S_2 to S_1 . Hence a sequence with the structure $S_1 S_2 S_3$ would extend lo_m . Let us now consider the ordering among events in each set. First, for S_3 , this set contains local events which are guaranteed to not come lo_m before any event of the LCA. An event in L_1^a will be concurrent with an event in L_2^a , and the linearization relation between them will depend upon the rc relation between the underlying operations (if the events don't commute). We now instantiate **BOTTOMUPTEMPLATE** for the case where both L_1^a and L_2^a are non-empty in the rule **BOTTOMUP-2-OP** in Fig. 12, so that the linearization needs to consider the rc relation between events in the two sets.

[BOTTOMUP-2-OP]	[BOTTOMUP-1-OP]	
$\frac{e_1 \neq e_2 \quad e_1 \xrightarrow{\text{rc}} e_2 \vee e_1 \rightleftharpoons e_2}{\text{merge}(l, e_1(a), e_2(b)) = e_2(\text{merge}(l, e_1(a), b))}$	$\frac{(e_\top \neq \epsilon \wedge e_1 \neq e_\top) \vee (e_\top = \epsilon \wedge l = b)}{\text{merge}(e_\top(l), e_1(a), e_\top(b)) = e_1(\text{merge}(e_\top(l), a, e_\top(b)))}$	
[BOTTOMUP-0-OP]	[MERGEIDEMPOTENCE]	[MERGECOMMUTATIVITY]
$\text{merge}(e_\top(l), e_\top(a), e_\top(b)) = e_\top(\text{merge}(l, a, b))$	$\text{merge}(a, a, a) = a$	$\text{merge}(l, a, b) = \text{merge}(l, b, a)$

Fig. 12. Bottom-up Linearization

Note that e_1, e_2 and l, a, b are all universally quantified. The **BOTTOMUP-2-OP** rule is an algebraic property of merge which needs to be separately shown for each MRDT implementation. For our case where we are trying to linearize $\text{merge}(\sigma_\top, \sigma_1, \sigma_2)$, we can apply **BOTTOMUP-2-OP** with $l = \sigma_\top$, $e_1(a) = \sigma_1$ and $e_2(b) = \sigma_2$. Note that since L_1^a and L_2^a are both non-empty, $e_1 \in L_1^a$, $e_2 \in L_2^b$ (in fact, e_1 and e_2 would be the maximal events in L_1^a and L_2^b according to lo_m). **BOTTOMUP-2-OP** would then linearize e_2 at the end of the sequence. If $e_1 \xrightarrow{\text{rc}} e_2$, then $e_1 \xrightarrow{\text{lo}_m} e_2$, and thus linearizing e_2 at the end obeys the lo_m ordering. Note that due to the **NO-RC-CHAIN** constraint, e_2 cannot come lo_m before another concurrent event e_3 . **BOTTOMUP-2-OP** can now be recursively applied on $\text{merge}(l, e_1(a), b)$, by considering e_1 and the last event leading to the state b . By repeatedly applying **BOTTOMUP-2-OP** all the remaining events in L_1^a and L_2^a can be linearized until one of the sets becomes empty.

Let us now consider the scenario where exactly one of L_1^a and L_2^a is empty. WLOG, let L_1^a be non-empty. We instantiate **BOTTOMUPTEMPLATE** for the case where L_1^a is non-empty and L_2^a is empty in the rule **BOTTOMUP-1-OP** in Fig. 12, so that the linearization orders all events of L_1^a after events of S_2 .

Let us consider the first clause in the premise where $e_\top \neq \epsilon$. To understand **BOTTOMUP-1-OP**, note that if L_2^a is empty, then all local events in L_2^a are linearized before the LCA events. In this case, the last event which leads to the state σ_2 must be an LCA event. **BOTTOMUP-1-OP** uses this observation, with $e_\top(l) = \sigma_\top$, $e_1(a) = \sigma_1$ and $e_\top(b) = \sigma_2$. Notice that the last event in both the LCA and the second argument to merge are exactly the same. e_\top will be the maximal event (according to lo_m relation) in L_\top^a , while e_1 will be the maximal event in L_1^a . **BOTTOMUP-1-OP** then linearizes e_1 at the end of the sequence, thus ensuring that all L_1^a events are linearized after events in S_1 and S_2 . It is possible that L_\top^a is empty, in which case L_2^a will be empty, which is covered by the second clause where $e_\top = \epsilon$ and $l = b$ since there is no local event in the second state.

Example 4.3. Referring to Example 4.1 for the execution in Fig. 6, recall that for the merge of v_3 and v_4 , we have $L_1^a = \{e_3\}$, $L_2^a = \emptyset$ and $L_\top = \{e_1\}$. **BOTTOMUP-1-OP** can be applied in this scenario to linearize e_3 at the end of the sequence.

BOTTOMUP-2-OP and BOTTOMUP-1-OP can thus be used to linearize all events in S_3 . Let us now consider S_2 , which contains both local events in $L_1^b \cup L_2^b$ and LCA events in L_\top^a . We first provide a more fine-grained structure of lo_m among events in the set S_2 . Let $L_\top^a = \{e_1^\top, \dots, e_m^\top\}$. For each e_i^\top , we collect all local events from L_1^b and L_2^b which need to be linearized before e_i^\top . For local events which need to be linearized before multiple e_i^\top s, we associate them with the smallest such i . We use $L_1^b(e_i^\top)$ and $L_2^b(e_i^\top)$ to denote these sets. Formally:

$$\forall e_i^\top \in L_\top^a. L_1^b(e_i^\top) = \{e \in L_1' \mid (\forall j. j < i \implies e \notin L_1^b(e_j^\top)) \wedge e \xrightarrow{\text{lo}_m} e_i^\top \vee \exists e' \in L_1'. e \xrightarrow{\text{lo}_m} e' \xrightarrow{\text{lo}_m} e_i^\top\}$$

$L_2^b(e_i^\top)$ is defined in a similar manner. We now prove the following lemma using NO-RC-CHAIN and RC-NON-COMM:

- LEMMA 4.4. (1) For all events $e_i^\top, e_j^\top \in L_\top^a$, where $L_\top^a = \{e_1^\top, \dots, e_m^\top\}$, $\neg(e_i^\top \xrightarrow{\text{lo}_m} e_j^\top)$
 (2) For events $e \in L_1^b(e_i^\top) \cup L_2^b(e_i^\top)$, $e' \in L_1^b(e_j^\top) \cup L_2^b(e_j^\top)$ where $j < i$, $\neg(e \xrightarrow{\text{lo}_m} e')$.

From (1) in the above lemma, since there is no lo_m relation among events in L_\top^a , consider the sequence $e_1^\top e_2^\top \dots e_m^\top$ as a starting point for the sequence of events in S_2 which extends lo_m . We then inject $L_1^b(e_i^\top) \cup L_2^b(e_i^\top)$ before each e_i^\top in the sequence $e_1^\top e_2^\top \dots e_m^\top$, as shown in Fig. 11. Note that in Fig. 11, we have only presented various segments of the sequence, with the ordering within those segments determined by vis and rc. By (2) in Lemma 4.4, we can show that such a sequence will extend lo_m among the events in S_2 .

To show that merge follows the sequence π for S_2 , we now instantiate BOTTOMUPTEMPLATE for the case where L_1^a and L_2^a are empty (i.e. S_3 has already been linearized) in the rule BOTTOM-0-OP in Fig. 12. Following the structure of π in Fig. 11, e_\top would be the event $e_m^\top \in L_\top^a$. Note that since e_m^\top is an LCA event, it will be present in both states being merged. BOTTOMUP-0-OP then allows this event to be linearized first at the end.

Example 4.5. Following on from Example 4.3 for the execution in Fig. 6 for the merge of v_3 and v_4 , after BOTTOMUP-1-OP is applied to linearize e_3 , the states to be merged would be the versions v_1 and v_4 (with LCA v_1), both of whose last operation is e_1 . Hence, BOTTOMUP-0-OP would be applicable, which would linearize e_1 .

After applying BOTTOMUP-0-OP to linearize the LCA event e_m^\top , we then need to linearize events in $L_1^b(e_m^\top) \cup L_2^b(e_m^\top)$. However, the event e_m^\top has already been linearized, so none of the events in $L_1^b(e_m^\top) \cup L_2^b(e_m^\top)$ appear lo_m after an LCA event. This scenario can now be handled using BOTTOMUP-2-OP (if both $L_1^b(e_m^\top)$ and $L_2^b(e_m^\top)$ are non-empty) or BOTTOMUP-1-OP (if one of 2 sets is empty). These rules will appropriately linearize the events in $L_1^b(e_m^\top) \cup L_2^b(e_m^\top)$ taking into account the rc relation for concurrent events and vis relation for non-concurrent events. Once $L_1^b(e_m^\top) \cup L_2^b(e_m^\top)$ becomes empty, we then encounter the next LCA event in L_\top^a , which can again be linearized using BOTTOMUP-0-OP.

The three instantiations of BOTTOMUPTEMPLATE can thus be repeatedly applied to linearize the rest of the events in S_2 . Following this, all the local events would have been linearized, leaving only the LCA events in S_1 . This would result in all three arguments to merge being equal, in which case we can use the MERGEIDEMPOTENCE rule in Fig. 12. Using MERGEIDEMPOTENCE, we can equate the output of merge to its argument, which has already been assumed to be appropriately linearized.

In order to avoid mirrored versions of BOTTOMUP-2-OP and BOTTOMUP-1-OP where the second and third arguments are swapped, we also require the MERGECOMMUTATIVITY property in Fig. 12. We now state our soundness theorem linking the various properties with RA-linearizability of MRDT:

THEOREM 4.6. *If an MRDT \mathcal{D} satisfies BOTTOMUP-2-OP, BOTTOMUP-1-OP, BOTTOMUP-0-OP, MERGEIDEMPOTENCE and MERGECOMMUTATIVITY, then \mathcal{D} is linearizable.*

The proof closely follows the informal arguments that we have presented in this sub-section, using induction on the size of the various sets $L_1^a, L_2^a, L_1^b \cup L_2^b, L_\top^a$.

4.2 Automated Verification

While we have identified the sufficient conditions to show RA-linearizability of an MRDT using bottom-up linearization, proving these conditions for arbitrary MRDTs is not straightforward. Further, while the BOTTOMUP-X-OP properties as shown in the previous sub-section had universal quantification over MRDT states l, a, b , in general, for proving RA-linearizability, we only need to show these properties for feasible states that may arise during an actual execution.

We now leverage the fact that the feasible states would have been obtained through linearization of the visible events at the respective versions. In particular, we can characterize the states on which merge can be invoked through the various events sets $L_1^a, L_2^a, L_1^b, L_2^b, L_\top^a, L_\top^b$ that we had defined in the previous sub-section. We only need to prove the BOTTOMUP-X-OP properties for states which have been obtained through linearizations of events in these event sets. For this purpose, we propose an induction scheme which establishes the required properties while traversing the event sets as depicted in Fig. 11 in a top-down fashion.

Table 1. Induction scheme for BOTTOMUPTEMPLATE. For clarity, we use \cdot for function composition, and μ for merge.

VC Name	Pre-condition	Post-condition
$\psi_{\text{base}}^{L_\top^b}$		$\mu(\pi_0(\sigma_0), \pi_1(\sigma_0), \pi_2(\sigma_0)) = \pi \cdot \mu(\pi'_0(\sigma_0), \pi'_1(\sigma_0), \pi'_2(\sigma_0))$
$\psi_{\text{ind}}^{L_\top^b}$	$\mu(\pi_0(l), \pi_1(l), \pi_2(l)) = \pi \cdot \mu(\pi'_0(l), \pi'_1(l), \pi'_2(l))$	$\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top(l), \pi_2 \cdot e_\top(l)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top(l), \pi'_2 \cdot e_\top(l))$
$\psi_{\text{ind}}^{L_\top^a}$	$\exists e. e \xrightarrow{rc} e_\top$	$\mu(\pi_0(l), \pi_1(a), \pi_2(b)) = \pi \cdot \mu(\pi'_0(l), \pi'_1(a), \pi'_2(b))$ $\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top(a), \pi_2 \cdot e_\top(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top(a), \pi'_2 \cdot e_\top(b))$
$\psi_{\text{ind}1}^{L_1^b}$	$e_b \xrightarrow{rc} e_\top$	$\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top(a), \pi_2 \cdot e_\top(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top(a), \pi'_2 \cdot e_\top(b))$ $\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top \cdot e_b(a), \pi_2 \cdot e_\top(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top \cdot e_b(a), \pi'_2 \cdot e_\top(b))$
$\psi_{\text{ind}2}^{L_1^b}$	$e_b \xrightarrow{rc} e_\top \wedge \neg e \not\Rightarrow e_b$	$\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top \cdot e_b(a), \pi_2 \cdot e_\top(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top \cdot e_b(a), \pi'_2 \cdot e_\top(b))$ $\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top \cdot e_b \cdot e(a), \pi_2 \cdot e_\top(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top \cdot e_b \cdot e(a), \pi'_2 \cdot e_\top(b))$
$\psi_{\text{ind}1}^{L_2^b}$	$e_b \xrightarrow{rc} e_\top$	$\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top(a), \pi_2 \cdot e_\top(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top(a), \pi'_2 \cdot e_\top(b))$ $\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top \cdot e_b(a), \pi_2 \cdot e_\top \cdot e_b(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top \cdot e_b(a), \pi'_2 \cdot e_\top \cdot e_b(b))$
$\psi_{\text{ind}2}^{L_2^b}$	$e_b \xrightarrow{rc} e_\top \wedge \neg e \not\Rightarrow e_b$	$\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top \cdot e_b(a), \pi_2 \cdot e_\top(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top \cdot e_b(a), \pi'_2 \cdot e_\top(b))$ $\mu(\pi_0 \cdot e_\top(l), \pi_1 \cdot e_\top \cdot e_b \cdot e(a), \pi_2 \cdot e_\top \cdot e_b \cdot e(b)) = \pi \cdot \mu(\pi'_0 \cdot e_\top(l), \pi'_1 \cdot e_\top \cdot e_b \cdot e(a), \pi'_2 \cdot e_\top \cdot e_b \cdot e(b))$
$\psi_{\text{ind}}^{L_1^a}$		$\mu(\pi_0(l), \pi_1(a), \pi_2(b)) = \pi \cdot \mu(\pi'_0(l), \pi'_1(a), \pi'_2(b))$ $\mu(\pi_0(l), \pi_1 \cdot e(a), \pi_2(b)) = \pi \cdot \mu(\pi'_0(l), \pi'_1 \cdot e(a), \pi'_2(b))$
$\psi_{\text{ind}}^{L_2^a}$		$\mu(\pi_0(l), \pi_1(a), \pi_2(b)) = \pi \cdot \mu(\pi'_0(l), \pi'_1(a), \pi'_2(b))$ $\mu(\pi_0(l), \pi_1(a), \pi_2 \cdot e(b)) = \pi \cdot \mu(\pi'_0(l), \pi'_1(a), \pi'_2 \cdot e(b))$

Here, we present the induction scheme for the generic BOTTOMUPTEMPLATE rule. The scheme can then be instantiated for all the three BOTTOMUP-X-OP rules. Table 1 contains the verification conditions corresponding to the base case and inductive case over the different event sets. Every VC has the form (pre-condition \implies post-condition), and all variables are universally quantified. Our goal is to show the BOTTOMUPTEMPLATE rule for all feasible MRDT states l, a, b , where l is the state of the LCA of a and b . Let L_\top, L_1, L_2 be the event sets corresponding to l, a, b respectively. We define the event sets $L_1^a, L_2^a, L_1^b, L_2^b, L_\top^a, L_\top^b$ in exactly the same manner as the previous sub-section, based on the linearization relation of the configuration obtained by the merge(l, a, b) transition.

Note that the events in π_0, π_1, π_2 (used in the `BOTTOMUPTEMPLATE` rule) would also come from the above event sets, but in the following discussion, we freeze these events, i.e. all our assertions about the events sets will be modulo these events.

We start with the VC $\psi_{\text{base}}^{L_\tau^b}$, which corresponds to the case where every event set is empty. There is no pre-condition, and the post-condition requires `BOTTOMUPTEMPLATE` to hold on the initial MRDT state σ_0 . For example, for the `BOTTOMUP-2-OP` rule, $\psi_{\text{base}}^{L_\tau^b}$ VC would be $\text{merge}(\sigma_0, e_1(\sigma_0), e_2(\sigma_0)) = e_2(\text{merge}(\sigma_0, e_1(\sigma_0), \sigma_0))$, where $e_1 \xrightarrow{\text{rc}} e_2$ or e_1 and e_2 commute. Notice that e_1 and e_2 would be events in L_1^a and L_2^a , and our assertion about all event sets being empty is modulo these events.

Next, the VC $\psi_{\text{ind}}^{L_\tau^b}$ corresponds to the inductive case on L_τ^b , where we assume every event set except L_τ^b to be empty. The pre-condition corresponds to the inductive hypothesis, where we assume the property to hold for some event set L_τ^b , and the post-condition asserts that the property holds while adding another event e_τ to L_τ^b . Recall that L_τ^b corresponds to the LCA events which come lo before all local events. Since all the other event sets are empty, this translates to the same state l for all the three arguments to merge in the pre-condition, and applying the LCA event e_τ to all three arguments in the post-condition.

Next, we induct on the set L_τ^a , i.e. the set of LCA events which occur lo after a local event. The base case, where $|L_\tau^a| = \emptyset$ exactly corresponds to the result of the induction on L_τ^b . The inductive case is covered by the VC $\psi_{\text{ind}}^{L_\tau^a}$, which adds an LCA event e_τ to all three arguments of merge from pre-condition to post-condition. Notice that we also have another pre-condition which requires the existence of some event e which should come rc-before e_τ , which is necessary for e_τ to be in L_τ^a . The post-condition just adds a new LCA event e_τ . The events in $L_1^b(e_\tau)$ and $L_2^b(e_\tau)$ will be added by the next 4 VCs.

$\psi_{\text{ind1}}^{L_1^b}$ and $\psi_{\text{ind2}}^{L_1^b}$ add an event in L_1^b from the pre-condition to the post-condition. $\psi_{\text{ind1}}^{L_1^b}$ considers an event e_b which occurs rc-before the LCA event e_τ . Notice that the pre-condition of $\psi_{\text{ind1}}^{L_1^b}$ is exactly the same as the post-condition of $\psi_{\text{ind}}^{L_\tau^a}$. In the post-condition of $\psi_{\text{ind1}}^{L_1^b}$, the event e_b is applied before e_τ on the argument a to merge, thus reflecting that this is an event in L_1^b . $\psi_{\text{ind2}}^{L_1^b}$ adds an event $e \in L_1^b$ which does not commute with an existing event $e_b \in L_1^b$ (see the definition of L_1^b). $\psi_{\text{ind1}}^{L_2^b}$ and $\psi_{\text{ind2}}^{L_2^b}$ are analogous and do the same thing for the argument b to merge.

Finally, $\psi_{\text{ind}}^{L_1^a}$ and $\psi_{\text{ind}}^{L_2^a}$ add events from L_1^a and L_2^a . The base cases for the two sets would exactly correspond to the result of the induction carried out so far on the rest of the event sets. For the inductive case, in $\psi_{\text{ind}}^{L_1^a}$ (resp. $\psi_{\text{ind}}^{L_2^a}$), a new event e is added on the second argument a (resp. third argument b) from the pre-condition to the post-condition. This establishes the rule `BOTTOMUPTEMPLATE` for any feasible input arguments to merge during any execution. We denote the set of VCs in Table 1 by $\psi^*(\text{BOTTOMUPTEMPLATE})$.

THEOREM 4.7. *If an MRDT \mathcal{D} satisfies the VCs $\psi^*(\text{BOTTOMUP-2-OP})$, $\psi^*(\text{BOTTOMUP-1-OP})$, $\psi^*(\text{BOTTOMUP-0-OP})$, `MERGEIDEMPOTENCE` and `MERGECOMMUTATIVITY`, then \mathcal{D} is linearizable.*

5 Experimental Evaluation

We have implemented our verification technique in the F^{*} programming language and verified several MRDTs using it. We also extracted OCaml code from the verified implementations and ran them as part of Irmin [9], a Git-like distributed database which follows the MRDT system model described in §3. This distinguishes our work from prior works in automated RDT verification [16] which focuses on verifying abstract models rather than actual implementations.

Our framework in F^* consists of an F^* interface that defines signatures for an MRDT implementation (Fig. 2) and the VCs described in Table 1; these are encoded as F^* lemmas. This interface contains 200 lines of F^* code. An MRDT developer instantiates the interface with their specific MRDT implementation and calls upon F^* to prove the lemmas (i.e., the VCs). Once this is done, our metatheory (Theorem 4.7) guarantees that the MRDT implementation is linearizable.

Table 2. Verified MRDTs. * denotes MRDT implementations not present in prior work.

MRDT	rc Policy	#LOC	Verification Time (s)
Increment-only counter [12]	none	6	0.72
PN counter [23]	none	10	1.64
Enable-wins flag*	disable \xrightarrow{rc} enable	30	29.80
Disable-wins flag*	enable \xrightarrow{rc} disable	30	37.91
Grows-only set [12]	none	6	0.45
Grows-only map [23]	none	11	4.65
OR-set [23]	rem _a \xrightarrow{rc} add _a	20	4.53
OR-set (efficient)*	rem _a \xrightarrow{rc} add _a	34	660.00
Remove-wins set*	add _a \xrightarrow{rc} rem _a	22	9.60
Set-wins map*	del _k \xrightarrow{rc} set _k	20	5.06
Replicated Growable Array [1]	none	13	1.51
Optional register*	unset \xrightarrow{rc} set	35	200.00
Multi-valued Register*	none	7	0.65
JSON-style MRDT*	Fig. 13	26	148.84

We instantiate the interface with MRDT implementations of several datatypes such as counter, flag, set, map, and list (Table 2). All the results were obtained on a Intel®Xeon®Gold 5120 x86-64 machine running Ubuntu 22.04 with 64GB of main memory. While some of the MRDTs have been taken from previous works [1, 12, 23] or translated from their CRDT counterparts, we also develop some new implementations, denoted by * in Table 2. We also uncovered bugs in previous MRDT implementations (Enable-wins flag and Efficient OR-set) from [23], which we fixed (more details in §5.2). We note that in all our experiments, all the VCs were automatically discharged by F^* in a reasonable amount of time.

While our approach ensures that the MRDT implementations are verified in the F^* framework, it is important to note that the user is obligated to trust the F^* language implementation, the extraction mechanism, the OCaml language implementation, the OCaml runtime, and the hardware.

We now highlight several notable features about our verified MRDTs. We have designed and developed the first correct implementations of both an enable-wins and disable-wins flag MRDT. Our implementation of efficient OR-set maintains a per-replica, per-element counter instead of adding different versions of the same element (as done by the OR-set implementation of Fig. 2), thus matching the theoretical lower bound in terms of space-efficiency for any OR-set CRDT implementation (as proved in [4]). We have developed the first known MRDT implementation of a remove-wins set datatype. Finally, as a demonstration of vertical compositionality, we have developed a JSON MRDT which is composed of several component MRDTs, with its correctness guarantee being directly derived from the correctness of the component MRDTs.

5.1 Case study: A verified polymorphic JSON-style MRDT

JSON is a notable example of a data type which is composed of several other datatypes. JSON is widely used as a data interchange format in many databases and web services [10]. Our JSON MRDT is modeled as an unordered collection of key/value pairs, where the values can be any primitive types, such as counter, list, etc., or they can be JSON type themselves. We assume that keys are

update-only; that is, key-value mappings can be added and modified, but once a key is added, it cannot be deleted. Previous works, such as Automerge [2], have developed similar JSON-style CRDT models. However, these models are monomorphic, which means that the data type of the values must be known in advance. Our goal is to develop a more generic JSON-style MRDT that supports polymorphic values, i.e. we leave the value data type as an abstract type which can be instantiated with any concrete MRDT.

```

1:  $\Sigma_{\text{json}} : (k : (\text{string} \times \Omega)) \rightarrow \Sigma_{\text{snd}(k)}$ 
2:  $O_{\text{json}} = \{\text{set}(k, o) \mid o \in O_{\text{snd}(k)}\}$ 
3:  $Q_{\text{json}} = \{\text{get}(k, q) \mid q \in Q_{\text{snd}(k)}\}$ 
4:  $\sigma_{0_{\text{json}}} = \lambda(k : \text{string} \times \Omega). \sigma_{0_{\text{snd}(k)}}$ 
5:  $\text{do}(\sigma, t, r, \text{set}(k, o)) = \sigma[k \mapsto o(\sigma(k), t, r)]$ 
6:  $\text{merge}_{\text{json}}(\sigma_1, \sigma_2) = \lambda(k : \text{string} \times \Omega). \text{merge}_{\text{snd}(k)}(\sigma_1(k), \sigma_2(k))$ 
7:  $\text{query}_{\text{json}}(\sigma, \text{get}(k, q)) = \text{query}_{\text{snd}(k)}(\sigma(k), q)$ 
8:  $\text{rc}_{\text{json}} = \{(\text{set}(k_1, o_1), \text{set}(k_2, o_2)) \in O_{\text{json}} \times O_{\text{json}} \mid k_1 = k_2 \wedge (o_1, o_2) \in \text{rc}_{\text{snd}(k_1)}\}$ 

```

Fig. 13. JSON-style MRDT implementation

Fig. 13 shows the implementation of the JSON MRDT. It uses a map to maintain the association between keys and values. Notice that the key is a tuple consisting of the identifier string and an MRDT type $\alpha \in \Omega$ which denotes the type of the value. The type α can be any arbitrary MRDT with implementation $\mathcal{D}_\alpha = (\Sigma_\alpha, \sigma_{0_\alpha}, \text{merge}_\alpha, \text{query}_\alpha, \text{rc}_\alpha)$. Different key strings can now map to different value MRDT types. We also allow overloading: the same key string can be associated with multiple values of different types. The JSON MRDT allows update operations of the form $\text{set}(k, o)$ where o is an operation of the underlying value MRDT associated with the key k . $\text{set}(k, o)$ simply applies the operation o on the value associated with k , leaving the other key-value pairs unchanged. The JSON merge calls the underlying MRDT merge on the values associated with each key. The query operation of the form $\text{get}(k, q)$ retrieves the value associated with k in σ and applies the query operation q of the underlying data type to it. The conflict resolution policy of JSON operations (rc_{json}) depends on the conflict resolution of the value types when two operations update the same key (i.e. same identifier and value type). Every other pair of JSON operations commute with each other.

Notably, the proof of RA-linearizability of the JSON MRDT is directly derived from the proofs of the underlying value MRDT types. If all the MRDTs in Ω are linearizable, then the JSON MRDT is also linearizable. We have proved all the VCs for the JSON MRDT in F^* by using the VCs of the underlying value MRDTs. We can now instantiate Ω with any set of verified MRDTs, thereby obtaining the verified JSON MRDT for free.

5.2 Buggy MRDT Implementation in [23]

We now present some details of one of the buggy MRDTs, Enable-wins flag, that we discovered using our framework in the work by Soundarapandian et al. [23]. The state of the enable-wins flag MRDT consists of a pair: a counter and a flag. The counter tracks the number of the enable events, while the flag is set to true on an enable event. The desired specification for this flag is that it should be true when there is at least one enable event not visible to any disable event. In our framework, we can express this specification as $\text{disable} \xrightarrow{\text{rc}} \text{enable}$, linearizing the enable operation after a concurrent disable. When we attempted to verify this implementation in our framework, we discovered that one of the VCs, $\psi_{\text{ind2-top}}^{L_2^b}$, was failing. Our investigation revealed that the implementation violated the specification. The bug appeared in an execution with intermediate merges.

Consider the execution depicted in Fig. 14. When merging versions v_3 and v_5 (with LCA v_1), since the counter value of v_5 is greater than v_1 , the flag in the merged version v_6 is set to true. However, this contradicts the Enable-wins flag specification, which states that the flag should be true only when there is an enable event that is not visible to any disable event. All enable events in the execution are disabled by subsequent disable events on their individual replicas, yet the flag is true at v_6 . Notice that the version v_5 is obtained due to an intermediate merge. We discovered that Soundarapandian et al. [23] had an implementation bug in the framework. The framework expects a simulation relation from the MRDT developer, in addition to the specification and the implementation. This simulation relation serves as a proof artefact. Soundarapandian et al. [23] check whether the developer-provided simulation relation is valid and the bug occurred during the validity-checking procedure. Due to this, Soundarapandian et al. [23] admitted the buggy enable-wins flag implementation⁵.

We further note that this buggy implementation does not even satisfy strong eventual consistency. In Fig. 14, merging v_3 and v_4 results in v_7 , where the flag is false. Note that both versions v_6 and v_7 have observed the same set of updates on both replicas, yet they lead to divergent states. This violates strong eventual consistency. We fixed this implementation by maintaining a counter-flag pair for every replica, i.e. changing the state to a map from replica-IDs to counter-flag pair.

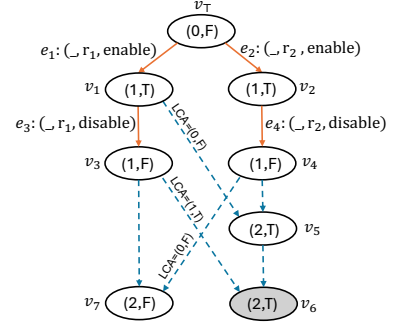


Fig. 14. An enable-wins flag execution

5.3 Verifying state-based CRDTs

Although the development in the paper so far has focused on verifying MRDTs, we note that our framework can also directly verify state-based CRDTs. The only difference between the two is that state-based CRDTs do not maintain the LCA, and merge is a binary function. Our VCs (Table 1) can be directly applied on state-based CRDTs, by simply ignoring the LCA argument for all merges. Note that while the merge function in state-based CRDTs does not use the LCA, our VCs still use the LCA to determine whether an event is local or common to both replicas, and appropriately linearize events taking into account both rc and vis relations. The entire set of VCs retrofitted for state-based CRDTs can be found in Table 4 of the extended version [25] of the paper. We have also successfully implemented and verified 7 state-based CRDTs in our framework: Increment-only counter, PN counter, Observed-Remove set, Two-Phase set, Grows-only set, Grows-only map and Multi-valued register.

5.4 Limitations

Our framework is currently unable to verify some MRDT implementations such as Queue from previous works [12, 23]. The Queue MRDT follows at-least-once semantics for dequeues, which allows concurrent dequeue operations to return the same element from the queue, thereby having the effect of a single dequeue. Such an implementation is clearly not linearizable as per our definition, since we cannot omit any event while constructing the linearization. It would be possible to modify our notion of linearization to also allow events to be omitted; we leave this investigation as part of future work. Our verification technique is also not complete, but in practice we have been able to successfully verify all MRDT implementations (except Queue) from earlier works.

⁵Buggy implementation can be found in Appendix §A.3 of the extended version [25] of the paper.

6 Related Work and Conclusion

Reconciling concurrent updates is a challenging problem in distributed systems. CRDTs [3, 20, 22] (and more recently MRDTs) have emerged as a principled approach for building correct and efficient replicated implementations. Numerous works have focused on specifying and verifying CRDTs [1, 4, 7, 8, 13, 15–18, 28, 29]. Op-based CRDTs have a considerably different system model than MRDTs, where every operation instance at a replica is individually sent to other replicas. Hence, verification efforts targeting them [7, 15–17, 28] are mostly orthogonal to our work.

The system model of state-based CRDTs is similar to MRDTs, as it also requires a merge function to be implemented for reconciling concurrent updates. However, state-based CRDTs have stricter requirements for convergence and consistency: CRDT states must form a join-semilattice, updates must be monotonic, and the merge function must be the lattice join operation. The three algebraic properties of a semilattice: idempotence, commutativity, and associativity guarantee convergence.

Some CRDT works focus solely on ensuring convergence without addressing functional correctness. For instance, Porre et al. [18] do not fully capture the user intent when verifying state-based CRDTs. Consider a Counter CRDT with only an increment operation and an *incorrect* merge function that ignores its input states and always returns 0. Such an implementation is still convergent. However, it clearly does not capture the developer intent, which is that the value of the counter should be equal to the number of increment operations. Functional correctness is as important as convergence for replicated data types. Our framework addresses both by couching both in terms of RA-linearizability. We will flag the above implementation as incorrect, since the state after merge cannot be obtained by linearizing the operations performed on both the replicas.

In the context of CRDTs, Wang et al. [28] proposed the notion of replication-aware linearizability, which requires all replicas to have a state which can be obtained by linearizing the update operations visible to the replica according to the sequential specification. However, they do not propose any automated verification methodology for RA-linearizability. Further, though the main paper Wang et al. [28] focuses on op-based CRDTs, the extended version Enea et al. [5] does address state-based CRDTs, but they also require a semi-lattice-based formulation of the CRDT states for proving RA-linearizability.

A few works [11, 23] have explored the problem of verifying MRDT implementations. Kaki et al. [11] only focus on verifying convergence, but not functional correctness. Moreover, they significantly restrict the underlying system model by synchronizing all merge operations, which as mentioned in the paper itself could lead to longer convergence times. Soundarapandian et al. [23] verify both convergence and functional correctness, without requiring synchronized merges. However, their approach is not fully automated, and requires developers to provide a simulation relation linking concrete MRDT states with an abstract state which is based on a event-based declarative model. Their specification language is also based on an event-based model and is not very intuitive or developer-friendly. A few MRDT implementations from [23] were found to be buggy, and these errors were due to faulty simulation relations.

To conclude, in this work, we present the first, fully-automated verification methodology for MRDTs. We introduce the notion of replication-aware linearizability for MRDTs, as well as a simple specification framework based on ordering non-commutative update operations. We identify certain restrictions on the specification to ensure existence of a consistent linearization. We then leverage the definition of replication-aware linearizability to propose an automated verification methodology based on induction on operation sequences. We have successfully applied the technique on a number of complex MRDTs. While the foundations have been laid in this work, we believe there is a lot of scope for enriching the technique even further by considering more complex linearization strategies.

7 Data-Availability Statement

The artifact supporting this paper is available on Zenodo [24] as well as in our GitHub repository [26]. It provides our framework in the F^{*} programming language that allows implementing MRDTs/CRDTs and automatically proving the verification conditions (VCs) required by our technique.

References

- [1] Hagit Attiya, Sebastian Burckhardt, Alexey Gotsman, Adam Morrison, Hongseok Yang, and Marek Zawirski. 2016. Specification and Complexity of Collaborative Text Editing. In *Proceedings of the 2016 ACM Symposium on Principles of Distributed Computing* (Chicago, Illinois, USA) (PODC '16). Association for Computing Machinery, New York, NY, USA, 259–268. doi:10.1145/2933057.2933090
- [2] Automerge. 2022. Automerge. <https://automerge.org/>
- [3] Annette Bieniusa, Marek Zawirski, Nuno M. Preguiça, Marc Shapiro, Carlos Baquero, Valter Balesgas, and Sérgio Duarte. 2012. An optimized conflict-free replicated set. *CoRR* abs/1210.3368 (2012). arXiv:1210.3368 <http://arxiv.org/abs/1210.3368>
- [4] Sebastian Burckhardt, Alexey Gotsman, Hongseok Yang, and Marek Zawirski. 2014. Replicated data types: specification, verification, optimality. In *Proceedings of the 41st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages* (San Diego, California, USA) (POPL '14). Association for Computing Machinery, New York, NY, USA, 271–284. doi:10.1145/2535838.2535848
- [5] Constantin Enea, Suha Orhun Mutluergil, Gustavo Petri, and Chao Wang. 2019. Replication-Aware Linearizability. *CoRR* abs/1903.06560 (2019). arXiv:1903.06560 <http://arxiv.org/abs/1903.06560>
- [6] Git. 2021. Git: A distributed version control system. <https://git-scm.com/>
- [7] Victor B. F. Gomes, Martin Kleppmann, Dominic P. Mulligan, and Alastair R. Beresford. 2017. Verifying strong eventual consistency in distributed systems. *Proc. ACM Program. Lang.* 1, OOPSLA, Article 109 (Oct. 2017), 28 pages. doi:10.1145/3133933
- [8] Alexey Gotsman, Hongseok Yang, Carla Ferreira, Mahsa Najafzadeh, and Marc Shapiro. 2016. 'Cause I'm strong enough: reasoning about consistency choices in distributed systems. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2016, St. Petersburg, FL, USA, January 20 - 22, 2016*, Rastislav Bodik and Rupak Majumdar (Eds.). ACM, 371–384. doi:10.1145/2837614.2837625
- [9] Irmin. 2021. Irmin: A distributed database built on the principles of Git. <https://irmin.org/>
- [10] Json. [n. d.]. Json: A lightweight data-interchange format. <https://www.json.org/>
- [11] Gowtham Kaki, Prasanth Prahladan, and Nicholas V. Lewchenko. 2022. RunTime-assisted convergence in replicated data types. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation* (San Diego, CA, USA) (PLDI 2022). Association for Computing Machinery, New York, NY, USA, 364–378. doi:10.1145/3519939.3523724
- [12] Gowtham Kaki, Swarn Priya, KC Sivaramakrishnan, and Suresh Jagannathan. 2019. Mergeable replicated data types. *Proc. ACM Program. Lang.* 3, OOPSLA, Article 154 (Oct. 2019), 29 pages. doi:10.1145/3360580
- [13] Shadaj Laddad, Conor Power, Mae Milano, Alvin Cheung, and Joseph M. Hellerstein. 2022. Katara: synthesizing CRDTs with verified lifting. *Proc. ACM Program. Lang.* 6, OOPSLA2, Article 173 (Oct. 2022), 29 pages. doi:10.1145/3563336
- [14] Leslie Lamport. 1978. Time, Clocks, and the Ordering of Events in a Distributed System. *Commun. ACM* 21, 7 (1978), 558–565. doi:10.1145/359545.359563
- [15] Yiyun Liu, James Parker, Patrick Redmond, Lindsey Kuper, Michael Hicks, and Niki Vazou. 2020. Verifying replicated data types with typeclass refinements in Liquid Haskell. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 216 (Nov. 2020), 30 pages. doi:10.1145/3428284
- [16] Kartik Nagar and Suresh Jagannathan. 2019. Automated Parameterized Verification of CRDTs. In *Computer Aided Verification - 31st International Conference, CAV 2019, New York City, NY, USA, July 15-18, 2019, Proceedings, Part II (Lecture Notes in Computer Science, Vol. 11562)*, Isil Dillig and Serdar Tasiran (Eds.). Springer, 459–477. doi:10.1007/978-3-030-25543-5_26
- [17] Sreeja S. Nair, Gustavo Petri, and Marc Shapiro. 2020. Proving the Safety of Highly-Available Distributed Objects. In *Programming Languages and Systems - 29th European Symposium on Programming, ESOP 2020, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2020, Dublin, Ireland, April 25-30, 2020, Proceedings (Lecture Notes in Computer Science, Vol. 12075)*, Peter Müller (Ed.). Springer, 544–571. doi:10.1007/978-3-030-44914-8_20
- [18] Kevin De Porre, Carla Ferreira, and Elisa Gonzalez Boix. 2023. VeriFx: Correct Replicated Data Types for the Masses. In *37th European Conference on Object-Oriented Programming, ECOOP 2023, July 17-21, 2023, Seattle, Washington, United States (LIPICs, Vol. 263)*, Karim Ali and Guido Salvaneschi (Eds.). Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 9:1–9:45. doi:10.4230/LIPICs.ECOOP.2023.9
- [19] Riak. 2021. Resilient NoSQL Databases. <https://riak.com/>

- [20] Hyun-Gul Roh, Myeongjae Jeon, Jin-Soo Kim, and Joonwon Lee. 2011. Replicated abstract data types: Building blocks for collaborative applications. *J. Parallel Distrib. Comput.* 71, 3 (March 2011), 354–368. doi:10.1016/j.jpdc.2010.12.006
- [21] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. 2011. *A comprehensive study of Convergent and Commutative Replicated Data Types*. Research Report RR-7506. Inria – Centre Paris-Rocquencourt ; INRIA. 50 pages. <https://inria.hal.science/inria-00555588>
- [22] Marc Shapiro, Nuno M. Preguiça, Carlos Baquero, and Marek Zawirski. 2011. Conflict-Free Replicated Data Types. In *Stabilization, Safety, and Security of Distributed Systems - 13th International Symposium, SSS 2011, Grenoble, France, October 10-12, 2011. Proceedings (Lecture Notes in Computer Science, Vol. 6976)*, Xavier Défago, Franck Petit, and Vincent Villain (Eds.). Springer, 386–400. doi:10.1007/978-3-642-24550-3_29
- [23] Vimala Soundarapandian, Adharsh Kamath, Kartik Nagar, and KC Sivaramakrishnan. 2022. Certified Mergeable Replicated Data Types. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation (San Diego, CA, USA) (PLDI 2022)*. Association for Computing Machinery, New York, NY, USA, 332–347. doi:10.1145/3519939.3523735
- [24] Vimala Soundarapandian, Kartik Nagar, Aseem Rastogi, and KC Sivaramakrishnan. 2025. Automatically Verifying Replication-aware Linearizability (Artifact). <https://doi.org/10.5281/zenodo.14922118>.
- [25] Vimala Soundarapandian, Kartik Nagar, Aseem Rastogi, and KC Sivaramakrishnan. 2025. Automatically Verifying Replication-aware Linearizability (Extended version). arXiv:2502.19967 [cs.PL] <https://arxiv.org/abs/2502.19967>
- [26] Vimala Soundarapandian, Kartik Nagar, Aseem Rastogi, and KC Sivaramakrishnan. 2025. Automatically Verifying Replication-aware Linearizability (GitHub Repository). <https://github.com/prismlab/certified-mrds>.
- [27] Nikhil Swamy, Cătălin Hrițcu, Chantal Keller, Aseem Rastogi, Antoine Delignat-Lavaud, Simon Forest, Karthikeyan Bhargavan, Cédric Fournet, Pierre-Yves Strub, Markulf Kohlweiss, Jean-Karim Zinzindohoue, and Santiago Zanella-Béguelin. 2016. Dependent types and multi-monadic effects in F*. *SIGPLAN Not.* 51, 1 (Jan. 2016), 256–270. doi:10.1145/2914770.2837655
- [28] Chao Wang, Constantin Enea, Suha Orhun Mutluergil, and Gustavo Petri. 2019. Replication-aware linearizability. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation (Phoenix, AZ, USA) (PLDI 2019)*. Association for Computing Machinery, New York, NY, USA, 980–993. doi:10.1145/3314221.3314617
- [29] Peter Zeller, Annette Bieniusa, and Arnd Poetzsch-Heffter. 2014. Formal Specification and Verification of CRDTs. In *Formal Techniques for Distributed Objects, Components, and Systems - 34th IFIP WG 6.1 International Conference, FORTE 2014, Held as Part of the 9th International Federated Conference on Distributed Computing Techniques, DisCoTec 2014, Berlin, Germany, June 3-5, 2014. Proceedings (Lecture Notes in Computer Science, Vol. 8461)*, Erika Ábrahám and Catuscia Palamidessi (Eds.). Springer, 33–48. doi:10.1007/978-3-662-43613-4_3

Received 2024-10-15; accepted 2025-02-18