

Specification and Complexity of Collaborative Text Editing

Hagit Attiya
Technion

Sebastian Burckhardt
Microsoft Research

Alexey Gotsman
IMDEA Software Institute

Adam Morrison
Technion

Hongseok Yang
University of Oxford

Marek Zawirski*
Inria & Sorbonne Universités,
UPMC Univ Paris 06, LIP6

ABSTRACT

Collaborative text editing systems allow users to concurrently edit a shared document, inserting and deleting elements (e.g., characters or lines). There are a number of protocols for collaborative text editing, but so far there has been no precise specification of their desired behavior, and several of these protocols have been shown not to satisfy even basic expectations. This paper provides a precise specification of a replicated *list* object, which models the core functionality of replicated systems for collaborative text editing. We define a *strong* list specification, which we prove is implemented by an existing protocol, as well as a *weak* list specification, which admits additional protocol behaviors.

A major factor determining the efficiency and practical feasibility of a collaborative text editing protocol is the space overhead of the metadata that the protocol must maintain to ensure correctness. We show that for a large class of list protocols, implementing *either the strong or the weak list specification* requires a metadata overhead that is at least *linear* in the number of elements deleted from the list. The class of protocols to which this lower bound applies includes all list protocols that we are aware of, and we show that one of these protocols almost matches the bound.

Keywords

Collaborative text editing; eventual consistency

1. INTRODUCTION

Collaborative text editing systems, like Google Docs [7, 8], Apache Wave [1], or wikis [16], allow users at multiple sites to concurrently edit the same document. To achieve high responsiveness and availability, such systems often replicate the document in geographically distributed sites or on user devices. A user can modify the document at a nearby replica, which propagates the modifications to other replicas asynchronously. This propagation can be done either via a *centralized server* or *peer-to-peer*. An essential feature of a collaborative editing system is that all changes eventually propagate to all replicas and get incorporated into the docu-

*Now at Google.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

PODC'16, July 25–28, 2016, Chicago, IL, USA.

© 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-3964-3/16/07...\$15.00

DOI: <http://dx.doi.org/10.1145/2933057.2933090>

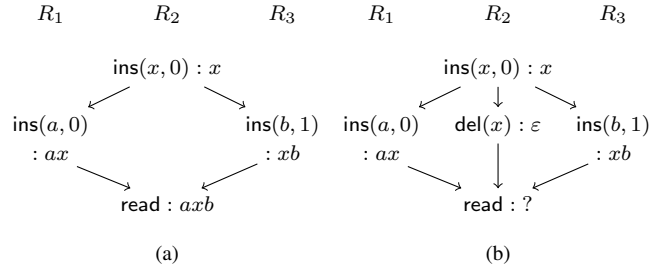


Figure 1: Example scenarios of collaborative text editing. Events are presented in format “operation : return value”. An arrow from an event e' to an event e expresses that the effects of e' get incorporated at e 's replica before e executes.

ment in a consistent way. In particular, such systems aim to guarantee *eventual consistency*: if users stop modifying the document, then the replicas will eventually converge to the same state [29,30].

Figure 1(a) gives an example scenario of a document edited at several replicas. First, replica R_2 inserts x at the first position (zero-indexed) into the empty list. This insertion then propagates to replica R_1 , which inserts a to the left of x , and to R_3 , which inserts b to the right of x . Later the modifications made by R_1 and R_3 propagate to all other replicas, including R_2 ; when the latter reads the list, it observes axb . In this scenario, the desired system behavior is straightforward, but sometimes this is not the case. To illustrate, consider the scenario in Figure 1(b), where R_2 deletes x from the list before the insertions of a and b propagate to it. One might expect the read by R_2 to return ab , given the orderings ax and xb established at other replicas. However, some implementations allow ba as a response; e.g., this is the case in a Jupiter protocol [20], used in public collaboration systems [31].

There have been a number of proposals of highly available collaborative editing protocols, using techniques such as *operation transformations* [13, 23, 27, 28] and *replicated data types* (aka CRDTs) [22, 24, 32]. However, specifications of their desired behavior [17,28] have so far been informal and imprecise, and several of the protocols have been shown not to satisfy even the basic expectation of eventual consistency [14]. To address this problem, we introduce a precise specification of a replicated list object, which allows its clients to insert and delete elements into the list at different replicas and thereby captures the core aspects of collaborative text editing [13] (Section 3). Our specification has two flavors. The *strong* specification ensures that orderings relative to deleted elements hold even after the deletion, thereby disallowing the response ba for the read in Figure 1(b). The *weak* specification provides no such guarantee, while still requiring the ordering between elements

that are not deleted to be consistent across the system. We show that both of these specifications ensure eventual consistency.

We prove that the strong specification is correctly implemented by a variant of an existing RGA (Replicated Growable Array) protocol [24], which is in the style of replicated data types [22] (Section 4). The protocol represents the list as a tree, with read operations traversing the tree in a deterministic order. Inserting an element a right after an element x (as in Figure 1(b)) adds a as a child of x in the tree. Deleting an element x just marks it as such; the node of x is left in the tree, creating a so-called *tombstone*. Keeping the tombstone enables the protocol to correctly incorporate insertions of elements received from other replicas that are ordered right after x (e.g., that of b in Figure 1(b)).

The simplicity of handling deletions via tombstones in the RGA protocol comes with a high space overhead. More precisely, the *metadata overhead* [5] of a list implementation is the ratio between the size of a replica’s state (in bits) and the size of the user-observable content of the state, i.e., the list that will be read in this state. As we show, the metadata overhead of the RGA protocol is $O(D \lg k)$, where D is the number of deletions issued by clients and k is the total number of operations (Section 4). The number of deletions can be high. For example a 2009 study [32] indicates that the “George W. Bush” Wikipedia page has about 500 lines. However, since modifications are usually handled as deleting the original line and then inserting the revised line, the page had accumulated about 1.6 million deletions.¹

Our main result is that this overhead is, in some sense, inherent. We prove that any protocol from a certain class which implements the list specification for $n \geq 3$ replicas incurs a metadata overhead of $\Omega(D)$, where D is the number of deletions. This result holds even for the *weak* list specification and even if the network guarantees causal atomic broadcast [10]. The result holds for all *push-based* protocols, where each replica propagates list updates to its peers as soon as possible, and merges remote updates into its state as soon as they arrive (we give a precise definition in Section 5). This assumption captures the operation of all highly available protocols that we are aware of.

We establish our lower bound for the peer-to-peer model. However, using the fact that it holds for a network with causal atomic broadcast, we extend it to show that, in a push-based client/server list protocol, the metadata overhead at the clients is still $\Omega(D)$.

We prove our lower bound using an information-theoretic argument. For every $d \approx D/2$ -bit string w , we construct a particular execution α_w of the protocol such that, at its end, the user-observable state σ_w of some replica is a list of size $O(1)$ bits. We then show that, given σ_w , we can decode w by exercising the protocol in a black-box manner. This implies that all states σ_w must be distinct and, since there are 2^d of them, one of these states must take at least d bits. The procedure that decodes w from σ_w is nontrivial and represents the key insight of our proof. It recovers w one bit at a time using a “feedback loop” between two processes: one performs a black-box experiment on the protocol to recover the next bit of w , and the other reconstructs the corresponding steps of the execution α_w ; the messages sent in the reconstructed part of α_w then form the basis for the experiment to decode the next bit of w .

2. SYSTEM MODEL

We are concerned with highly available implementations of a replicated object [3, 5], which supports a set of operations Op . Such an implementation consists of *replicas* that receive and re-

spond to user operations on the object and use message passing to communicate changes to the object’s state. The *high availability* property sets this model apart from standard message-passing models: we require that replicas respond to user operations *immediately*—without performing any communication—so that user operations complete regardless of network latency and network partitions (e.g., device disconnection).

Replicas. We model a replica as a state machine $R = (Q, M, \Sigma, \sigma_0, E, \Delta)$, where Q is a set of *internal states*, M is a set of possible *messages*, $\Sigma = Q \times (M \cup \{\perp\})$ is a set of *replica states*, $\sigma_0 = (q_0, \perp) \in \Sigma$ is the *initial state*, E is a set of possible *events*, and $\Delta : \Sigma \times E \rightarrow \Sigma$ is a (partial) *transition function*. Note that a replica state explicitly includes a *send buffer*, containing the message pending transmission or \perp , which indicates that no message is pending. If $\Delta(\sigma, e)$ is defined, we say that event e is *enabled* in state σ . Transitions determined by Δ describe local steps of a replica in which it interacts with users and other replicas. These interactions are modeled by three kinds of events:

- $do(op, v)$: a user invokes an operation $op \in \text{Op}$ on the replicated object and immediately receives a response v from the replica;
- $send(m)$: the replica broadcasts a message $m \in M$; and
- $receive(m)$: the replica receives a message $m \in M$.

A *protocol* is a collection \mathcal{R} of replicas.

We require that a $send(m)$ event is enabled in state σ if and only if $\sigma = (q, m)$ for $m \neq \perp$, and in this case $\Delta((q, m), send(m)) = (q', \perp)$ for some q' . We further require that a replica can execute any operation with its return values computed deterministically: for any operation $op \in \text{Op}$, exactly one $do(op, v)$ event is enabled in σ . We also require that a replica can accept any message: for any message m , $receive(m)$ is enabled in σ . We assume that messages are unique and that a message’s sender is uniquely identifiable (e.g., messages are tagged with the sender id and a sequence number). We also assume that a replica broadcasts messages to all replicas, including itself²; replicas can implement point-to-point communication by ignoring messages for which they are not the intended recipient.

Executions. An *execution* of a protocol \mathcal{R} is a (possibly infinite) sequence of events occurring at the replicas in \mathcal{R} .³ For each event e , we let $\text{repl}(e) \in \mathcal{R}$ be the replica at which it occurs, and for each do event $e = do(op, v)$ we let $\text{op}(e) = op$ and $\text{rval}(e) = v$. A (finite or infinite) sequence of events e_1, e_2, \dots occurring at a replica $R = (Q, M, \Sigma, \sigma_0, E, \Delta)$ is *well-formed* if there is a sequence of states $\sigma_1, \sigma_2, \dots$ such that $\sigma_i = \Delta(\sigma_{i-1}, e_i)$ for all i . If the sequence is of length n , we refer to σ_n as the *state of R at the end of the sequence*.

We consider only *well-formed* executions, in which for every replica $R \in \mathcal{R}$: (1) the subsequence of events at R , denoted $\alpha|_R$, is well-formed; and (2) every $receive(m)$ event at R is preceded by a $send(m)$ event in α .

Let α be an execution. Event $e \in \alpha$ *happens before* event $e' \in \alpha$ [15] (written $e \xrightarrow{\text{hb}(\alpha)} e'$, or simply $e \xrightarrow{\text{hb}} e'$ if the context is clear) if one of the following conditions holds: (1) *Thread of execution*: $\text{repl}(e) = \text{repl}(e')$ and e precedes e' in α . (2) *Message delivery*: $e = send(m)$ and $e' = receive(m)$. (3) *Transitivity*: There is an event $f \in \alpha$ such that $e \xrightarrow{\text{hb}} f$ and $f \xrightarrow{\text{hb}} e'$.

²The latter is used to support atomic broadcast [10], defined later.

³Formally, an execution consists of events instrumented with unique event ids and replicas. In the paper we do not use this more accurate formulation so as to avoid clutter.

¹Wikipedia stores this information also to track the document’s edit history.

Network model. To ensure that every operation *eventually* propagates to all the replicas, we require that the network does not remain partitioned indefinitely. A replica R has a *message pending in event e of execution α* if R 's has a $\text{send}(m)$ event enabled in the state at the end of $\alpha'|_R$, where α' is the prefix of α ending with e .

DEFINITION 1. *The network is sufficiently connected in an infinite well-formed execution α of a protocol \mathcal{R} if the following conditions hold for all replicas $R \in \mathcal{R}$: (1) Eventual transmission: if R has a message pending infinitely often in α , then R also sends a message infinitely often in α , and (2) Eventual delivery: if R sends a message m , then every replica $R' \neq R$ eventually receives m .*

Collaborative editing protocols generally assume causal message delivery [24,27]. We model this by considering only executions that satisfy *causal broadcast* [6]:

DEFINITION 2. *An execution α of a protocol \mathcal{R} satisfies causal broadcast if for any messages m, m' , whenever $\text{send}(m) \xrightarrow{\text{hb}} \text{send}(m')$, any replica R can receive m' only after it receives m .*

In fact, our results hold even under a more powerful *atomic broadcast* [10] model, which delivers all messages to all replicas in the exact same order.

DEFINITION 3. *An execution α of protocol \mathcal{R} satisfies causal atomic broadcast if the following conditions hold: (1) Causal broadcast: α satisfies causal broadcast. (2) No duplicate delivery: each $\text{send}(m)$ event in α is followed by at most one $\text{receive}(m)$ event per replica $R' \in \mathcal{R}$. (3) Consistent order: if R receives m before m' , then any other replica R' receives m before m' .*

These broadcast primitives can be implemented when not provided by the network [6]; by providing them “for free,” we strengthen our lower bounds and ensure their independence from the complexity of implementing the broadcast primitive.

3. COLLABORATIVE TEXT EDITING

Following Ellis and Gibbs [13], we model the collaborative text editing problem (henceforth, simply collaborative editing) as the problem of implementing a highly available replicated list object whose elements are from some universe U . Users can insert elements, remove elements and read the list using the following operations, which form Op:

- $\text{ins}(a, k)$ for $a \in U$ and $k \in \mathbb{N}$: inserts a at position k in the list (starting from 0) and returns the updated list. For k exceeding the list size, we assume an insertion at the end. We assume that users pass identifiers a that are globally unique.
- $\text{del}(a)$ for $a \in U$: deletes the element a and returns the updated list. We assume that users pass only identifiers a that appear in the return value of the preceding operation on the same replica.
- read : returns the contents of the list.

The definition above restricts user behavior to simplify our technical development. Note that these restrictions are insignificant from a practical viewpoint, because they can be easily enforced: (1) identifiers can be made unique by attaching replica identifiers and sequence numbers; and (2) before each deletion, we can read the state of the list and skip the deletion if the deleted element does not appear in it.

3.1 Preliminaries: Replicated Data Types

We cannot specify the list object with a standard sequential specification, since replicas may observe only subsets of operations executed in the system, as a result of remote updates being delayed by the network. We address this difficulty by specifying the response of a list operation based on operations that are *visible* to it. Intuitively, these are the prior operations executed at the same replica and remote operations whose effects have propagated to the replica through the network. Formally, we use a variant of a framework by Burckhardt et al. [5] for specifying replicated data types [26]. We specify the list object by a set of *abstract executions*, which record the operations performed by users (represented by *do* events) and visibility relationships between them. Since collaborative editing systems generally preserve causality between operations [27], here we consider only *causal* abstract executions, where the visibility relation is transitive.

DEFINITION 4. *A causal abstract execution is a pair (H, vis) , where H is a sequence of *do* events⁴, and $\text{vis} \subseteq H \times H$ is an acyclic visibility relation (with $(e_1, e_2) \in \text{vis}$ denoted by $e_1 \xrightarrow{\text{vis}} e_2$) such that: (1) if e_1 precedes e_2 in H and $\text{repl}(e_1) = \text{repl}(e_2)$, then $e_1 \xrightarrow{\text{vis}} e_2$; (2) if $e_1 \xrightarrow{\text{vis}} e_2$, then e_1 precedes e_2 in H ; and (3) vis is transitive (if $e_1 \xrightarrow{\text{vis}} e_2$ and $e_2 \xrightarrow{\text{vis}} e_3$, then $e_1 \xrightarrow{\text{vis}} e_3$).*

Figure 1 graphically depicts abstract executions, where vis is the transitive closure of arrows in the figure and H is the result of some topological sort of vis . An abstract execution $A' = (H', \text{vis}')$ is a *prefix* of abstract execution A if: (1) H' is a prefix of H ; and (2) $\text{vis}' = \text{vis} \cap (H' \times H')$. A *specification* of an object is a prefix-closed set of abstract executions. A protocol correctly implements a specification when the outcomes of operations that it produces in any (concrete) execution can be justified by some abstract execution allowed by the specification.

DEFINITION 5. *An execution α of a protocol \mathcal{R} complies with an abstract execution $A = (H, \text{vis})$ if for every replica $R \in \mathcal{R}$, $H|_R = \alpha|_R^{\text{do}}$, where $\alpha|_R^{\text{do}}$ denotes the subsequence of *do* events by replica R in α .*

DEFINITION 6. *A protocol \mathcal{R} satisfies a specification S if every execution α of \mathcal{R} complies with some abstract execution $A \in S$.*

3.2 Specifying the List Object

We present two list specifications: strong and weak. Conceptually, the *strong* specification ensures that orderings relative to deleted elements hold even after the deletion, thereby disallowing the response ba for the read in Figure 1(b). The *weak* specification does not guarantee this property, allowing both ba and ab as responses.

We denote by $\text{elems}(A)$ the set of all elements inserted into the list in an abstract execution $A = (H, \text{vis})$:

$$\text{elems}(A) = \{a \mid \text{do}(\text{ins}(a, _), _) \in H\}.$$

Recall that we assume all inserted elements to be unique, and so there is a one-to-one correspondence between inserted elements and insert operations. For brevity, we write $e_1 \leq_{\text{vis}} e_2$ for $e_1 = e_2 \vee e_1 \xrightarrow{\text{vis}} e_2$.

DEFINITION 7. *An abstract execution $A = (H, \text{vis})$ belongs to the strong list specification $\mathcal{A}_{\text{strong}}$ if and only if there is a relation $\text{lo} \subseteq \text{elems}(A) \times \text{elems}(A)$, called the list order, such that:*

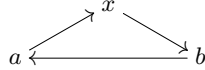
⁴Formally, H consists of *do* events instrumented with unique event ids and replicas, as in the case of an execution α . To avoid clutter, we do not use this more accurate presentation.

1. Each event $e = do(op, w) \in H$ returns a sequence of elements $w = a_0 \dots a_{n-1}$, where $a_i \in \text{elems}(A)$, such that
 - (a) w contains exactly the elements visible to e that have been inserted, but not deleted:

$$\forall a. a \in w \iff (do(\text{ins}(a, _), _) \leq_{\text{vis}} e) \wedge \neg(do(\text{del}(a, _), _) \leq_{\text{vis}} e).$$
 - (b) The order of the elements is consistent with the list order:

$$\forall i, j. (i < j) \implies (a_i, a_j) \in \text{lo}.$$
 - (c) Elements are inserted at the specified position: if $op = \text{ins}(a, k)$, then $a = a_{\min\{k, n-1\}}$.
2. The list order lo is transitive, irreflexive and total, and thus determines the order of all insert operations in the execution.

For example, the strong list specification is satisfied by the abstract execution in Figure 1(a) and the one in Figure 1(b) with the read returning ab ; this is justified by the list order $a \rightarrow x \rightarrow b$. On the other hand, the specification is not satisfied by the execution in Figure 1(b) with the read returning ba : for the outcomes of operations in this execution to be consistent with item 1 of Definition 7, the list order would have to be as shown above; but this order contains a cycle, contradicting item 2. In Section 4 we prove that the strong specification is implemented by an existing protocol, RGA [24]. However, some protocols, such as Jupiter [20], provide weaker guarantees and, in particular, allow the outcome ba in Figure 1(b). We therefore introduce the following weak list specification, to which our lower bound result applies (Section 6)⁵.



DEFINITION 8. An abstract execution $A = (H, \text{vis})$ belongs to the weak list specification $\mathcal{A}_{\text{weak}}$ if and only if there exists a relation $\text{lo} \subseteq \text{elems}(A) \times \text{elems}(A)$ such that:

1. Condition 1 in Definition 7 is satisfied.
2. lo is irreflexive and, for all events $e = do(op, w) \in H$, it is transitive and total on $\{a \mid a \in w\}$.

Unlike the strong specification, the weak one allows the list order lo to have cycles; the order is required to be acyclic only on the elements returned by some operation. In particular, the weak specification allows the execution in Figure 1(b) with the read returning ba , which is justified using the above cyclic list order. Since at the time of the read, x is deleted from the list, the specification permits us to decide how to order a and b without taking into account the orderings involving x : $a \rightarrow x$ and $x \rightarrow b$.

Eventual consistency. A desirable property of highly available replicated objects is *eventual consistency*. Informally, this guarantees that, if users stop issuing update requests, then the replicas will eventually converge to the same state [29, 30]. Our specifications imply a related *convergence* property: in an abstract execution satisfying $\mathcal{A}_{\text{strong}}$ or $\mathcal{A}_{\text{weak}}$, two read operations that see the same sets of list updates return the same response. This is because such operations will return the same elements (Definition 7, item 1a) and in the same order (Definition 7, item 1b). From the convergence property we can establish that our specifications imply eventual consistency for a class of protocols that guarantee the following property of *eventual visibility*.

DEFINITION 9. An abstract execution $A = (H, \text{vis})$ satisfies eventual visibility if for every event $e \in H$, there are only finitely many events $e' \in H$ such that $\neg(e \xrightarrow{\text{vis}} e')$.

⁵We conjecture that Jupiter satisfies the weak specification.

DEFINITION 10. A protocol \mathcal{R} satisfying the weak (resp., strong) list specification guarantees eventual visibility if every execution α of \mathcal{R} complies with some abstract execution $A \in \mathcal{A}_{\text{weak}}$ (respectively, $A \in \mathcal{A}_{\text{strong}}$) that satisfies eventual visibility.

Informally, eventual consistency holds for a protocol guaranteeing eventual visibility because: in an abstract execution with finitely many list updates, eventual visibility ensures that all but finitely many reads will see all the updates; then convergence ensures that they will return the same list. To guarantee eventual visibility, a protocol would rely on the network being sufficiently connected (Definition 1).

3.3 Metadata Overhead

In addition to the user-observable list contents, the replica state in a list protocol typically contains user-unobservable metadata that is used internally to provide correct behavior. The metadata overhead is the proportion of metadata relative to the user-observable list content.

Formally, let the *size* of an internal replica state q or a list $w \in U^*$ be the number of bits required to represent it in a standard encoding; we denote the size of x by $|x|$. The *metadata overhead* [5] of a state $\sigma = (q, m)$ is $|q|/|w|$ for the unique w such that $do(\text{read}, w)$ is enabled in σ ; here w represents the user-observable contents of σ . Note that the contents of the send buffer is not part of the metadata.

DEFINITION 11 ([5]). The worst-case metadata overhead of a protocol over a given subset of its executions is the largest metadata overhead of the state of any replica in any of these executions.

4. AN IMPLEMENTATION OF THE STRONG LIST SPECIFICATION

We now present an implementation of the list object, which is a reformulation of the RGA (Replicated Growable Array) protocol [24], and prove that it implements the strong list specification.

4.1 Timestamped Insertion Trees

Our representation of the list at a replica uses a *timestamped insertion (TI) tree* data structure. It stores both the list content and timestamp metadata used for deterministically resolving the order between elements concurrently inserted at the same position.

Formally, a tree is a finite set N of *nodes*, each corresponding to an element inserted into the list. A node is a tuple $n = (a, t, p)$, where $a \in U$ is the element, $t \in C$ is a *timestamp* for the insertion, and $p \in (C \cup \{\circ\})$ is either the parent node (identified by its timestamp) or the symbol \circ representing the tree root. We define the set C of timestamps and a total order on them later (Section 4.2). For a node $n = (a, t, p)$ we let $n.a = a$, $n.t = t$, and $n.p = p$. For two nodes n, n' with $n'.p = n.t$, we say n is the parent of n' and write $n \xrightarrow{\text{pa}} n'$.

DEFINITION 12. A set of nodes N is a TI tree if (1) timestamps uniquely identify nodes: $\forall n, n' \in N : n.t = n'.t \implies n = n'$; (2) all parents are present: if $n \in N$ and $n.p \neq \circ$, then $n' \xrightarrow{\text{pa}} n$ for some $n' \in N$; and (3) parents are older than their children: $n \xrightarrow{\text{pa}} n' \implies n.t < n'.t$.

Figure 2 shows an example of a TI tree and illustrates the read and insert operations explained below.

Read. To read the list, we traverse the tree N by depth-first search, starting at the root. We assemble the visited elements into a sequence $s(N)$ using prefix order (the parent precedes its children) and visit the children in decreasing timestamp order.

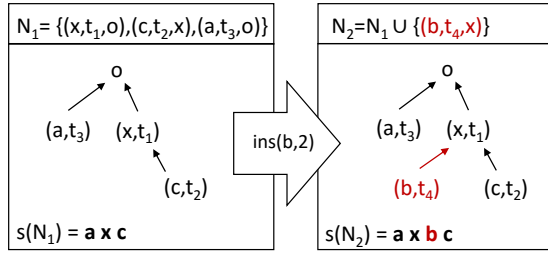


Figure 2: Illustration of TI trees. Each box shows a tree, with the set of nodes that define it, its graphical representation, and the sequence of elements it denotes. The tree on the right results from an insert operation for element b at position 2. The order on the timestamps is $t_1 < t_2 < t_3 < t_4$.

Insert. To insert a new element a at position k into the list, let $s(N) = a_0 \dots a_{n-1}$ and pick a new timestamp t that is larger than any of the timestamps appearing in N . Then let p be the element to the left of the insertion position: $p = a_{k-1}$ (if $k > 0$) or $p = \circ$ (if $k = 0$). We now add a new node (a, t, p) to N . Note that a newly inserted node is the child of the immediately preceding element with the highest timestamp. Thus, it is visited immediately after that element during a read, which makes it appear at the correct position in the list.

4.2 The RGA Protocol

We now define the RGA protocol $\mathcal{R}_{\text{rga}}^n$ for n replicas. Each replica stores a TI tree, as well as a set of elements that represent *tombstones*, used to handle deletions (Section 1). Insertions and deletions are recorded in a send buffer, which is periodically transmitted to other replicas by causal broadcast.

State and messages. Timestamps are pairs (x, i) , where $x \in N$ and $i \in \{1, \dots, n\}$ is a replica identifier. They are ordered lexicographically:

$$(x, i) < (x', i') \iff (x < x') \vee ((x = x') \wedge (i < i')).$$

Messages are of the form (A, K) , where A is a set of nodes (representing insert operations) and K is a set of elements (representing delete operations), and either A or K is non-empty. The state of a replica is $(N, T, (A, K))$, where: N is a TI tree, representing the replica-local view of the list; $T \subseteq U$ is the set of tombstones; and (A, K) is a send buffer, containing the message to send next. A pair (\emptyset, \emptyset) indicates that no message is pending (thus corresponding to \perp in Section 2). The initial state is $(\emptyset, \emptyset, (\emptyset, \emptyset))$.

do transitions. To execute an insert operation at a replica i in a state $(N, T, (A, K))$, we construct a node as described in the “Insert” procedure of Section 4.1 and add it to both N and A . As the timestamp of the node we take $((1 + (\text{the largest timestamp in } N)), i)$, or $(1, i)$ if $N = \emptyset$. This timestamp is guaranteed to be globally unique. To execute a delete operation, we add the deleted element to both T and K . All operations return the local view of the list, which is obtained by traversing N as described in the “Read” procedure of Section 4.1, and then removing all elements belonging to T .

send transition is enabled whenever either A or K is nonempty. It sends (A, K) as the message and sets both A and K to empty.

receive transition for a message (A_m, K_m) adds A_m to N and K_m to T . The protocol relies on causal delivery of messages, which ensures that no parents can be missing from N . In particular, N stays well-formed after adding A_m .

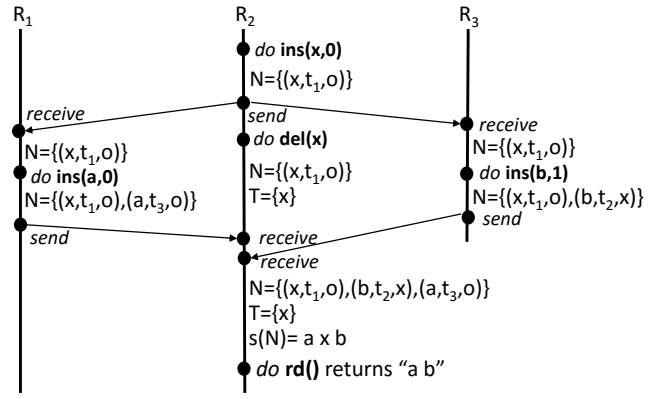


Figure 3: Illustration of an RGA execution, with time proceeding from top to bottom. Bullets show the transitions of the replicas R_1 , R_2 , and R_3 , and in some places we indicate the current state of the tree N and the tombstone set T . The order on the timestamps is $t_1 < t_2 < t_3$.

We show an example execution in Figure 3, which matches the example in Figure 1(b) and complies with the strong list specification.

4.3 Guarantees

The following theorems state the correctness and asymptotic complexity bounds of RGA. We provide full proofs in [2, §A] and discuss the key insights (convergence and stability) below.

THEOREM 1. *The protocol $\mathcal{R}_{\text{rga}}^n$ satisfies the strong list specification.*

THEOREM 2. *The worst-case metadata overhead of $\mathcal{R}_{\text{rga}}^n$ over executions with k operations and D deletions is $O(D \lg k)$.*

Convergence. Each replica maintains a TI tree that grows over time, meaning that nodes are added to the set, but never modified or removed. Because set union is associative and commutative, the order in which nodes are added does not matter. For example, changing the order of message delivery to R_2 in Figure 3 does not change the final tree. As a consequence, if the same set of nodes is delivered to any two replicas in any order, their trees are guaranteed to match, which ensures convergence.

Stability. The following lemma (proved in [2, §A]) shows that when we add more elements to a TI tree, the order of existing elements remains stable. This implies the strong list specification, because all replicas order all insertions the same way at all times.

LEMMA 3. *Let A, B be two TI trees such that $A \subseteq B$. Then $s(A)$ is a subsequence of $s(B)$.*

Trees vs. Lists. In the standard RGA implementation [25], TI trees are represented as lists (corresponding to the tree traversal). We show in [2, §A] that these representations are functionally equivalent. Lists are convenient to implement, but offer little insight as to why the algorithm guarantees convergence. Not surprisingly, the reason why RGA actually works has been a bit of a mystery, and we are not aware of any prior correctness proofs.

5. PUSH-BASED PROTOCOLS

Our lower bound results hold for *push-based* protocols, a class of protocols that contains the protocols of several collaborative editing

systems [20, 22, 24, 27], including the RGA protocol of Section 4. Informally, a replica in a push-based protocol propagates list updates to its peers as soon as possible and merges remote updates into its state as soon as they arrive (as opposed to using a more sophisticated mechanism, such as a consensus protocol). We define this class of protocols assuming that the network provides causal broadcast; when this is not the case, a protocol may need to delay merging arriving updates to enforce causality. Formally, we require that in a push-based protocol, every operation observe all operations that happen before it, that list insertions always generate a message, and that a deletion—which, unlike an insertion, may not be unique—generates a message if it does not already observe another deletion of the same element.

DEFINITION 13. *A protocol \mathcal{R} satisfying the weak (strong) list specification is push-based if the following hold:*

- For any execution α of \mathcal{R} and $e = do(\text{ins}(a, _), _) \in \alpha$, replica $\text{repl}(e)$ has a message pending after e .
- For any execution α of \mathcal{R} and $e = do(\text{del}(a, _), _) \in \alpha$, if there does not exist event $e' = do(\text{del}(a, _), _) \in \alpha$ that happens before e , then replica $\text{repl}(e)$ has a message pending after e .
- For every execution α of \mathcal{R} there exists an abstract execution $A = (H, \text{vis}) \in \mathcal{A}_{\text{weak}}$ ($A \in \mathcal{A}_{\text{strong}}$) that α complies with, such that $\forall e', e \in H. e' \xrightarrow{\text{vis}} e \iff e' \xrightarrow{\text{hb}} e$.

The class of push-based protocols contains both *op-based* protocols [5], in which a message carries a description of the latest operations that the sender has performed (e.g., RGA), and *state-based* [5] protocols, in which a message describes all operations the sender knows about (i.e., its state). We also show ([2, §B]) that the class of push-based protocols contains eventually consistent write-propagating protocols [3]—which model many deployed highly available eventually consistent protocols [4, 5, 9, 11, 12, 18, 26, 33]—under the natural assumption that sending a message does not affect the state of the list at the sending replica.

6. LOWER BOUNDS ON METADATA OVERHEAD

Here we show a lower bound on the worst-case metadata overhead (Definition 11) of a push-based protocol satisfying the weak or strong list specification.

THEOREM 4. *Let \mathcal{R} be a push-based protocol that satisfies the weak or strong list specification for $n \geq 3$ replicas. Then the worst-case metadata overhead of \mathcal{R} over executions with D deletions is $\Omega(D)$.*

This follows from the following theorem, because any execution consistent with the strong list specification is also consistent with the weak one.

THEOREM 5. *Let \mathcal{R} be a push-based protocol that satisfies the weak list specification for $n \geq 3$ replicas. Then for every integer $D \geq 4$, there exists an execution α_D of \mathcal{R} with D deletions such that: (1) the metadata overhead of some state σ of some replica R in α_D is $\Omega(D)$; (2) α_D satisfies causal atomic broadcast; and (3) R does not receive any message before σ in α_D .*

PROOF. Let $d = \lfloor (D - 2)/2 \rfloor$. We show that there exists an execution of \mathcal{R} with D deletions that satisfies the desired conditions, in which the user-observable contents of some internal state is a list

with a single element, and yet the size of this state is at least d bits. It follows that the metadata overhead of this state is $\Omega(D)$.

We show the existence of this execution using an information-theoretic argument. Namely, for every d -bit string w we construct an execution α_w that satisfies causal atomic broadcast and in which: (1) replica R_1 performs D deletions and receives no messages; (2) at the end of α_w , the user-observable list at R_1 contains the single element “*” and R_1 has no messages pending; and yet (3) we can decode w given only σ_w , the state of R_1 at the end of α_w (this decoding process exercises the protocol \mathcal{R} in a black-box manner). Hence, all states σ_w must be distinct. Since there are 2^d of them, one of these states σ_{w_0} must take at least d bits. Since this state has no messages pending, its metadata overhead is $\Omega(D)$, and thus, α_{w_0} is the desired execution.

Encoding w . Given a d -bit string $w = w_1 \dots w_d$, we construct an execution α_w of \mathcal{R} that builds a list encoding the path from the root of a binary tree of height d to the w -th leaf (when w is interpreted as the binary representation of an integer). Figure 4(a) details the construction: it shows pseudocode which, as it executes, constructs the execution; instructions of the form e_i correspond to a state transition e at replica R_i . We abuse notation by writing *op* instead of *do(op, _)*, by specifying inserts of whole strings instead of element by element, and by specifying positions relative to prior insertions rather than with integers. Figure 5(a) depicts α_w for $w = 10$.

Only replica R_1 participates in the encoding execution α_w . We start by inserting the string $[0]_0$ (i.e., the root). Because \mathcal{R} is a push-based protocol, R_1 has a message m_1 pending following these insertions. We then proceed with a series of steps, for $i = 1, \dots, d$. Each step i begins with R_1 in state σ_i having a message m_i pending. R_1 first broadcasts m_i . We then insert the string $[i]_i$ immediately to the left or to the right of $[i-1]_{i-1}$, depending on whether the i -th bit of w is set. Because \mathcal{R} is a push-based protocol, R_1 has a message pending following these insertions, and we proceed to step $i+1$. When we are done, we broadcast the current pending message and insert the element * between $[_d]$ and $]_d$, and broadcast the message m_{d+2} that is pending following this insertion. For example, if $w = 10$, the state of the list at R_1 at this point is $[0]_0 [2*]_2 [1]_1$. We then delete all the $[i]$ and $]_i$ elements, for $i = 0, \dots, d$, and if D is odd, we insert and delete an additional element, so that the number of deletions in α_w is exactly D . Because \mathcal{R} is a push-based protocol, R_1 has a message pending following these deletions, which we broadcast to empty R_1 's send buffer. Finally, we read the list at R_1 , observing that it is *. This follows because for any abstract execution $A = (H, \text{vis})$ that the encoding execution α_w complies with, all ins and del events are visible to the read, due to Condition (1) of Definition 4. The read's response must thus be *, since by assumption one of such executions A is consistent with the weak list specification.

The output of the encoding procedure is σ_w , the state of R_1 at the end of the encoding execution α_w . It is easy to check that α_w is well-formed; furthermore, it vacuously satisfies causal atomic broadcast.

Decoding w from σ_w . We reconstruct w one bit at a time by “replaying” the execution α_w . To replay iteration i of α_w , we rely on a procedure `Recover()` that recovers w_i from σ_w and m_1, \dots, m_i . (We describe `Recover()` in the next paragraph; for now, assume it is an oracle.) Knowing w_i , in turn, determines the next event of R_1 in α_w , and hence provides us with m_{i+1} . The decoding process thus only uses messages from R_1 that it reconstructs with the bits of w already known. Figure 4(b) shows the pseudocode which, as it executes, decodes w . We start with R_1 in its initial state and recon-

Input: $w = w_1, \dots, w_d$
// R_1 starts in its initial state σ_0 .
 $\text{ins}_1([0]_0, 0)$
for $i = 1, \dots, d$
// The state of R_1 here is σ_i .
 $\text{send}_1(m_i)$
if $w_i = 1$ **then**
 $\text{ins}_1([i]_i \text{ just after } [i-1])$
else
 $\text{ins}_1([i]_i \text{ just before } [i-1])$
// The state of R_1 here is σ_{i+1} .
 $\text{send}_1(m_{d+1})$
 $\text{ins}_1(* \text{ between } [d \text{ and } d])$
 $\text{send}_1(m_{d+2})$
for $i = 0, \dots, d$
 $\text{del}_1([i])$
 $\text{del}_1([i])$
if $D = 2(d+1) + 1$ **then**
 $\text{ins}_1(b, 0)$
 $\text{del}_1(b)$
 $\text{send}_1(m_{d+3})$
 $\text{read}_1 = *$
// The state of R_1 here is σ_w .

(a) Execution α_w , at the end of which the state of R_1 encodes w

Input: σ_w
// R_1 starts in its initial state σ_0 ,
// which does not depend on w .
 $\text{ins}_1([0]_0, 0)$
for $i = 1, \dots, d$
// We now know σ_i ,
// so can generate m_i .
 $\text{send}_1(m_i)$
 $w_i \leftarrow \text{Recover}(\sigma_w, m_1, \dots, m_i)$
if $w_i = 1$ **then**
 $\text{ins}_1([i]_i \text{ just after } [i-1])$
else
 $\text{ins}_1([i]_i \text{ just before } [i-1])$
// We now know σ_{i+1} .
output w

(b) Decoding w given σ_w

Input: $\sigma_w, m_1, \dots, m_i$
// We perform the state transitions
// below on copies of the state
// machines, with R_2 starting in its
// initial state and R_1 in state σ_w .
// The replica state in the outer
// decoding procedure is unaffected.
for $j = 1, \dots, i$
 $\text{receive}_2(m_j)$
 $\text{receive}_1(m_j)$
 $\text{read}_2 = \dots [i-1]_{i-1} \dots$
 $\text{ins}_2(x \text{ between } [i-1] \text{ and } [i-1])$
 $\text{send}_2(m_x)$
 $\text{receive}_2(m_x)$
 $\text{receive}_1(m_x)$
if $\text{read}_1 = x *$
return 1
else // $\text{read}_1 = * x$
return 0

(c) $\text{Recover}(\sigma_w, m_1, \dots, m_i)$

Figure 4: Encoding and decoding procedures. Events are subscripted with the id of their replica.

struct m_1 , which does not depend on w . We then proceed in steps, for $i = 1, \dots, d$. In step i we know m_1, \dots, m_i , and we recover bit w_i from σ_w and m_1, \dots, m_i . Having recovered w_i , we replay the insertion that R_1 performs at step i of the encoding and reconstruct m_{i+1} .

Recovering w_i from σ_w and m_1, \dots, m_i . The $\text{Recover}()$ procedure determines w_i by performing state transitions on fresh copies of R_1 and R_2 ; the transitions that an execution of $\text{Recover}()$ performs have no effect on the state of the replicas in the “replayed” execution constructed by the decoding process, or on other $\text{Recover}()$ executions. Figure 4(c) shows these state transitions, and Figures 5(b)–5(c) illustrate the overall decoding of $w = 10$ (the use of the replica R_3 is explained below). We start off with R_2 in its initial state and R_1 in state σ_w . We deliver the messages m_1, \dots, m_i to both replicas in the same order. We then read at R_2 and receive response v_w^i ; we will show that $[i-1]_{i-1} \in v_w^i$. Next, R_2 inserts element x between $[i-1]$ and $[i-1]$ and broadcasts a message m_x , which we deliver to both replicas. Finally, we read at R_1 and observe the list in state y_w^i . We will show that y_w^i contains only x and $*$, and if x precedes $*$ then $w_i = 1$; otherwise, $w_i = 0$.

Validity of $\text{Recover}(\sigma_w, m_1, \dots, m_i)$ state transitions. Assuming that m_1, \dots, m_i are the first i messages sent by R_1 in α_w , we show that the state transitions performed by an execution of $\text{Recover}(\sigma_w, m_1, \dots, m_i)$ in Figure 4(c) occur in an extension β_w^i of α_w of the form:

$$\beta_w^i = \alpha_w$$

$$\text{receive}_2(m_1) \text{receive}_1(m_1) \dots \text{receive}_2(m_i) \text{receive}_1(m_i)$$

$$\text{do}_2(\text{read}_2, v_w^i) \text{do}_2(\text{ins}_2(x, k_w^i, _)) \text{send}_2(m_x)$$

$$\text{receive}_2(m_x) \text{receive}_1(m_x) \text{do}_1(\text{read}_1, y_w^i),$$

where k_w^i is the position at which R_2 inserts x into the list. In the following, we prove that the execution β_w^i is well-formed (Claim 6), that it satisfies causal atomic broadcast (Claim 7), that $[i-1]_{i-1} \in v_w^i$ (Claim 8), and that $y_w^i = x *$ or $y_w^i = * x$

(Claim 9). To this end, we exploit the fact that σ_w is R_1 ’s state at the end of α_w , which allows $\text{Recover}()$ to perform the same state transitions at R_1 that occur in β_w^i , without having access to the entire execution α_w that leads R_1 to state σ_w .

CLAIM 6. *Execution β_w^i is well-formed.*

PROOF. By assumption, $\text{Recover}()$ is passed the first i messages sent by R_1 in α_w . The claim thus follows from the following: (1) α_w is well-formed; (2) at the end of α_w , R_1 is in state σ_w ; (3) because R_2 does not participate in α_w , the state of R_2 at the end of α_w is its initial state; (4) a replica always accepts any sent message (by definition); and (5) \mathcal{R} is push-based, and so R_2 has a message pending following its ins operation. \square

CLAIM 7. *Execution β_w^i satisfies causal atomic broadcast.*

PROOF. Immediate from inspection of the message delivery order in β_w^i . \square

CLAIM 8. $[i-1]_{i-1} \in v_w^i$.

PROOF. For $j = 1, \dots, d + 1$, let $e_j, f_j \in \alpha_w$ be the *do* events in which R_1 inserts $[j-1]$ and $[j-1]$ into the list. Let $r \in \beta_w^i$ be the *do* event at which R_2 reads v_w^i . Because \mathcal{R} is a correct push-based protocol and β_w^i satisfies causal atomic broadcast, by Definition 13, β_w^i complies with some abstract execution $A = (H, \text{vis}) \in \mathcal{A}_{\text{weak}}$ such that $e_j \xrightarrow{\text{vis}} r$ and $f_j \xrightarrow{\text{vis}} r$ if and only if $j \leq i$, and no del operation is visible to r . This holds because in β_w^i , R_2 receives only the messages m_1, \dots, m_i before r , and each m_j is the first message sent by R_1 after e_j and f_j . It thus follows from the definition of the weak list specification (Definition 8) that $v_w^i = \dots [i-1]_{i-1} \dots$. \square

CLAIM 9. $y_w^i = x *$ or $y_w^i = * x$.

PROOF. Let $f \in \beta_w^i$ be the *do* event at which R_2 inserts x into the list, and $r \in \beta_w^i$ be the *do* event at which R_1 reads y_w^i . Because

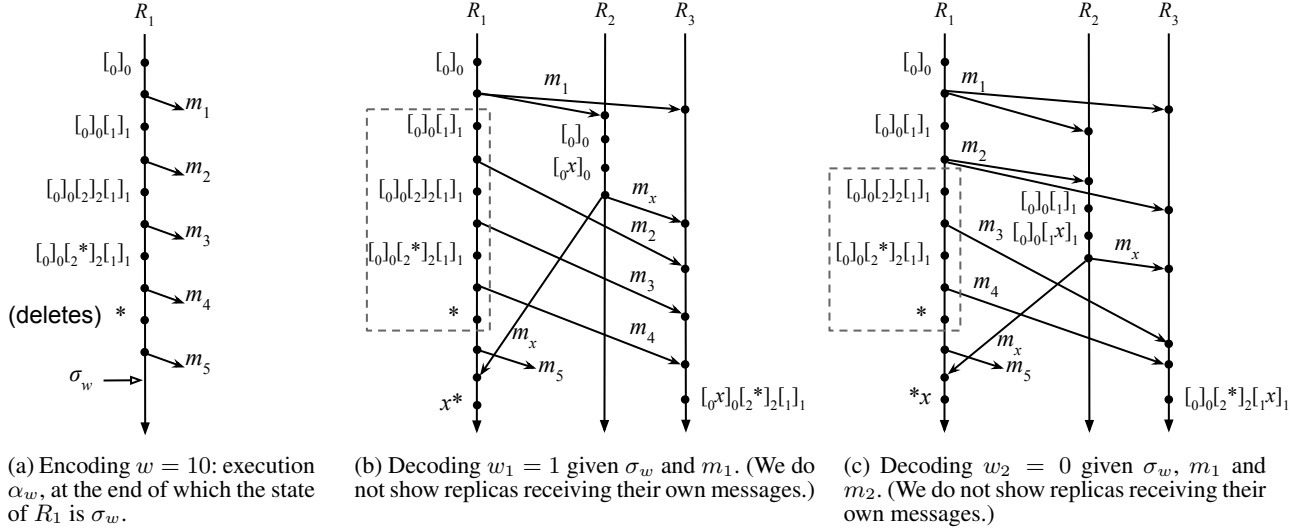


Figure 5: Examples of the encoding and decoding procedures from Theorem 5 applied to $w = 10$. Figure 5(a) shows the execution α_w constructed by the encoding procedure, whose output is σ_w , the state of R_1 at the end of α_w . Figure 5(b) shows the first step of decoding w from σ_w . The decoding procedure performs the state transitions of the events at R_1 that are outside of the dashed rectangle and of the events at R_2 ; these transitions are valid because they occur in the depicted execution, in which the events inside the dashed rectangle lead R_1 to state σ_w . The relative order of x and $*$ read at R_1 therefore recovers bit $w_1 = 1$. Figure 5(c) shows the second step of decoding w : Recovering $w_1 = 1$ allows the decoding procedure to perform the state transitions by R_1 in α_w that depend on w_1 .

\mathcal{R} is a correct push-based protocol, β_w^i compiles with some abstract execution $A = (H, \text{vis}) \in \mathcal{A}_{\text{weak}}$ such that $f \xrightarrow{\text{vis}} r$. Now, let $e \in \beta_w^i$ be the *do* event at which R_1 inserts $*$ into the list. Then $e \xrightarrow{\text{vis}} r$ by definition of an abstract execution (Definition 4). Because all other elements inserted in β_w^i are deleted by R_1 before r , but x and $*$ are not deleted in β_w^i , the claim follows. \square

Correctness of recovering w_i . Having shown that the state transitions performed by Recover() yield the lists $y_w^i = x*$ or $y_w^i = *x$, it remains to show that we correctly recover w_i from y_w^i : ($y_w^i = x*$) \iff ($w_i = 1$).

In principle, the weak list specification allows R_1 's read to order x and $*$ arbitrarily, since $]_{i-1}$ and $]_{i-1}$ are deleted from the list by the time the read occurs. We show, however, that R_1 cannot do this, because it cannot rule out the possibility that another replica has already observed $]_{i-1} x]_{i-1}$ and $*$ together, and therefore their order is fixed. Consider the following extension of β_w^i , in which R_3 receives the messages generated after each insertion and then reads the list (it is easy to see that this execution satisfies causal atomic broadcast):

$$\gamma_w^i = \beta_w^i \text{ receive}_3(m_1) \dots \text{receive}_3(m_i) \text{ receive}_3(m_x) \\ \text{receive}_3(m_{i+1}) \dots \text{receive}_3(m_{d+2}) \text{ do}_3(\text{read}_3, z_w^i).$$

We show that the list z_w^i contains $*$ after x if and only if $w_i = 1$. Informally, this follows because every element inserted from iteration i onwards in the encoding procedure (and hence in γ_w^i), including $*$, goes after $]_{i-1}$ if and only if $w_i = 1$, and no del events are visible to R_3 , so its read response must order x before $]_{i-1}$ before $*$.

Formally, consider the following events in γ_w^i : w_* , the ins of $*$ by R_1 ; w_x , the ins of x by R_2 ; and r_z , the read by R_3 , whose response is z_w^i . Because \mathcal{R} is a correct push-based protocol and γ_w^i satisfies causal broadcast, by Definition 13, γ_w^i complies with some abstract execution $A = (H, \text{vis}) \in \mathcal{A}_{\text{weak}}$ such that for any $e, e' \in H$, $e' \xrightarrow{\text{vis}} e$ if and only if $e' \xrightarrow{\text{hb}} e$. Therefore, no del event

is visible to w_* , w_x or r_z . Let lo be a list order that A is consistent with (Definition 8). We proceed to show that $(x, *) \in \text{lo}$ if and only if $w_i = 1$. Observe that if $w_i = 1$, every element inserted from iteration i onwards of the encoding process is inserted after $]_{i-1}$, and if $w_i = 0$, every element inserted from iteration i onwards is inserted before $]_{i-1}$. Therefore, the response of w_* establishes that $(]_{i-1}, *) \in \text{lo}$ if and only if $w_i = 1$. The response of w_x establishes that $(]_{i-1}, x) \in \text{lo}$ and $(x,]_{i-1}) \in \text{lo}$. It follows that $(x, *) \in \text{lo}$ if and only if $w_i = 1$, since $]_{i-1}, x,]_{i-1}, * \in z_w^i$ and lo is total and transitive on $\{a \mid a \in \text{rval}(r_z) = z_w^i\}$.

We conclude by noting that $y_w^i = x*$ or $y_w^i = *x$ (Claim 9); recall that y_w^i is the response to the read at R_1 performed by Recover(). Since $(x, *) \in \text{lo}$ if and only if $w_i = 1$, then $y_w^i = x*$ if and only if $w_i = 1$.

6.1 Extension to a Client/Server Model

In a *client/server protocol*, replicas communicate only with a central server and not directly with each other. (The motivation is to maintain state on the server instead of on the replicas, and so the server usually does more than merely relay messages between replicas [20].) To model such protocols in our framework, which assumes a broadcast transport, we require replicas to process only messages to/from the server:

DEFINITION 14. A protocol $\mathcal{R} = \{R_1, \dots, R_n, S\}$ is a client/server protocol if for every replica $R_i = (Q^i, M, \Sigma^i, \sigma_0^i, E, \Delta^i)$, and $\sigma \in \Sigma^i$, if $\Delta^i(\sigma, \text{receive}(m)) \neq \sigma$, then m was sent by S . We call S the server.

In practice, users do not interact directly with the server, and so we consider only executions in which *do* events do not occur at the server.

Assuming atomic broadcast, a broadcast protocol can *simulate* a client/server protocol using state machine replication [15].

PROPOSITION 10. Let $\mathcal{R} = \{R_1, \dots, R_n, S\}$ be a client/server protocol. Then there exists a protocol $\mathcal{R}' = \{R'_1, \dots, R'_n\}$

State	Event e	New state
$((r, s), m)$	$do(op, v)$	$((r', s), m')$, where $(r', m') = \Delta^i((r, m), e)$
$((r, s), m)$	$send(m)$	$((r', s), \perp)$, where $(r', \perp) = \Delta^i((r, m), e)$
$((r, s), m)$	$receive(m)$	$((r', s'), m')$, where if $\Delta^S((s, \perp), e) = (s^*, \perp)$, then $s' = s^*$ and $(r', m') = (r, m) =$ $\Delta^i((r, m), receive(m))$; and if $\Delta^S((s, \perp), e) = (s^*, m^*)$, then $(s', \perp) = \Delta^S((s^*, m^*), send(m^*))$, $(r', m') = \Delta^i((\hat{r}, \hat{m}), receive(m^*))$, where $(\hat{r}, \hat{m}) = (r, m) =$ $\Delta^i((r, m), receive(m))$

Figure 6: State machine of replica $R'_i \in \mathcal{R}'$ simulating replica $R_i = (Q^i, M, \Sigma^i, \sigma_0^i, E, \Sigma^i) \in \mathcal{R}$. The initial internal state is (q_0^i, q_0^S) , where q_0^i and q_0^S are the initial internal states of R_i and of S , respectively.

that simulates \mathcal{R} in the following sense: (1) for any execution α' of \mathcal{R}' that satisfies causal atomic broadcast, there exists an execution α of \mathcal{R} such that $\alpha|_{R_i}^{do} = \alpha'|_{R_i}^{do}$ for $i = 1..n$; (2) the set of internal states of each R'_i is $Q^i \times Q^S$, where Q^i and Q^S are respective sets of R_i and S ; and (3) until R'_i receives a message, its state is (\perp, q_0^S) , where (q_0^S, \perp) is the initial state of S .

PROOF. For $1 = 1..n$, replica $R'_i \in \mathcal{R}'$ maintains two state machines, of R_i and of S . R'_i broadcasts exactly the messages broadcast by the replica R_i it is simulating. We use the fact that messages are delivered to all replicas in \mathcal{R}' in the same order to simulate the server S using state machine replication.

Figure 6 shows the state machine of replica $R'_i \in \mathcal{R}'$. Upon a *do* event, R'_i performs the corresponding transition on R_i 's state machine and broadcasts any message m' that R_i would send to S . Upon receiving a message m , R'_i delivers m to the two state machines it maintains. (However, because m corresponds to a message sent by some R_j , the R_i state machine ignores it, by Definition 14.) If, as a result of receiving m , S broadcasts a message m^* , then R'_i (locally) delivers m^* to R_i 's state machine and broadcasts any message m' that R_i sends as a result of receiving m^* .

In any execution α' of \mathcal{R}' that satisfies causal atomic broadcast, all messages are delivered to all replicas in the same order. Therefore, each replica R'_i performs the same state transitions at S , and (locally) delivers the same messages from S to its R_i state machine. The claim follows. \square

Client/server lower bound. Since the executions constructed in the proof of Theorem 5 satisfy causal atomic broadcast, they can also be viewed as executions of a protocol simulating a push-based client/server protocol (Proposition 10). We therefore obtain

COROLLARY 11. *Let \mathcal{R} be a push-based client/server protocol that satisfies the weak or strong list specification for $n \geq 3$ replicas. Then the worst-case metadata overhead of \mathcal{R} on the clients over executions with D deletions is $\Omega(D)$.*

PROOF. Let \mathcal{R} be a push-based client/server protocol that satisfies the weak or strong list specification. Let \mathcal{R}' be the protocol simulating \mathcal{R} from Proposition 10. Take $D \geq 4$. By Theorem 5, there exists an execution α_D of \mathcal{R}' with D deletions such that: (1) the metadata overhead of some state σ of some replica $R' \in \mathcal{R}'$ is $\Omega(D)$; (2) α_D satisfies causal atomic broadcast; and

(3) R' does not receive any message before σ . Because \mathcal{R}' simulates \mathcal{R} , we have that $\sigma = ((q^R, q^S), \perp)$, where q^R and q^S are, respectively, internal states of the replica $R \in \mathcal{R}$ that R' is simulating and of the server. Moreover, it follows from Proposition 10 that q^S is the initial internal state of the server. Therefore, $|(q^R, q^S)| = O(|q^R| + |q^S|) = O(|q^R|)$, because $|q^S|$ is a constant. Since the user-observable content at R' and R is the same, it follows that the metadata overhead at R is $\Omega(D)$. \square

7. RELATED WORK

Previous attempts at specifying the behavior of replicated list objects [17, 28] have been informal and imprecise: they typically required the execution of an operation at a remote replica to *preserve the effect* of the operation at its original replica, but they have not formally defined the notions of the effect and its preservation.

Burckhardt et al. [5] have previously proposed a framework for specifying replicated data types (on which we base our list specifications) and proved lower bounds on the metadata overhead of several data types. In contrast to us, they handle much simpler data types than a list. Thus, our specifications have to extend theirs with an additional relation, defining the order of elements in the list. Similarly, their proof strategy (and its extension in [3]) for establishing lower bounds would not be applicable to lists; obtaining a lower bound in this case requires a more delicate decoding argument, recovering information incrementally.

There are more protocols implementing a highly available replicated list than the RGA protocol we considered. Treedoc [22] and Logoot [32] are other implementations of the strong list specification using the approach of replicated data types [26]. As in RGA, the state of a replica can be viewed as a tree, where a deterministic traversal defines the order of the list. The replication protocol represents position of a node in the tree as a sequence of edge labels on the path from the root of the tree (in RGA, it is a relative position to an existing node). Like RGA, these protocols have worst-case metadata overhead linear in the number of deletions. WOOT [21] is a graph-based list implementation: its main component is a representation of a partial list order, i.e., ordering restrictions inferred at the time user performs operations. The total list order is computed as a view of the graph based on a non-declarative specification of intended ordering.

Another class of protocols is based on operational transformations (OT) [13], which apply certain transformation functions to pairs of concurrent updates. If applying the transformation function allows commuting two operations (TP1) and three operations (TP2) then OT ensures that the list state converges, regardless of the order in which the operations are received [23]. However, it was shown [14] that several OT protocols do not satisfy TP1 and TP2 and do not converge. OT protocols store a log of updates at each replica, so their metadata overhead is also at least linear in the number of updates.

8. CONCLUSION

This paper provides a precise specification of the list replicated object, which models the core functionality of collaborative text editing systems. We define a *strong* list—and show that it is implemented by an existing system [24]—as well as a *weak* list, which we conjecture describes the behavior of the Jupiter protocol [20], underlying public collaboration systems [31].

We prove a lower bound of $\Omega(D)$, where D is the number of deletions, on the metadata overhead of push-based list protocols, which model the implementation of all highly available list protocols that we are aware of. Our lower bound applies for both weak

and strong semantics. Exploring client/server systems in future research is therefore of practical interest, as some client/server systems do not implement the strong semantics, and our results suggest that this might not offer a complexity advantage.

We also show a simple list protocol whose metadata overhead is $O(D \lg k)$, where k is the number of operations. Closing the gap between the upper and lower bound is left for future work, as is the question of relaxing the restriction to push-based protocols.

Our work is a first step towards specifying and analyzing general collaborative editing systems, providing more features than those captured by the list object. This includes systems for sharing structured documents, such as XML [19]. While our lower bound would hold for more general systems, it is possible that the additional features induce additional complexity.

Acknowledgments

We would like to thank Marc Shapiro and Pascal Urso for comments that helped improve the paper. Attiya and Morrison were supported by the Israel Science Foundation (grant 1749/14) and by Yad Hanadiv Foundation. Gotsman was supported by an EU project ADVENT. Yang was supported by an Institute for Information & Communication Technology Promotion (IITP) grant funded by the Korea government (MSIP, No. R0190-15-2011). Zawirski was supported by an EU project SyncFree.

9. REFERENCES

- [1] Apache Wave. <https://incubator.apache.org/wave/>.
- [2] H. Attiya, S. Burckhardt, A. Gotsman, A. Morrison, H. Yang, and M. Zawirski. Specification and Complexity of Collaborative Text Editing (Extended Version). Available from <http://www.software.imdea.org/~gotsman>.
- [3] H. Attiya, F. Ellen, and A. Morrison. Limitations of Highly-Available Eventually-Consistent Data Stores. In *PODC*, 2015.
- [4] P. Bailis, A. Ghodsi, J. M. Hellerstein, and I. Stoica. Bolt-on Causal Consistency. In *SIGMOD*, 2013.
- [5] S. Burckhardt, A. Gotsman, H. Yang, and M. Zawirski. Replicated Data Types: Specification, Verification, Optimality. In *POPL*, 2014.
- [6] C. Cachin, R. Guerraoui, and L. Rodrigues. *Introduction to Reliable and Secure Distributed Programming*. Springer Publishing Company, Incorporated, 2nd edition, 2011.
- [7] J. Day-Richter. What's different about the new Google Docs: Conflict resolution. http://googledrive.blogspot.com/2010/09/whats-different-about-new-google-docs_22.html, 2010.
- [8] J. Day-Richter. What's different about the new Google Docs: Making collaboration fast. <http://googledrive.blogspot.com/2010/09/whats-different-about-new-google-docs.html>, 2010.
- [9] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Vosshall, and W. Vogels. Dynamo: Amazon's Highly Available Key-value Store. In *SOSP*, 2007.
- [10] X. Défago, A. Schiper, and P. Urbán. Total Order Broadcast and Multicast Algorithms: Taxonomy And Survey. *ACM CSUR*, 36(4), 2004.
- [11] J. Du, S. Elnikety, A. Roy, and W. Zwaenepoel. Orbe: Scalable Causal Consistency Using Dependency Matrices and Physical Clocks. In *SoCC*, 2013.
- [12] J. Du, C. Iorgulescu, A. Roy, and W. Zwaenepoel. Closing The Performance Gap between Causal Consistency and Eventual Consistency. In *PaPEC*, 2014.
- [13] C. A. Ellis and S. J. Gibbs. Concurrency Control in Groupware Systems. In *SIGMOD*, 1989.
- [14] A. Imine, M. Rusinowitch, G. Oster, and P. Molli. Formal Design and Verification of Operational Transformation Algorithms for Copies Convergence. *TCS*, 351(2), 2006.
- [15] L. Lamport. Time, Clocks, and the Ordering of Events in a Distributed System. *CACM*, 21(7), 1978.
- [16] B. Leuf and W. Cunningham. *The Wiki Way: Quick Collaboration on the Web*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2001.
- [17] D. Li and R. Li. Preserving Operation Effects Relation in Group Editors. In *CSCW*, 2004.
- [18] W. Lloyd, M. J. Freedman, M. Kaminsky, and D. G. Andersen. Don't Settle for Eventual Consistency. *CACM*, 57(5), 2014.
- [19] S. Martin, P. Urso, and S. Weiss. Scalable XML collaborative editing with undo. In *OTM*, 2010.
- [20] D. A. Nichols, P. Curtis, M. Dixon, and J. Lamping. High-latency, low-bandwidth windowing in the Jupiter collaboration system. In *UIST*, 1995.
- [21] G. Oster, P. Urso, P. Molli, and A. Imine. Data consistency for P2P collaborative editing. In *CSCW*, 2006.
- [22] N. Preguiça, J. M. Marqués, M. Shapiro, and M. Letia. A commutative replicated data type for cooperative editing. In *ICDCS*, 2009.
- [23] M. Ressel, D. Nitsche-Ruhland, and R. Gunzenhäuser. An integrating, transformation-oriented approach to concurrency control and undo in group editors. In *CSCW*, 1996.
- [24] H.-G. Roh, M. Jeon, J.-S. Kim, and J. Lee. Replicated abstract data types: Building blocks for collaborative applications. *JPDC*, 71(3), 2011.
- [25] M. Shapiro, N. Preguiça, C. Baquero, and M. Zawirski. A Comprehensive study of Convergent and Commutative Replicated Data Types. Technical Report RR-7506, INRIA, Jan. 2011.
- [26] M. Shapiro, N. Preguiça, C. Baquero, and M. Zawirski. Conflict-free Replicated Data Types. In *SSS*, 2011.
- [27] C. Sun and C. Ellis. Operational transformation in real-time group editors: issues, algorithms, and achievements. In *CSCW*, 1998.
- [28] C. Sun, X. Jia, Y. Zhang, Y. Yang, and D. Chen. Achieving convergence, causality preservation, and intention preservation in real-time cooperative editing systems. *ACM TOCHI*, 5(1), 1998.
- [29] D. B. Terry, M. M. Theimer, K. Petersen, A. J. Demers, M. J. Spreitzer, and C. H. Hauser. Managing Update Conflicts in Bayou, a Weakly Connected Replicated Storage System. In *SOSP*, 1995.
- [30] W. Vogels. Eventually Consistent. *CACM*, 52(1), 2009.
- [31] D. Wang, A. Mah, and S. Lassen. Google Wave Operational Transformation. <https://wave-protocol.googlecode.com/hg/whitepapers/operational-transform/operational-transform.html>, 2010.
- [32] S. Weiss, P. Urso, and P. Molli. Logoot: A Scalable Optimistic Replication Algorithm for Collaborative Editing on P2P Networks. In *ICDCS*, 2009.
- [33] M. Zawirski, N. Preguiça, S. Duarte, A. Bieniusa, V. Balesgas, and M. Shapiro. Write fast, read in the past: Causal consistency for client-side applications. In *Middleware*, 2015.