

UNISTORE: A fault-tolerant marriage of causal and strong consistency

Manuel Bravo

Alexey Gotsman

Borja de Régil

Hengfeng Wei *

IMDEA Software Institute

Nanjing University

Abstract

Modern online services rely on data stores that replicate their data across geographically distributed data centers. Providing strong consistency in such data stores results in high latencies and makes the system vulnerable to network partitions. The alternative of relaxing consistency violates crucial correctness properties. A compromise is to allow multiple consistency levels to coexist in the data store. In this paper we present UNISTORE, the first fault-tolerant and scalable data store that combines causal and strong consistency. The key challenge we address in UNISTORE is to maintain liveness despite data center failures: this could be compromised if a strong transaction takes a dependency on a causal transaction that is later lost because of a failure. UNISTORE ensures that such situations do not arise while paying the cost of durability for causal transactions only when necessary. We evaluate UNISTORE on Amazon EC2 using both microbenchmarks and a sample application. Our results show that UNISTORE effectively and scalably combines causal and strong consistency.

1 Introduction

Many of today’s Internet services rely on geo-distributed data stores, which replicate data in different geographical locations. This improves user experience by allowing accesses to the closest site and ensures disaster-tolerance. However, geo-distribution also makes it more challenging to keep the data consistent. The classical approach is to make replication transparent to clients by providing strong consistency models, such as linearizability [33] or serializability [70]. The downside is that this approach requires synchronization between data centers in the critical path. This significantly increases latency [1] and makes the system unavailable during network partitionings [28]. Thus, even though several commercial geo-distributed systems follow this approach [18, 20, 26, 63, 71], the associated cost has prevented it from being adopted more widely.

An alternative approach is to relax synchronization: the data store executes an operation at a single data center, without any communication with others, and propagates updates to other data centers in the background [21, 66]. This minimizes the latency and makes the system *highly available*, i.e., operational even during network partitionings. But on the downside, the systems following this approach provide weaker consistency models: e.g., eventual consistency [66, 69] or *causal consistency* [2]. The latter is particularly appealing: it guarantees that clients see updates in an order that respects the potential causality between them. For example, assume that in a banking application Alice deposits \$100 into Bob’s account (u_1) and then posts a notification about it into Bob’s inbox (u_2). Under causal consistency, if Bob sees the notification (u_3), and then checks his account balance (u_4), he will see the deposit. This is not guaranteed under eventual consistency, which does not respect causality relationships, such as those between u_1 and u_2 . In some settings, causal consistency has been shown to be the strongest model that allows availability during network partitionings [7, 45]. It has been a subject of active research in recent years, with scalable implementations [3, 43, 48] and some industrial deployments [55, 67].

However, even causal consistency is often too weak to preserve critical application invariants. For example, consider a banking application that disallows overdrafts and thus maintains an invariant that an account balance is always non-negative. Assume that the balance of an account stored at two replicas is 100, and clients concurrently issue two `withdraw(100)` operations (u_5 and u_6) at different replicas. Since causal consistency executes operations without synchronization, both withdrawals will succeed, and once the replicas exchange the updates, the balance will go negative. To ensure integrity invariants in examples such as this, the programmer has to introduce synchronization between replicas, and, since synchronization is expensive, it pays off to do this sparingly. To this end, several research [9, 41, 42, 65] and commercial [5, 6, 29, 49, 58] data stores allow the programmer to choose whether to execute a particular operation under weak or strong consistency. For example, to preserve the integrity

*Also with the State Key Laboratory for Novel Software Technology, Software Institute.

invariant in our banking application, only withdrawals need to use strong consistency, and hence, synchronize; deposits may use weaker consistency and proceed without synchronization.

Given the benefits of causal consistency, it is particularly appealing to marry it with strong consistency in a geo-distributed data store. But like real-life marriages, to be successful this one needs to hold together both in good times and in bad – when data centers fail due to catastrophic events or power outages. Unfortunately, none of the existing data stores meant for geo-replication combine causal and strong consistency while providing such fault tolerance [9, 41, 42]. In this paper we present UNISTORE – the first fault-tolerant and scalable data store that combines causal and strong consistency. More precisely, UNISTORE implements a transactional variant of *Partial Order-Restrictions consistency (PoR consistency)* [31, 42]. This guarantees transactional causal consistency by default [3] and allows the programmer to additionally specify which pairs of transactions *conflict*, i.e., have to synchronize. For instance, to preserve the integrity invariant in our previous example, the programmer should declare that withdrawals from the same account conflict. Then one of the withdrawals u_5 and u_6 will observe the other and will fail.

The key challenge we have to address in UNISTORE is to maintain liveness despite data center failures. Just adding a Paxos-based commit protocol for strong transactions [19, 20, 35] to an existing causally consistent protocol does not yield a fault-tolerant data store. In such a data store, a committed strong transaction t_2 may never become visible to clients if a causal transaction t_1 on which it depends is lost due to a failure of its origin data center. This compromises the liveness of the system, because no transaction t_3 conflicting with t_2 can commit from now on: according to the PoR model, one of the transactions t_2 and t_3 has to observe the other, but t_2 will never be visible and t_2 did not observe t_3 .

UNISTORE addresses this problem by ensuring that, before a strong transaction commits, all its causal dependencies are *uniform*, i.e., will eventually become visible at all correct data centers. This adapts the classical notion of uniformity in distributed computing to causal consistency [16]. UNISTORE does so without defeating the benefits of causal consistency. Causal transactions remain highly available at the cost of increasing the latency of strong transactions: a strong transaction may have to wait for some of its dependencies to become uniform before committing. To minimize this cost, UNISTORE executes causal transactions on a snapshot that is slightly in the past, such that a strong transaction will mostly depend on causal transactions that are already uniform before committing. Furthermore, UNISTORE reuses the mechanism for tracking uniformity to let clients make causal transaction durable on demand and to enable consistent client migration.

In addition to being fault tolerant, UNISTORE scales horizontally, i.e., with the number of machines in each data center; this also goes beyond previous proposals [9, 41, 42]. To this end, UNISTORE builds on Cure [3] – a scalable implementa-

tion of transactional causal consistency. Our protocol extends Cure with a novel mechanism that distributes the task of tracking the set of uniform transactions among the machines of a data center. We also add the ability for data centers to forward transactions they receive from others, so that a transaction can propagate through the system even if its origin data center fails. Finally, we carefully integrate an existing fault-tolerant atomic commit for strong transactions [19] into the protocol for causal consistency.

We have rigorously proved the correctness of the UNISTORE protocol (§7 and [12, §D]). We have also evaluated it on Amazon EC2 using both microbenchmarks and a more realistic RUBiS benchmark. Our evaluation demonstrates that UNISTORE scalably combines causal and strong transactions, with the latter not affecting the performance of the former. Under the RUBiS mix workload, causal transactions exhibit a low latency (1.2ms on average), and the overall average latency is $3.7\times$ lower than that of a strongly consistent system.

2 System Model

We consider a geo-distributed system consisting of a set of data centers $\mathcal{D} = \{1, \dots, D\}$ that manage a large set of data items. A data item is uniquely identified by its *key*. For scalability, the key space is split into a set of logical partitions $\mathcal{P} = \{1, \dots, N\}$. Each data center stores replicas of all partitions, scattered among its servers. We let p_d^m be the replica of partition m at data center d , and we refer to replicas of the same partition as *sibling* replicas. As is standard, we assume that $D = 2f + 1$ and at most f data centers may fail. We call a data center that does not fail *correct*. If a data center fails, all partition replicas it stores become unavailable. For simplicity, we do not consider the failures of individual replicas within a data center: these can be masked using standard state-machine replication protocols executing within a data center [38, 56].

Replicas have physical clocks, which are loosely synchronized by a protocol such as NTP. The correctness of UNISTORE does not depend on the precision of clock synchronization, but large drifts may negatively impact its performance. Any two replicas are connected by a reliable FIFO channel, so that messages between correct data centers are guaranteed to be delivered. As is standard, to implement strong consistency we require the network to be *eventually synchronous*, so that message delays between sibling replicas in correct data centers are eventually bounded by some constant [24].

3 Consistency Model

A client interacts with UNISTORE by executing a stream of *transactions* at the data center it is connected to. A transaction consists of a sequence of operations, each on a single data item, and can be *interactive*: the data items it accesses are not known a priori. A transaction that modifies at least one data item is an *update* transaction; otherwise it is *read-only*.

A *consistency model* defines a contract between the data store and its clients that specifies which values the data store is allowed to return in response to client operations. UNISTORE implements a transactional variant of *Partial Order Restrictions consistency (PoR consistency)* [31, 42], which we now define informally; we give a formal definition in [12, §B]. The PoR model enables the programmer to classify transactions as either *causal* or *strong*. Causal transactions satisfy transactional *causal consistency*, which guarantees that clients see transactions in an order that respects the potential causality between them [2, 3]. However, clients can observe causally independent transactions in arbitrary order. Strong transactions give the programmer more control over their visibility. To this end, the programmer provides a symmetric *conflict relation* \bowtie on operations that is lifted to strong transactions as follows: two transactions conflict if they perform conflicting operations on the same data item. Then the PoR model guarantees that, out of two conflicting strong transactions, one has to observe the other.

More precisely, a transaction t_1 precedes a transaction t_2 in the *session order* if they are executed by the same client and t_1 is executed before t_2 . A set of transactions T committed by the data store satisfies PoR consistency if there exists a *causal order relation* \prec on T such that the following properties hold:

Causality Preservation. The relation \prec is transitive, irreflexive, and includes the session order.

Return Value Consistency. Consider an operation u on a data item k in a transaction $t \in T$. The return value of u can be computed from the state of k obtained as follows: first execute all operations on k by transactions preceding t in \prec in an order consistent with \prec ; then execute all operations on k that precede u in t .

Conflict Ordering. For any distinct strong transactions $t_1, t_2 \in T$, if $t_1 \bowtie t_2$, then either $t_1 \prec t_2$ or $t_2 \prec t_1$.

Eventual Visibility. A transaction $t \in T$ that is either strong or originates at a correct data center eventually becomes visible at all correct data centers: from some point on, t precedes in \prec all transactions issued at correct data centers.

If all transactions are causal, then the above definition specializes to transactional causal consistency [3, 17]. If all transactions are strong and all pairs of operations conflict, then we obtain (non-strict) serializability.

When $t_1 \prec t_2$, we say that t_1 is a *causal dependency* of t_2 . Return Value Consistency ensures that all operations in a transaction t execute on a snapshot consisting of its causal dependencies (as well as prior operations by t). Transactions are atomic, so that either all of their operations are included into the snapshot or none at all. The transitivity of \prec , mandated by Causality Preservation, ensures that the snapshot a transaction executes on is *causally consistent*: if a transaction t_1 is included into the snapshot, then so is any other transaction t_2 on which t_1 depends (i.e., $t_2 \prec t_1$). The inclusion of the session order into \prec , also mandated by Causality Preservation,

ensures session guarantees such as *read your writes* [64]. The consistency model disallows the causality violation anomaly from §1. Indeed, since \prec includes the session order, we have $u_1 \prec u_2$ and $u_3 \prec u_4$. Moreover, Bob sees Alice’s message, and by Return Value Consistency this can only happen if $u_2 \prec u_3$. Then since \prec is transitive, $u_1 \prec u_4$, and by Return Value Consistency, Bob has to see Alice’s deposit.

Causal consistency nevertheless allows the overdraft anomaly from §1: the withdrawals u_5 and u_6 may not be related by \prec , and thus may both execute on the balance 100 and succeed. The Conflict Ordering property can be used to disallow this anomaly by declaring that `withdraw` operations on the same account conflict and labeling transactions containing these as strong. Then one of the withdrawal transactions will be guaranteed to causally precede the other. The latter will be executed on the account balance 0 and will fail.

Finally, Eventual Visibility ensures that strong transactions and those causal ones that originate at correct data centers are durable, i.e., will eventually propagate through the system.

To facilitate the use of causal transactions, UNISTORE includes *replicated data types* (aka CRDTs), which implement policies for merging concurrent updates to the same data item [57]. Each data item in the store is associated with a type (e.g., counter, set), which is backed by a CRDT implementation managing updates to it. For example, the programmer can use a counter CRDT to represent an account balance. Then if two clients concurrently deposit 100 and 200 into an empty account using causal transactions, eventually the balance at all replicas will be 300. Using ordinary writes here would yield 100 or 200, depending on the order in which the writes are applied. More generally, CRDTs ensure that two replicas receiving the same set of updates are in the same state, regardless of the receipt order. Together with Eventual Visibility, this implies the expected guarantee of eventual consistency [66]. Due to space constraints, we omit details about the use of CRDTs from our protocol descriptions.

4 Key Design Decisions in UNISTORE

Baseline causal consistency. A causal transaction in UNISTORE first executes at a single data center on a causally consistent snapshot. After this it immediately commits, and its updates are replicated to all other data centers in the background. This minimizes the latency of causal transactions and makes them highly available, i.e., they can be executed even when the network connections between data centers fail.

As is common in causally consistent data stores [3, 23, 43], to ensure that causal transactions execute on consistent snapshots, a data center exposes a remote transaction to clients only after exposing all its dependencies. Then to satisfy the Eventual Visibility property under failures, a data center receiving a remote causal transaction may need to forward it to other data centers, as in reliable broadcast [11] and anti-entropy protocols for replica reconciliation [54].

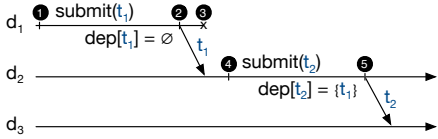


Figure 1: Why UNISTORE may need to forward remote causal transactions.

Figure 1 depicts a scenario that demonstrates how Eventual Visibility could be violated in the absence of this mechanism. Let t_1 be a causal transaction submitted at a data center d_1 (event 1). Assume that d_1 replicates t_1 to a correct data center d_2 (event 2) and then fails (event 3), so that t_1 does not get replicated anywhere else. Let t_2 be a transaction submitted at d_2 after t_1 becomes visible there, so that t_2 depends on t_1 (event 4). Transaction t_2 will eventually be replicated to all correct data centers (event 5). But it will never be exposed at any of them, because its dependency t_1 is missing. If data centers can forward remote causal transactions, then d_2 can eventually replicate t_1 to all correct data centers, preventing this problem.

On-demand strong consistency. UNISTORE uses optimistic concurrency control for strong transactions: they are first executed speculatively and the results are then *certified* to determine whether the transaction can commit, or must abort due to a conflict with a concurrent strong transaction [70]. Certifying a strong transaction requires synchronization between the replicas of partitions it accessed, located in different data centers. UNISTORE implements this using an existing fault-tolerant protocol that combines two-phase commit and Paxos [19] while minimizing commit latency. However, just using such a protocol is not enough to make the overall system fault tolerant: for this, before a strong transaction commits, all its causal dependencies must be uniform in the following sense.

DEFINITION 1. A transaction is *uniform* if both the transaction and its causal dependencies are guaranteed to be eventually replicated at all correct data centers.

This adapts the classical notion of uniformity in distributed computing to causal consistency [16]. UNISTORE considers a transaction to be uniform once it is visible at $f + 1$ data centers, because at least one of these must be correct, and data centers can forward causal transactions to others.

The following scenario, depicted in Figure 2, demonstrates why committing a strong transaction before its dependencies become uniform can compromise the liveness of the system. Assume that a causal transaction t_1 and a strong transaction t_2 are submitted at a data center d_1 in such a way that t_1 becomes a dependency of t_2 (events 1 and 2). Assume also that t_2 is certified, committed and delivered to all relevant replicas (events 3 and 4) before t_1 is replicated to any data center, and thus before it is uniform. Now if d_1 fails before

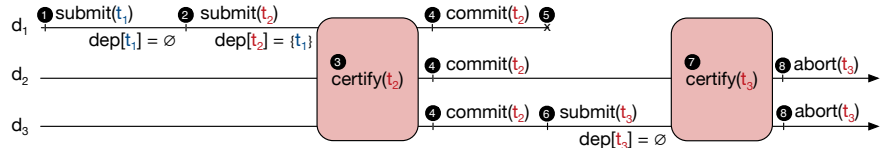


Figure 2: Why UNISTORE needs to ensure that the dependencies of a strong transaction are uniform before committing it.

replicating t_1 (event 5), no remote data center will be able to expose t_2 , because its dependency t_1 is missing. This violates the Eventual Visibility property, and even worse, no strong transaction conflicting with t_2 can commit from now on. For instance, let t_3 be such a transaction, submitted at d_3 (event 6). Because d_3 cannot expose t_2 , transaction t_3 executes on a snapshot excluding t_2 . Hence, t_3 will abort during certification (events 7 and 8): committing it would violate the Conflict Ordering property, since transactions t_2 and t_3 conflict, but neither of them is visible to the other. Ensuring that t_1 is uniform before committing t_2 prevents this problem, because it guarantees that t_1 will eventually be replicated at d_3 . After this t_2 will be exposed to conflicting transactions at this data center, which will allow them to commit.

Minimizing the latency of strong transactions. Ensuring that all the causal dependencies of a strong transaction are uniform before committing it may significantly increase its latency, since this requires additional communication between data centers. UNISTORE mitigates this problem by executing causal transactions on a snapshot that is slightly in the past, which is allowed by causal consistency. Namely, UNISTORE makes a remote causal transaction visible to the clients only after it is uniform. This minimizes the latency of a strong transaction, since to commit it only needs to wait for causal transactions originating at the local data center to become uniform. We cannot delay the visibility of the latter transactions due to the need to guarantee *read your writes* to local clients.

On-demand durability of causal transactions. Client applications interacting with the external world require hard durability guarantees: e.g., a banking application has to ensure that a withdrawal is durably recorded before authorizing the operation. UNISTORE guarantees that, once a strong transaction commits, the transaction and its dependencies are durable. However, UNISTORE returns from a causal transaction before it is replicated, and thus the transaction may be lost if its origin data center fails. Ensuring the durability of every single causal transaction would require synchronization between data centers on its critical path, defeating the benefits of causal consistency. Instead, UNISTORE reuses the mechanism for tracking uniformity to let the clients pay the cost of durability only when necessary. Even though UNISTORE replicates causal transactions asynchronously, it allows clients to execute a *uniform barrier*, which ensures that the transactions they have observed so far are uniform, and thus durable.

Client migration. Clients may need to migrate between data centers, e.g., because of roaming or for load balancing. UNISTORE also uses the uniformity mechanism to preserve session guarantees during migration. A client wishing to migrate from its local data center d to another data center i first invokes a uniform barrier at d . This guarantees that the transactions the client has observed or issued at d are durable and will eventually become visible at i , even if d fails. The client then makes an *attach* call at the destination data center i that waits until i stores all the above transactions. After this, the client can operate at i knowing that the state of the data center is consistent with the client’s previous actions.

Currently UNISTORE does not support consistent client migration in response to a data center failure: if the data center a client is connected to fails, the client will have to restart its session when connecting to a different data center. As shown in [72], this limitation can be lifted without defeating the benefits of causal consistency. We leave integrating the corresponding mechanisms into UNISTORE for future work.

5 Fault-Tolerant Causal Consistency Protocol

We first describe the UNISTORE protocol for the case when all transactions are causal. We give its pseudocode in Algorithms 1 and 2; for now the reader should ignore highlighted lines, which are needed for strong transactions. For simplicity, we assume that each handler in the algorithms executes atomically (although our implementation is parallelized). We reference pseudocode lines using the format `algorithm#:line#`.

5.1 Metadata

Most metadata in our protocol are represented by vectors with an entry per each data center, where each entry stores a scalar timestamp. However, different pieces of metadata use the vectors in different ways, which we now describe.

Tracking causality. The first use of the vectors is as vector clocks [27, 47], to track causality between transactions. Given vectors V_1 and V_2 , we write $V_1 < V_2$ if each entry of V_1 is no greater than the corresponding entry of V_2 , and at least one is strictly smaller. Each update transaction is tagged with a *commit vector* $commitVec$. The order on these vectors is consistent with the causal order \prec from §3: if $commitVec_1$ and $commitVec_2$ are the commit vectors of two update transactions t_1 and t_2 such that $t_1 \prec t_2$, then $commitVec_1 < commitVec_2$. For a transaction originating at a data center d with a commit vector $commitVec$, we call $commitVec[d]$ its *local timestamp*.

Each replica p_d^m maintains a log `opLog[k]` of update operations performed on each data item k stored at the replica. Each log entry stores, together with the operation, the commit vector of the transaction that performed it. This allows reconstructing different versions of a data item from its log.

Representing causally consistent snapshots. The protocol also uses a vector to represent a snapshot of the data store

on which a transaction operates: a snapshot vector V represents all transactions with a commit vector $\leq V$. This snapshot is causally consistent. Indeed, consider a transaction t_1 included into it, i.e., $commitVec_1 \leq V$. Since any causal dependency t_0 of t_1 is such that $commitVec_0 < commitVec_1$, we have $commitVec_0 < V$, so that t_0 is also included into the snapshot. A client also maintains a vector `pastVec` that represents its *causal past*: a causally consistent snapshot including the update transactions the client has previously observed.

Tracking what is replicated where. Each replica p_d^m maintains three vectors that are used to compute which transactions are uniform. These respectively track the sets of transactions replicated at p_d^m , the local data center d , and $f + 1$ data centers. Each of these vectors V represents the set of update transactions originating at a data center i with a local timestamp $\leq V[i]$. Note that this set may not form a causally consistent snapshot. The first vector maintained by p_d^m is `knownVec`. For each data center i , it defines the prefix of update transactions originating at i (in the order of local timestamps) that p_d^m knows about.

PROPERTY 1. For each data center i , the replica p_d^m stores the updates to partition m by all transactions originating at i with local timestamps $\leq \text{knownVec}[i]$.

Our protocol ensures that $\text{knownVec}[d] \leq \text{clock}$ at any replica in data center d . The vector `knownVec` at p_d^m records whether the updates to partition m by a given transaction are stored at this replica. In contrast, the next vector `stableVec` records whether the updates to *all* partitions by a transaction are stored at the local data center d .

PROPERTY 2. For each data center i , the data center d stores the updates by all transactions originating at i with local timestamps $\leq \text{stableVec}[i]$. More precisely, we are guaranteed that $\text{knownVec}[i]$ at any replica of $d \geq \text{stableVec}[i]$ at any p_d^m .

Finally, the last vector `uniformVec` defines the set of update transactions that p_d^m knows to have been replicated at $f + 1$ data centers, including d .

PROPERTY 3. Consider `uniformVec[i]` at p_d^m . All update transactions originating at i with local timestamps $\leq \text{uniformVec}[i]$ are replicated at $f + 1$ datacenters including d . More precisely: $\text{knownVec}[i]$ at any replica of these data centers $\geq \text{uniformVec}[i]$ at p_d^m .

When `uniformVec` is reinterpreted as a causally consistent snapshot, it defines transactions that p_d^m knows to be uniform according to Definition 1:

PROPERTY 4. Consider `uniformVec` at p_d^m . All update transactions with commit vectors $\leq \text{uniformVec}$ are uniform.

Proof sketch. Consider a transaction t_1 that originates at a data center i with a commit vector $commitVec_1 \leq \text{uniformVec}$ at p_d^m . In particular, $commitVec_1[i] \leq \text{uniformVec}[i]$, and by Property 3, t_1 is replicated at $f + 1$ data centers. We as-

sume at most f failures. Then the transaction forwarding mechanism of our protocol (§4) guarantees that t_1 will eventually be replicated at all correct data centers. Consider now any causal dependency t_2 of t_1 with a commit vector $commitVec_2$. Since commit vectors are consistent with causality, $commitVec_2 < commitVec_1 \leq uniformVec$. Then as above, we can again establish that t_2 will be replicated at all correct data centers, as required by Definition 1. \square

5.2 Causal Transaction Execution

Starting a transaction. A client can submit a transaction to any replica in its local data center by calling $START_TX(V)$, where V is the client’s causal past $pastVec$ (line 1:1, for brevity, we omit the pseudocode of the client). A replica p_d^m receiving such a request acts as the transaction *coordinator*. It generates a unique transaction identifier tid , computes a snapshot $snapVec[tid]$ on which the transaction will execute, and returns tid to the client (we explain lines 1:2-3 and similar ones later). The snapshot is computed by combining uniform transactions from $uniformVec$ (line 1:5) with the transactions from the client’s causal past originating at d (line 1:6). The former is crucial to minimize the latency of strong transactions (§4), while the latter ensures *read your writes*.

Transaction execution. The client proceeds to execute the transaction tid by issuing a sequence of operations at its coordinator via DO_OP (line 1:9). When the coordinator receives an operation op on a data item k , it sends a $GET_VERSION$ message with the transaction’s snapshot $snapVec[tid]$ to the local replica responsible for k (line 1:11). Upon receiving the message (line 1:18), the replica first ensures that it is as up-to-date as required by the snapshot (line 1:21). It then computes the latest version of k within the snapshot by applying the operations from $opLog[k]$ by all transactions with commit vectors $\leq snapVec[tid]$. The result is sent to the coordinator in a $VERSION$ message. After receiving it (line 1:12), the coordinator further applies the operations on k previously executed by the transaction, which are stored in a buffer $wbuff[tid]$; this ensures *read your writes* within the transaction. If the operation is an update, the coordinator then appends it to $wbuff[tid]$. Finally, the coordinator executes the desired operation op and forwards its return value to the client.

Commit. A client commits a causal transaction by calling $COMMIT_CAUSAL$ (line 1:26). This returns immediately if the transaction is read-only, since it already read a consistent snapshot (line 1:28). To commit an update transaction, $UNISTORE$ uses a variant of two-phase commit protocol (recall that for simplicity we only consider whole-data center failures, not those of individual replicas, §2). The coordinator first sends a $PREPARE$ message to the replicas in the local data center storing the data items updated by the transaction (line 1:29). The message to each replica contains the part of the write buffer relevant to that replica. When a replica receives the message (line 1:36), it computes the transaction’s

Algorithm 1 Transaction execution at p_d^m .

```

1: function START_TX( $V$ )
2:   for  $i \in \mathcal{D} \setminus \{d\}$  do
3:      $uniformVec[i] \leftarrow \max\{V[i], uniformVec[i]\}$ 
4:    $var\ tid \leftarrow generate\_tid()$ 
5:    $snapVec[tid] \leftarrow uniformVec$ 
6:    $snapVec[tid][d] \leftarrow \max\{V[d], uniformVec[d]\}$ 
7:    $snapVec[tid][strong] \leftarrow \max\{V[strong], stableVec[strong]\}$ 
8:   return  $tid$ 

9: function DO_OP( $tid, k, op$ )
10:   $var\ l \leftarrow partition(k)$ 
11:  send GET_VERSION( $snapVec[tid], k$ ) to  $p_d^l$ 
12:  wait receive VERSION( $state$ ) from  $p_d^l$ 
13:  for all  $\langle k, op' \rangle \in wbuff[tid][l]$  do  $state \leftarrow apply(op', state)$ 
14:   $rset[tid] \leftarrow rset[tid] \cup \{\langle k, op \rangle\}$ 
15:  if  $op$  is an update then
16:     $wbuff[tid][l] \leftarrow wbuff[tid][l] \cdot \langle k, op \rangle$ 
17:  return  $retval(op, state)$ 

18: when received GET_VERSION( $snapVec, k$ ) from  $p$ 
19:  for  $i \in \mathcal{D} \setminus \{d\}$  do
20:     $uniformVec[i] \leftarrow \max\{snapVec[i], uniformVec[i]\}$ 
21:  wait until  $knownVec[d] \geq snapVec[d] \wedge$ 
     $knownVec[strong] \geq snapVec[strong]$ 
22:   $var\ state \leftarrow \perp$ 
23:  for all  $\langle op', commitVec \rangle \in opLog[k].commitVec \leq snapVec$  do
24:     $state \leftarrow apply(op', state)$ 
25:  send VERSION( $state$ ) to  $p$ 

26: function COMMIT_CAUSAL( $tid$ )
27:   $var\ L \leftarrow \{l \mid wbuff[tid][l] \neq \emptyset\}$ 
28:  if  $L = \emptyset$  then return  $snapVec[tid]$ 
29:  send PREPARE( $tid, wbuff[tid][l], snapVec[tid]$ ) to  $p_d^l, l \in L$ 
30:   $var\ commitVec \leftarrow snapVec[tid]$ 
31:  for all  $l \in L$  do
32:    wait receive PREPARE_ACK( $tid, ts$ ) from  $p_d^l$ 
33:     $commitVec[d] \leftarrow \max\{commitVec[d], ts\}$ 
34:  send COMMIT( $tid, commitVec$ ) to  $p_d^l, l \in L$ 
35:  return  $commitVec$ 

36: when received PREPARE( $tid, wbuff, snapVec$ ) from  $p$ 
37:  for  $i \in \mathcal{D} \setminus \{d\}$  do
38:     $uniformVec[i] \leftarrow \max\{snapVec[i], uniformVec[i]\}$ 
39:   $var\ ts \leftarrow clock$ 
40:   $preparedCausal \leftarrow preparedCausal \cup \{\langle tid, wbuff, ts \rangle\}$ 
41:  send PREPARE_ACK( $tid, ts$ ) to  $p$ 

42: when received COMMIT( $tid, commitVec$ )
43:  wait until  $clock \geq commitVec[d]$ 
44:   $\langle tid, wbuff, \_ \rangle \leftarrow find(tid, preparedCausal)$ 
45:   $preparedCausal \leftarrow preparedCausal \setminus \{\langle tid, \_, \_ \rangle\}$ 
46:  for all  $\langle k, op \rangle \in wbuff$  do
47:     $opLog[k] \leftarrow opLog[k] \cdot \langle op, commitVec \rangle$ 
48:   $committedCausal[d] \leftarrow committedCausal[d] \cup$ 
     $\{\langle tid, wbuff, commitVec \rangle\}$ 

49: function UNIFORM_BARRIER( $V$ )
50:  wait until  $uniformVec[d] \geq V[d]$ 

51: function ATTACH( $V$ )
52:  wait until  $\forall i \in \mathcal{D} \setminus \{d\}. uniformVec[i] \geq V[i]$ 

```

prepare time ts from its local clock and adds the transaction to preparedCausal, which stores the set of causal transactions that are prepared to commit at the replica. The replica then returns ts to the coordinator in a PREPARE_ACK message.

When the coordinator receives replies from all replicas updated by the transaction, it computes the transaction's commit vector $commitVec$: it sets the local timestamp $commitVec[d]$ to the maximum among the prepare times proposed by the replicas (line 1:33), and it copies the other entries of $commitVec$ from the snapshot vector $snapVec[tid]$ (line 1:30). The latter reflects the fact that the transactions in the snapshot become causal dependencies of tid .

After computing $commitVec$, the coordinator sends it in a COMMIT message to the relevant replicas at the local data center (line 1:34) and returns it to the client (line 1:35). The client then sets its causal past pastVec to the commit vector. When a replica receives the COMMIT message (line 1:42), it removes the transaction from preparedCausal, adds the transaction's updates to opLog, and adds the transaction to a set committedCausal[d], which stores transactions waiting to be replicated to sibling replicas at other data centers.

5.3 Transaction Replication

Each replica p_d^m periodically replicates locally committed update transactions to sibling replicas in other data centers by executing PROPAGATE_LOCAL_TXS (line 2:1). Transactions are replicated in the order of their local timestamps. The prefix of transactions that is ready to be replicated is determined by knownVec[d]: according to Property 1, p_d^m stores updates to m by all transactions originating at d with local timestamps \leq knownVec[d]. Thus, the replica first updates knownVec[d] while preserving Property 1.

There are two cases of this update. If the replica does not have any prepared transactions (preparedCausal = \emptyset), it sets knownVec[d] to the current value of the clock (line 2:2). This preserves Property 1 because in this case a new transaction originating at d and updating m will get a prepare time at m higher than the current clock (line 1:39), and thus also a higher local timestamp (line 1:33). If the replica has some prepared transactions, then they may end up getting local timestamps lower than the current clock. In this case, the replica sets knownVec[d] to just below the smallest prepared time (line 2:3). This preserves Property 1 because: (i) currently prepared transactions will get a local timestamp no lower than their prepare time; and (ii) as we argued above, new transactions will get a prepare time higher than the current clock and, hence, than the smallest prepare time.

After updating knownVec[d], the replica sends a REPLICATE message to its siblings with the transactions in committedCausal[d] such that $commitVec[d] \leq$ knownVec[d], and then removes them from committedCausal[d]. In other words, the replica sends all transactions from the prefix determined by knownVec[d] that it has not yet replicated.

When a replica p_d^m receives a REPLICATE message with

Algorithm 2 Transaction replication at p_d^m .

```

1: function PROPAGATE_LOCAL_TXS()      ▷ Run periodically
2:   if preparedCausal =  $\emptyset$  then knownVec[ $d$ ]  $\leftarrow$  clock
3:   else knownVec[ $d$ ]  $\leftarrow$  min{ $ts \mid \langle \_, \_, ts \rangle \in$  preparedCausal} - 1
4:   var txs  $\leftarrow$  { $\langle \_, \_, commitVec \rangle \in$  committedCausal[ $d$ ] |
                                      $commitVec[d] \leq$  knownVec[ $d$ ]}
5:   if txs  $\neq \emptyset$  then
6:     send REPLICATE( $d$ , txs) to  $p_i^m, i \in \mathcal{D} \setminus \{d\}$ 
7:     committedCausal[ $d$ ]  $\leftarrow$  committedCausal[ $d$ ] \setminus txs
8:   else send HEARTBEAT( $d$ , knownVec[ $d$ ]) to  $p_i^m, i \in \mathcal{D} \setminus \{d\}$ 
9:   when received REPLICATE( $i$ , txs)
10:  for all  $\langle tid, wbuff, commitVec \rangle \in$  txs in  $commitVec[i]$  order do
11:    if  $commitVec[i] >$  knownVec[ $i$ ] then
12:      for all  $\langle k, op \rangle \in$  wbuff do
13:        opLog[ $k$ ]  $\leftarrow$  opLog[ $k$ ]  $\cdot \langle op, commitVec \rangle$ 
14:        committedCausal[ $i$ ]  $\leftarrow$  committedCausal[ $i$ ]  $\cup$ 
                                     { $\langle tid, wbuff, commitVec \rangle$ }
15:        knownVec[ $i$ ]  $\leftarrow$   $commitVec[i]$ 
16:  when received HEARTBEAT( $i$ , ts)
17:  pre:  $ts >$  knownVec[ $i$ ]
18:  knownVec[ $i$ ]  $\leftarrow$  ts
19:  function FORWARD_REMOTE_TXS( $i$ ,  $j$ )
20:  var txs  $\leftarrow$  { $\langle \_, \_, commitVec \rangle \in$  committedCausal[ $j$ ] |
                                      $commitVec[j] >$  globalMatrix[ $i$ ][ $j$ ]}
21:  if txs  $\neq \emptyset$  then send REPLICATE( $j$ , txs) to  $p_i^m$ 
22:  else send HEARTBEAT( $j$ , knownVec[ $j$ ]) to  $p_i^m$ 
23:  function BROADCAST_VECS()          ▷ Run periodically
24:  send KNOWNVEC_LOCAL( $m$ , knownVec) to  $p_l^l, l \in \mathcal{P}$ 
25:  send STABLEVEC( $d$ , stableVec) to  $p_i^m, i \in \mathcal{D}$ 
26:  send KNOWNVEC_GLOBAL( $d$ , knownVec) to  $p_i^m, i \in \mathcal{D}$ 
27:  when received KNOWNVEC_LOCAL( $l$ , knownVec)
28:  localMatrix[ $l$ ]  $\leftarrow$  knownVec
29:  for  $i \in \mathcal{D}$  do stableVec[ $i$ ]  $\leftarrow$  min{localMatrix[ $n$ ][ $i$ ] |  $n \in \mathcal{P}$ }
30:  stableVec[strong]  $\leftarrow$  min{localMatrix[ $n$ ][strong] |  $n \in \mathcal{P}$ }
31:  when received STABLEVEC( $i$ , stableVec)
32:  stableMatrix[ $i$ ]  $\leftarrow$  stableVec
33:   $G \leftarrow$  all groups with  $f + 1$  replicas that include  $p_d^m$ 
34:  for  $j \in \mathcal{D}$  do
35:    var  $ts \leftarrow$  max{min{stableMatrix[ $h$ ][ $j$ ] |  $h \in g$ } |  $g \in G$ }
36:    uniformVec[ $j$ ]  $\leftarrow$  max{uniformVec[ $j$ ],  $ts$ }
37:  when received KNOWNVEC_GLOBAL( $l$ , knownVec)
38:  globalMatrix[ $l$ ]  $\leftarrow$  knownVec

```

a set of transactions txs originating at a sibling replica p_i^m (line 2:9), it iterates over txs in $commitVec[i]$ order. For each new transaction in txs with commit vector $commitVec$, the replica adds the transaction's operations to its log and sets knownVec[i] = $commitVec[i]$. Since communication channels are FIFO, p_d^m processes all transactions from p_i^m in their local timestamp order. Hence, the above update to knownVec[i] preserves Property 1: p_d^m stores updates originating at p_i^m by all transactions with $commitVec[i] \leq$ knownVec[i]. Finally, the replica adds the transactions to committedCausal[i], which is used to implement transaction forwarding (§4). Due to

the forwarding, p_d^m may receive the same transaction from different data centers. Thus, when processing transactions in the REPLICATE message, it checks for duplicates (line 2:11).

5.4 Advancing the Uniform Snapshot

Replicas run a background protocol that refreshes the information about uniform transactions. This proceeds in two stages. First, a replica keeps track of which transactions have been replicated at the replicas of other partitions in the same data center. To this end, replicas in the same data center periodically exchange KNOWNVEC_LOCAL messages with their knownVec vectors, which they store in a matrix localMatrix (lines 2:24 and 2:27); in our implementation this is done via a dissemination tree. This matrix is then used to compute the vector stableVec, which represents the set of transactions that have been fully replicated at the local data center as per Property 2. To ensure this, a replica computes an entry stableVec[i] as the minimum of knownVec[i] it received from the replicas of other partitions in the same data center (line 2:29).

In the second stage of the background protocol, sibling replicas periodically exchange STABLEVEC messages with their stableVec vectors, which they store in a matrix stableMatrix (lines 2:25 and 2:31). This matrix is then used by a replica to compute uniformVec, which characterizes the update transactions that are replicated at $f + 1$ data centers as per Property 3. To this end, a replica first enumerates all groups G of $f + 1$ data centers that include its local data center (line 2:33). For each data center j the replica performs the following computation. First, for each group $g \in G$, it computes the minimum j -th entry in the stable vectors of all data centers $h \in g$: $\min\{\text{stableMatrix}[h][j] \mid h \in g\}$. By Property 2 all update transactions originating at j with local timestamp \leq the minimum have been replicated at all data centers in g . The replica then sets uniformVec[j] to the maximum of the resulting values computed for all groups $g \in G$, to cover transactions that are replicated at any such group. According to Property 4, the transactions with commit vectors \leq uniformVec are uniform, and now become visible to transactions coordinated by p_d^m (§5.2).

Replicas also update uniformVec in lines 1:2-3, 1:19-20 and 1:37-38 by incorporating snapVec[i] for remote data centers i . This is safe because a transaction executes on a snapshot that only includes uniform remote transactions.

Finally, if a replica does not receive new transactions for a long time, it sends the value of its knownVec[d] as a heartbeat (lines 2:8 and 2:16). This allows advancing stableVec and uniformVec even under skewed load distributions.

5.5 Transaction Forwarding

As we explained in §4, to guarantee that a transaction originating at a correct data center eventually becomes exposed at all correct data centers despite failures (Eventual Visibility), replicas may have to forward remote update transactions. To determine which transactions to forward, each replica keeps track

of the update transactions that have been replicated at sibling replicas in other data centers. To this end, sibling replicas periodically exchange KNOWNVEC_GLOBAL messages with their knownVec vectors, which they store in a matrix globalMatrix (lines 2:26 and 2:37). Thus, p_i^m has received all update transactions from p_j^m with $\text{commitVec}[j] \leq \text{globalMatrix}[i][j]$.

A replica p_d^m only forwards transactions when it suspects that a data center j may have failed before replicating all the update transactions originating at it to a data center i (this information is provided by a separate module). In this case, p_d^m executes FORWARD_REMOTE_TXS(i, j) (line 2:19). The function forwards the set of transactions txs received from p_j^m that have not been replicated at p_i^m according to $\text{globalMatrix}[i][j]$. For example, in Figure 1, UNISTORE will eventually invoke FORWARD_REMOTE_TXS(d_1, d_3) at replicas in d_2 to forward t_1 . The replica p_d^m sends the transactions in txs to p_i^m in a REPLICATE message. If there are no update transactions to forward, p_d^m sends a heartbeat to p_i^m with knownVec[j].

UNISTORE periodically deletes from committedCausal transactions that have been replicated at every data center (omitted from the pseudocode for brevity).

5.6 On-Demand Durability and Client Migration

A client may wish to ensure that the transactions it has observed so far are durable. To this end, the client can call UNIFORM_BARRIER(V) at any replica in its local data center d , where V is the client's causal past pastVec (line 1:49). The replica returns to the client only when all the transactions from pastVec that originate at d are uniform, and thus durable. Then the same holds for all transactions from pastVec, because the protocol only exposes remote transactions to clients when they are already uniform (§5.2).

A client wishing to migrate from its local data center d to another data center i first calls UNIFORM_BARRIER(V) at any replica in d with $V = \text{pastVec}$, to ensure that the transactions the client has observed or issued at d will eventually become visible at i . The client then calls ATTACH(V) at any replica in i (line 1:51). The replica returns when its uniformVec includes all remote transactions from V (line 1:52). The client can then be sure that its transactions at i will operate on snapshots including all the transactions it has observed before.

6 Adding Strong Transactions

We now describe the full UNISTORE protocol with both causal and strong transactions. It is obtained by adding the highlighted lines to Algorithms 1-2 and a new Algorithm 3.

6.1 Metadata

The Conflict Ordering property of our consistency model requires any two conflicting strong transactions to be related one way or another by the causal order \prec (§3). To ensure this, the protocol assigns to each strong transaction a scalar

strong timestamp, analogous to those used in optimistic concurrency control for serializability [70]. Several vectors used as metadata in the causal consistency protocol (§5.1) are then extended with an extra *strong* entry.

First, we extend commit vectors and those representing causally consistent snapshots. Commit vectors are compared using the previous order $<$, but considering all entries; as before, this order is consistent with the causal order \prec . Furthermore, conflicting strong transactions are causally ordered according to their strong timestamps.

PROPERTY 5. For any conflicting strong transactions t_1 and t_2 with commit vectors $commitVec_1$ and $commitVec_2$, we have: $t_1 \prec t_2 \iff commitVec_1[strong] < commitVec_2[strong]$.

A consistent snapshot vector V now defines the set of transactions with a commit vector $\leq V$, according to the new $<$. The vectors $knownVec$ and $stableVec$ maintained by a replica p_d^m are also extended with a *strong* entry. The entries $knownVec[strong]$ and $stableVec[strong]$ define the prefix of strong transactions that have been replicated at p_d^m and the local data center d , respectively:

PROPERTY 6. Replica p_d^m stores the updates to m by all strong transactions with $commitVec[strong] \leq knownVec[strong]$.

PROPERTY 7. Data center d stores the updates by all strong transactions with $commitVec[strong] \leq stableVec[strong]$.

To ensure Property 7, the *strong* entry of $stableVec$ is updated at line 2:30 similarly to its other entries. We do not extend $uniformVec$, because our commit protocol for strong transactions automatically guarantees their uniformity.

6.2 Transaction Execution

UNISTORE uses optimistic concurrency control for strong transactions, with the same protocol for executing causal and speculatively executing strong transactions. To this end, Algorithm 1 is modified as follows. First, the computation of the snapshot vector $snapVec[tid]$ is extended to compute the *strong* entry (line 1:7), which is now taken into account when checking that a replica state is up to date (line 1:21). The *strong* entry of the snapshot vector is computed so as to include all strong transactions known to be fully replicated in the local data center, as defined by $stableVec[strong]$. To ensure *read your writes*, the snapshot additionally includes strong transactions from the client’s causal past, as defined by $V[strong]$. Finally, the coordinator of a transaction now maintains not only its write set, but also its read set $rset$ that records all operations by the transaction, including read-only ones (line 1:14). The latter is used to certify strong transactions.

After speculatively executing a strong transaction, the client tries to commit it by calling `COMMIT_STRONG` at its coordinator (line 3:1). The coordinator first waits until the snapshot on which the transaction operated becomes uniform by calling `UNIFORM_BARRIER` (line 3:2): as we argued in §4, this is crucial for liveness. The coordinator next submits the trans-

Algorithm 3 Committing strong transactions at p_d^m .

```

1: function COMMIT_STRONG(tid)
2:   UNIFORM_BARRIER(snapVec[tid])
3:   return CERTIFY(tid, wbuff[tid], rset[tid], snapVec[tid])

4: upon DELIVER_UPDATES(W)
5:   for all  $\langle wbuff, commitVec \rangle \in W$  in  $commitVec[strong]$  order do
6:     for all  $\langle k, op \rangle \in wbuff$  do
7:       opLog[k]  $\leftarrow$  opLog[k]  $\cdot \langle op, commitVec \rangle$ 
8:       knownVec[strong]  $\leftarrow$  commitVec[strong]

9: function HEARTBEAT_STRONG() ▷ Run periodically
10: return CERTIFY( $\perp, \emptyset, \emptyset, \vec{0}$ )

```

action to a *certification service*, which determines whether the transaction commits or aborts (line 3:3, see §6.3). In the former case, the service also determines its commit vector, which the coordinator returns to the client. If the transaction commits, the client sets its causal past $pastVec$ to the commit vector; otherwise, it re-executes the transaction.

The certification service also notifies replicas in all data centers about updates by strong transactions affecting them via `DELIVER_UPDATES` upcalls, invoked in an order consistent with strong timestamps of the transactions (line 3:4). A replica receiving an upcall adds the new operations to its log and refreshes $knownVec[strong]$ to preserve Property 6.

Finally, a replica p_d^m that has not seen any strong transactions updating its partition m for a long time submits a dummy strong transaction that acts as a heartbeat (line 3:9). Similarly to heartbeats for causal transactions, this allows coping with skewed load distributions.

6.3 Certification Service

We implement the certification service using an existing fault-tolerant protocol from [19], with transaction commit vectors computed using the techniques from [30]. The protocol integrates two-phase commit across partitions accessed by the transaction and Paxos among the replicas of each partition. It furthermore uses white-box optimizations between the two protocols to minimize the commit latency. The use of Paxos ensures that a committed strong transaction is durable and its updates will eventually be delivered at all correct data centers (line 3:4). For each partition, a single replica functions as the Paxos leader. The protocol is described and formally specified elsewhere [19], and here we discuss it only briefly. Its pseudocode and formal specification are given in [12, §A] and [12, §C], respectively.

The certification service accepts the read and write sets of a transaction and its snapshot vector (line 3:3). Even though the service is distributed, it guarantees that commit/abort decisions are computed like in a centralized database with optimistic concurrency control – in a total *certification order*. To ensure Conflict Ordering, the decisions are computed using a concurrency-control policy similar to that for serializability [70]: a transaction commits if its snapshot includes all

conflicting transactions that precede it in the certification order. The certification service also computes a commit vector for each committed transaction by copying its per-data center entries from the transaction’s snapshot vector and assigning a strong timestamp consistent with the certification order.

7 Proof of Correctness

We have rigorously proved that UNISTORE correctly implements the specification of PoR consistency for the case when the data store manages last-writer-wins registers. The proof uses the formal framework from [14, 15, 17] and establishes Properties 1-7 stated earlier. Due to space constraints, we defer the proof to [12, §D].

8 Evaluation

We have implemented UNISTORE and several other protocols (listed in the following) in the same codebase, consisting of 10.3K SLOC of Erlang. We evaluate the protocols on Amazon EC2 using m4.2xlarge VMs from 5 different regions. Each VM has 8 virtual cores and 32GB of RAM. The RTT between regions ranges from 26ms to 202ms. Unless otherwise stated, our experiments deploy 3 data centers, thus tolerating a single data center failure: Virginia (US-East), California (US-West) and Frankfurt (EU-FRA). All Paxos leaders are located in Virginia. By default we use 4 replica machines per data center. Each machine stores replicas of 8 partitions, matching the number of cores. Clients are hosted on separate machines in each data center. We run each experiment for at least 5 minutes, with the first and the last minute ignored. Replicas propagate local update transactions (line 2:1) and broadcast vectors (line 2:23) every 5ms.

8.1 Does UNISTORE combine causal and strong consistency effectively?

We start by analyzing the performance of UNISTORE using RUBiS – a popular benchmark that emulates an online auction website such as eBay [41, 42]. It defines 11 read-only transactions and 5 update transactions, e.g., selling items, bidding on items, and consulting outstanding auctions. As in previous work [42], to make the benchmark more challenging, we add an extra update transaction `closeAuction` to declare the winner of an auction. We also borrow from [42] a conflict relation between RUBiS transactions that preserves key integrity invariants in the PoR consistency model. This marks four transactions as strong (`registerUser`, `storeBuyNow`, `storeBid` and `closeAuction`) and declares three conflicts between them. For example, `storeBid`, which places a bid on a item, conflicts with `closeAuction` if both act on the same item: this is needed to preserve the invariant that the winner of an auction is the highest bidder. Our RUBiS database is configured according to the benchmark specification: it is populated with 33,000 items for sale and 1 million users; client

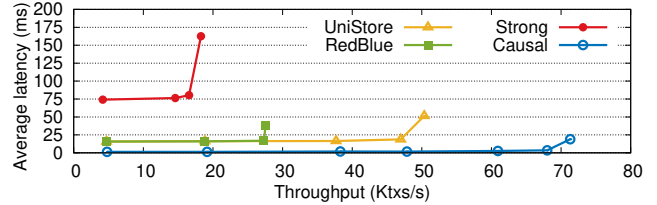


Figure 3: RUBiS benchmark: throughput vs. average latency.

think times are 500ms. We run the bidding mix workload of RUBiS with 15% of update transactions, which yields 10% of strong transactions.

We compare UNISTORE with STRONG, REDBLUE and CAUSAL. STRONG implements serializability [70] as a special case of UNISTORE where all transactions are strong and all pairs of operations conflict. REDBLUE implements red-blue consistency [41], which like PoR, combines causal and strong consistency. However, it declares conflicts between all strong transactions. REDBLUE certifies strong transactions at a centralized replicated service, with a replica at each data center. CAUSAL implements causal consistency as a special case of UNISTORE where all transactions are causal. It cannot preserve the integrity invariants of RUBiS, but gives an upper bound on the expected performance.

Throughput and average latency. Figure 3 evaluates average transaction latency and throughput. As the figure shows, UNISTORE exhibits a high throughput: 72% and 183% higher than REDBLUE and STRONG respectively at their saturation point. This is expected, as UNISTORE implements the consistency model that enables the most concurrency. STRONG classifies all transactions as strong. This impacts performance because executing a strong transaction is significantly more expensive than executing a causal one. REDBLUE uses a centralized certification service that saturates before the UNISTORE’s distributed service, creating a bottleneck. UNISTORE exhibits an average latency of 16.5ms, lower than 80.4ms of STRONG. The latency of REDBLUE is comparable to that of UNISTORE. This is because both systems mark the same set of transactions as strong. Still, REDBLUE declares conflicts between all strong transactions and thus aborts more transactions than UNISTORE: 0.12% vs 0.027%. The clients whose transactions abort have to retry them, thus increasing latency. Since the abort rate remains low in both cases, the difference in latency is negligible in our experiment. We expect a more significant difference in workloads with higher contention. Finally, in comparison to CAUSAL, UNISTORE penalizes throughput by 45%. This is the unavoidable price to pay to preserve application-specific invariants.

Latency of each transaction type. In UNISTORE, the latency of strong transactions is dominated by the RTT between Virginia (the leader’s region) and California (Virginia’s closest data center) – 61ms. Strong transactions exhibit a latency of 73.9ms on average. The latency varies depending on the

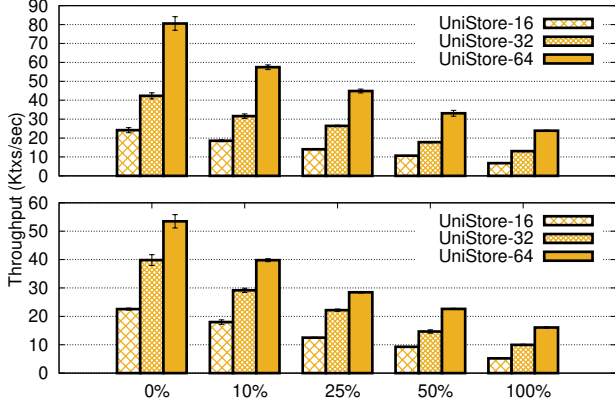


Figure 4: Scalability when varying the ratio of strong transactions with uniform data access (top) and under contention (bottom).

client’s location: from 65.4ms on average at the leader’s site to 93.2ms at the site furthest from the leader (Frankfurt). Since causal transactions do not require coordination between data centers, they exhibit a very low latency – 1.2ms on average, which is comparable to that of CAUSAL. This demonstrates that UNISTORE is able to mix causal and strong consistency effectively, as the latency of causal transactions remains low regardless of concurrently executing strong transactions.

8.2 How does UNISTORE scale with the number of machines?

We evaluate the peak throughput of UNISTORE as we increase the number of machines per each data center from 2 to 8, i.e., the number of partitions from 16 to 64. We use a microbenchmark with 100% of update transactions, where each transaction accesses three data items. We vary the ratio of strong transactions from 0% to 100% to understand their impact on scalability.

Scalability under low contention. For this set of experiments, the data items accessed by each transaction are picked uniformly at random. This yields a very low contention: e.g., with 16 partitions, the probability of two transactions accessing the same partition is 0.031. As shown by the top plot of Figure 4, UNISTORE is able to scale almost linearly even when the workload includes strong transactions: a 9.76% throughput drop compared to the optimal scalability. This is because, with uniform accesses, the task of committing transactions is balanced among partitions. Thus, when the number of partitions increases, so does the system’s capacity. The scalability is not perfect due to the cost of the background protocol that computes stableVec, which grows logarithmically with the number of partitions. The plot also shows that strong transactions are expensive: 25.72% of throughput drop on average with 10% of strong transactions. The performance is dominated by the number of strong transactions that a partition can certify per second.

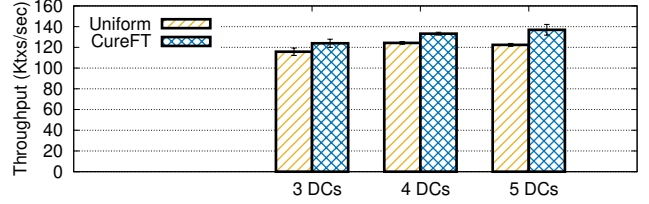


Figure 5: Throughput penalty of tracking uniformity.

Impact of contention. For this set of experiments, we set the ratio of strong transactions that access a designated partition to 20% to create contention. As shown by the bottom plot of Figure 4, UNISTORE is still able to scale fairly well under contention. But, as expected, contention has an impact on scalability: a 17.15% throughput drop from the optimal scalability compared to the 9.76% throughput drop in the experiments without contention.

8.3 What is the cost of uniformity?

We compare CUREFT to UNIFORM. CUREFT implements Cure [3], a causally consistent data store, and makes it fault tolerant by adding transaction forwarding (§4). UNIFORM is a simplified version of UNISTORE that removes all the mechanisms related to strong transactions. UNIFORM tracks uniformity and makes remote transactions visible only when these are uniform; CUREFT does not. We use a microbenchmark with only causal transactions and 15% of update transactions. Each transaction accesses three data items.

Throughput penalty. Figure 5 evaluates the cost of tracking uniformity. It shows the peak throughput when the number of data centers increases from 3 to 5. We first add Ireland and then Brazil. As we do this, the throughput remains almost constant. This is because each data center stores replicas of all partitions and the computational power gained when adding a data center is offset by the cost of replicating update transactions. As the figure shows, the cost of tracking uniformity is small: a 7.97% drop on average. The gap grows as we increase the number of data centers: a 10.61% drop on average with 5 data centers. This is because, to track uniform transactions, sibling replicas exchange messages: the more data centers, the more messages exchanged. The penalty can be reduced by decreasing the frequency at which sibling replicas exchange their stableVec (line 2:25), at the expense of an extra delay in the visibility of remote transactions.

Reading from a uniform snapshot. Figure 6 evaluates the delay on the visibility of remote transactions when reading from a uniform snapshot. We deploy four data centers: Virginia, California, Frankfurt and Brazil. We set $f = 2$ to tolerate 2 data center failures (when $f = 1$, UNIFORM shows no delay). Under such a configuration, a data center makes a transaction visible when it knows that 3 data centers store the transaction and its dependencies (§5.4). The figure shows the cumulative distribution of the delay before updates from

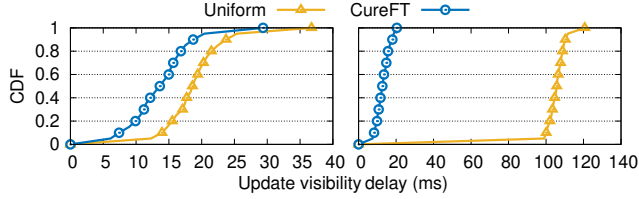


Figure 6: Left: California to Brazil (best case). Right: California to Virginia (worst case).

California are visible in Brazil and Virginia.

The extra delay at Brazil is only of 5ms at the 90th percentile. This is the best case scenario for UNIFORM because Brazil learns that Virginia stores a transaction originating at California only 2ms after receiving it. The worst case scenario for UNIFORM is when the origin and the destination data center are the closest ones. This is why the extra delay at Virginia is 92ms at the 90th percentile: Brazil learns that Frankfurt stores a transaction originating at California 88ms after receiving it. Note that when clients communicate only via the data store, the delay is unnoticeable. Even if clients communicate out of band, as the maximum extra delay is less than 100ms, it is unlikely that a client will miss an update.

9 Related Work

Systems with multiple consistency levels. A number of data stores have combined weak and strong consistency, including several commercial and academic systems that combine eventual and strong consistency [5, 6, 29, 49, 58, 65, 73]. Several academic data stores combined causal and strong consistency [9, 37, 41, 42, 59, 65]. Pileus [65] funnels all updates through a single data center. In the fault-tolerant version of lazy replication [37], causal operations require synchronization between replicas on its critical path. In both cases, causal operations are not highly available, defeating the benefits of causal consistency. Walter [59] restricts causal operations to a specific type and lacks fault tolerance due to the use of two-phase commit across data centers. The remaining works [9, 41, 42] support highly available causal operations, but are not fault tolerant. First, they do not make causal operations uniform on demand to guarantee the liveness of strong operations. Thus, they suffer from the liveness issue we explained in §4 (Figure 2). In addition, these systems do not use fault-tolerant mechanisms even for strong transactions. They guard the use of strong transactions using mechanisms similar to locks; if the lock holder fails before releasing it, no other data center can execute a strong transaction requiring the same lock. This occurs even when the service handing locks is fault-tolerant, as in [42]. Finally, the above systems either do not include mechanisms for partitioning the key space among different machines in a data center or include per-data center centralized services, which limits their scalability (§8.2).

Some group communication systems mix causal and atomic

broadcast [10, 68]. However, these systems do not provide mechanisms for maintaining transactional data consistency.

Several papers have proposed tools that use formal verification technology to ensure that consistency choices do not violate application invariants [9, 31, 34, 40, 51, 52]. Such tools can make it easier for programmers to use our system.

Causal consistency implementations. Our subprotocol for causal consistency belongs to a family of highly scalable protocols that avoid using any centralized components or dependency check messages [3, 23, 60–62]; other alternatives are less scalable [4, 8, 13, 22, 32, 43, 44, 48, 72]. While we base our causal consistency subprotocol on an existing one, Cure [3], we have extended it in nontrivial ways, by integrating mechanisms for tracking uniformity (§5.4) and for transaction forwarding (§5.5). Some of the above protocols [32, 61] use hybrid clocks instead of real time [36] to improve performance with large clock skews; this technique can also be integrated into UNISTORE.

SwiftCloud [72] implements *k-stability* [46], a notion similar to uniformity, to enable client migration. SwiftCloud relies on a single per-data center sequencer, which makes tracking *k-stability* easy, but the data store less scalable. Our protocol is more sophisticated, since we distribute the responsibility of tracking uniformity among the replicas in a data center.

Paxos variants. Several Paxos variants [25, 39, 50, 53] lower the latency by allowing commutative operations to execute at replicas in arbitrary orders. In contrast to them, UNISTORE implements PoR consistency, which allows causal transactions to execute without any synchronization at all.

10 Conclusion

This paper presented UNISTORE, the first fault-tolerant and scalable data store that combines causal and strong consistency. UNISTORE carefully integrates state-of-the-art scalable protocols and extends them in nontrivial ways. To maintain liveness despite data center failures, unlike previous work, UNISTORE commits a strong transaction only when all its causal dependencies are uniform. Our results show that UNISTORE combines causal and strong consistency effectively: 3.7× lower latency on average than a strongly consistent system with 1.2ms latency on average for causal transactions. We expect that the key ideas in UNISTORE will pave the way for practical systems that combine causal and strong consistency.

Acknowledgements. We thank our shepherd, Heming Cui, as well as Gregory Chockler, Vitor Enes, Luís Rodrigues and Marc Shapiro for comments and suggestions. This work was partially supported by an ERC Starting Grant RACCOON, the Juan de la Cierva Formación funding scheme (FJC2018-036528-I), the CCF-Tencent Open Fund (CCF-Tencent RAGR20200124) and the AWS Cloud Credit for Research program.

References

- [1] D. Abadi. Consistency tradeoffs in modern distributed database system design: CAP is only part of the story. *IEEE Computer*, 45(2), 2012.
- [2] M. Ahamad, G. Neiger, J. E. Burns, P. Kohli, and P. W. Hutto. Causal memory: Definitions, implementation, and programming. *Distributed Comput.*, 9(1), 1995.
- [3] D. D. Akkoorath, A. Z. Tomsic, M. Bravo, Z. Li, T. Crain, A. Bieniusa, N. Preguiça, and M. Shapiro. Cure: Strong semantics meets high availability and low latency. In *International Conference on Distributed Computing Systems (ICDCS)*, 2016.
- [4] S. Almeida, J. Leitão, and L. Rodrigues. ChainReaction: A causal+ consistent datastore based on chain replication. In *European Conference on Computer Systems (EuroSys)*, 2013.
- [5] Amazon. Read consistency. <https://docs.aws.amazon.com/amazondynamodb/latest/developerguide/HowItWorks.ReadConsistency.html>, 2020.
- [6] Apache Cassandra. Read repair. https://cassandra.apache.org/doc/latest/operating/read_repair.html, 2020.
- [7] H. Attiya, F. Ellen, and A. Morrison. Limitations of highly-available eventually-consistent data stores. *IEEE Trans. Parallel Distributed Syst.*, 28(1), 2017.
- [8] P. Bailis, A. Ghodsi, J. M. Hellerstein, and I. Stoica. Bolt-on causal consistency. In *International Conference on Management of Data (SIGMOD)*, 2013.
- [9] V. Balesar, N. Preguiça, R. Rodrigues, S. Duarte, C. Ferreira, M. Najafzadeh, and M. Shapiro. Putting the consistency back into eventual consistency. In *European Conference on Computer Systems (EuroSys)*, 2015.
- [10] K. Birman, A. Schiper, and P. Stephenson. Lightweight causal and atomic group multicast. *ACM Trans. Comput. Syst.*, 9(3), 1991.
- [11] K. P. Birman and T. A. Joseph. Reliable communication in the presence of failures. *ACM Trans. Comput. Syst.*, 5(1), 1987.
- [12] M. Bravo, A. Gotsman, B. de Régil, and H. Wei. UniStore: A fault-tolerant marriage of causal and strong consistency (extended version). *arXiv CoRR*, abs/2106.00344, 2021.
- [13] M. Bravo, L. Rodrigues, and P. Van Roy. Saturn: A distributed metadata service for causal consistency. In *European Conference on Computer Systems (EuroSys)*, 2017.
- [14] S. Burckhardt. *Principles of Eventual Consistency*. Now Publishers, 2014.
- [15] S. Burckhardt, A. Gotsman, H. Yang, and M. Zawirski. Replicated data types: specification, verification, optimality. In *Symposium on Principles of Programming Languages (POPL)*, 2014.
- [16] C. Cachin, R. Guerraoui, and L. E. T. Rodrigues. *Introduction to Reliable and Secure Distributed Programming (2nd ed.)*. Springer, 2011.
- [17] A. Cerone, G. Bernardi, and A. Gotsman. A framework for transactional consistency models with atomic visibility. In *International Conference on Concurrency Theory (CONCUR)*, 2015.
- [18] Y. L. Chen, S. Mu, J. Li, C. Huang, J. Li, A. Ogus, and D. Phillips. Giza: Erasure coding objects across global data centers. In *USENIX Annual Technical Conference (USENIX ATC)*, 2017.
- [19] G. Chockler and A. Gotsman. Multi-shot distributed transaction commit. In *Symposium on Distributed Computing (DISC)*, 2018.
- [20] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, W. C. Hsieh, S. Kanthak, E. Kogan, H. Li, A. Lloyd, S. Melnik, D. Mwaura, D. Nagle, S. Quinlan, R. Rao, L. Rolig, Y. Saito, M. Szymaniak, C. Taylor, R. Wang, and D. Woodford. Spanner: Google’s Globally-Distributed Database. In *Symposium on Operating Systems Design and Implementation (OSDI)*, 2012.
- [21] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Voshall, and W. Vogels. Dynamo: Amazon’s highly available key-value store. In *Symposium on Operating Systems Principles (SOSP)*, 2007.
- [22] J. Du, S. Elnikety, A. Roy, and W. Zwaenepoel. Orbe: Scalable causal consistency using dependency matrices and physical clocks. In *Symposium on Cloud Computing (SoCC)*, 2013.
- [23] J. Du, C. Iorgulescu, A. Roy, and W. Zwaenepoel. Gentlerain: Cheap and scalable causal consistency with physical clocks. In *Symposium on Cloud Computing (SoCC)*, 2014.
- [24] C. Dwork, N. Lynch, and L. Stockmeyer. Consensus in the Presence of Partial Synchrony. *Journal of the ACM*, 35(2), 1988.

- [25] V. Enes, C. Baquero, T. F. Rezende, A. Gotsman, M. Perin, and P. Sutra. State-machine replication for planet-scale systems. In *European Conference on Computer Systems (EuroSys)*, 2020.
- [26] FaunaDB. What is FaunaDB? <https://docs.fauna.com/fauna/current/introduction.html>, 2020.
- [27] C. Fidge. Timestamps in message-passing systems that preserve the partial ordering. In *Australian Computer Science Conference (ASCS)*, 1988.
- [28] S. Gilbert and N. Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *SIGACT News*, 33(2), 2002.
- [29] Google. Balancing strong and eventual consistency with datastore. <https://cloud.google.com/datastore/docs/articles/balancing-strong-and-eventual-consistency-with-google-cloud-datastore>, 2020.
- [30] A. Gotsman, A. Lefort, and G. Chockler. White-box atomic multicast. In *International Conference on Dependable Systems and Networks (DSN)*, 2019.
- [31] A. Gotsman, H. Yang, C. Ferreira, M. Najafzadeh, and M. Shapiro. ‘Cause I’m strong enough: reasoning about consistency choices in distributed systems. In *Symposium on Principles of Programming Languages (POPL)*, 2016.
- [32] C. Gunawardhana, M. Bravo, and L. Rodrigues. Unobtrusive deferred update stabilization for efficient geo-replication. In *USENIX Annual Technical Conference (USENIX ATC)*, 2017.
- [33] M. P. Herlihy and J. M. Wing. Linearizability: A correctness condition for concurrent objects. *ACM Trans. Program. Lang. Syst.*, 12(3), 1990.
- [34] G. Kaki, K. Earanky, K. C. Sivaramakrishnan, and S. Jagannathan. Safe replication through bounded concurrency verification. *Proc. ACM Program. Lang.*, 2(OOPSLA), 2018.
- [35] T. Kraska, G. Pang, M. J. Franklin, S. Madden, and A. Fekete. MDCC: Multi-data center consistency. In *European Conference on Computer Systems (EuroSys)*, 2013.
- [36] S. S. Kulkarni, M. Demirbas, D. Madappa, B. Avva, and M. Leone. Logical physical clocks. In *International Conference on Principles of Distributed Systems (OPODIS)*, 2014.
- [37] R. Ladin, B. Liskov, L. Shrira, and S. Ghemawat. Providing high availability using lazy replication. *ACM Trans. Comput. Syst.*, 10(4), 1992.
- [38] L. Lamport. The part-time parliament. *ACM Trans. Comput. Syst.*, 16(2), 1998.
- [39] L. Lamport. Generalized consensus and Paxos. Technical Report MSR-TR-2005-33, Microsoft Research, 2005.
- [40] C. Li, J. Leitão, A. Clement, N. M. Preguiça, R. Rodrigues, and V. Vafeiadis. Automating the choice of consistency levels in replicated systems. In *USENIX Annual Technical Conference (USENIX ATC)*, 2014.
- [41] C. Li, D. Porto, A. Clement, R. Rodrigues, N. Preguiça, and J. Gehrke. Making geo-replicated systems fast if possible, consistent when necessary. In *Symposium on Operating Systems Design and Implementation (OSDI)*, 2012.
- [42] C. Li, N. Preguiça, and R. Rodrigues. Fine-grained consistency for geo-replicated systems. In *USENIX Annual Technical Conference (USENIX ATC)*, 2018.
- [43] W. Lloyd, M. J. Freedman, M. Kaminsky, and D. G. Andersen. Don’t settle for eventual: Scalable causal consistency for wide-area storage with cops. In *Symposium on Operating Systems Principles (SOSP)*, 2011.
- [44] W. Lloyd, M. J. Freedman, M. Kaminsky, and D. G. Andersen. Stronger semantics for low-latency geo-replicated storage. In *Conference on Networked Systems Design and Implementation (NSDI)*, 2013.
- [45] P. Mahajan, L. Alvisi, and M. Dahlin. Consistency, availability, and convergence. Technical Report TR-11-22, University of Texas at Austin, 2011.
- [46] P. Mahajan, S. Setty, S. Lee, A. Clement, L. Alvisi, M. Dahlin, and M. Walfish. Depot: Cloud storage with minimal trust. *ACM Trans. Comput. Syst.*, 29(4), 2011.
- [47] F. Mattern. Virtual time and global clocks in distributed systems. In *International Workshop on Parallel and Distributed Algorithms*, 1988.
- [48] S. A. Mehdi, C. Littlely, N. Crooks, L. Alvisi, N. Bronson, and W. Lloyd. I can’t believe it’s not causal! Scalable causal consistency with no slowdown cascades. In *Conference on Networked Systems Design and Implementation (NSDI)*, 2017.
- [49] Microsoft. Consistency levels in Azure Cosmos DB. <https://docs.microsoft.com/en-us/azure/cosmos-db/consistency-levels>, 2020.

- [50] I. Moraru, D. G. Andersen, and M. Kaminsky. There is more consensus in egalitarian parliaments. In *Symposium on Operating Systems Principles (SOSP)*, 2013.
- [51] S. S. Nair, G. Petri, and M. Shapiro. Proving the safety of highly-available distributed objects. In *European Symposium on Programming (ESOP)*, 2020.
- [52] M. Najafzadeh, A. Gotsman, H. Yang, C. Ferreira, and M. Shapiro. The CISE tool: proving weakly-consistent applications correct. In *Workshop on the Principles and Practice of Consistency for Distributed Data (PaPoC)*, 2016.
- [53] F. Pedone and A. Schiper. Generic broadcast. In *International Symposium on Distributed Computing (DISC)*, 1999.
- [54] K. Petersen, M. J. Spreitzer, D. B. Terry, M. M. Theimer, and A. J. Demers. Flexible update propagation for weakly consistent replication. In *Symposium on Operating Systems Principles (SOSP)*, 1997.
- [55] Redis Labs. Causal consistency in an active-active database. <https://docs.redislabs.com/latest/rs/administering/database-operations/causal-consistency-crdb/>, 2020.
- [56] F. B. Schneider. Implementing fault-tolerant services using the state machine approach: A tutorial. *ACM Comput. Surv.*, 22(4), 1990.
- [57] M. Shapiro, N. M. Preguiça, C. Baquero, and M. Zawirski. Conflict-free replicated data types. In *International Symposium on Stabilization, Safety, and Security of Distributed Systems (SSS)*, 2011.
- [58] D. Shukla, S. Thota, K. Raman, M. Gajendran, A. Shah, S. Ziuzin, K. Sundaram, M. G. Guajardo, A. Wawrzyniak, S. Boshra, R. Ferreira, M. Nassar, M. Koltachev, J. Huang, S. Sengupta, J. J. Levandoski, and D. B. Lomet. Schema-agnostic indexing with azure documentdb. *Proc. VLDB Endow.*, 8(12), 2015.
- [59] Y. Sovran, R. Power, M. K. Aguilera, and J. Li. Transactional storage for geo-replicated systems. In *Symposium on Operating Systems Principles (SOSP)*, 2011.
- [60] K. Spirovska, D. Didona, and W. Zwaenepoel. Wren: Nonblocking reads in a partitioned transactional causally consistent data store. In *International Conference on Dependable Systems and Networks (DSN)*, 2018.
- [61] K. Spirovska, D. Didona, and W. Zwaenepoel. Paris: Causally consistent transactions with non-blocking reads and partial replication. In *International Conference on Distributed Computing Systems (ICDCS)*, 2019.
- [62] K. Spirovska, D. Didona, and W. Zwaenepoel. Optimistic causal consistency for geo-replicated key-value stores. *IEEE Transactions on Parallel and Distributed Systems*, 32(3), 2020.
- [63] R. Taft, I. Sharif, A. Matei, N. VanBenschoten, J. Lewis, T. Grieger, K. Niemi, A. Woods, A. Birzin, R. Poss, P. Bardea, A. Ranade, B. Darnell, B. Gruneir, J. Jaffray, L. Zhang, and P. Mattis. CockroachDB: The Resilient Geo-Distributed SQL Database. In *International Conference on Management of Data (SIGMOD)*, 2020.
- [64] D. B. Terry, A. J. Demers, K. Petersen, M. Spreitzer, M. Theimer, and B. W. Welch. Session guarantees for weakly consistent replicated data. In *International Conference on Parallel and Distributed Information Systems (PDIS)*, 1994.
- [65] D. B. Terry, V. Prabhakaran, R. Kotla, M. Balakrishnan, M. K. Aguilera, and H. Abu-Libdeh. Consistency-based service level agreements for cloud storage. In *Symposium on Operating Systems Principles (SOSP)*, 2013.
- [66] 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 *Symposium on Operating Systems Principles (SOSP)*, 1995.
- [67] M. Tyulenev, A. Schwerin, A. Kamsky, R. Tan, A. Cabral, and J. Mulrow. Implementation of cluster-wide logical clock and causal consistency in MongoDB. In *International Conference on Management of Data (SIGMOD)*, 2019.
- [68] R. van Renesse, K. P. Birman, and S. Maffei. Horus: A flexible group communication system. *Commun. ACM*, 39(4), 1996.
- [69] W. Vogels. Eventually consistent. *CACM*, 52(1), 2009.
- [70] G. Weikum and G. Vossen. *Transactional Information Systems: Theory, Algorithms, and the Practice of Concurrency Control and Recovery*. Morgan Kaufmann Publishers Inc., 2001.
- [71] YugabyteDB. Replication. <https://docs.yugabyte.com/latest/architecture/docdb-replication/replication/>, 2020.
- [72] 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 *International Middleware Conference (Middleware)*, 2015.
- [73] I. Zhang, N. K. Sharma, A. Szekeres, A. Krishnamurthy, and D. R. K. Ports. Building consistent transactions with inconsistent replication. *ACM Trans. Comput. Syst.*, 35(4), 2018.