

# Rigorous design and implementation of full-featured database backends

SAALIK HATIA, Sorbonne-Université (LIP6), France

ANNETTE BIENIUSA, TU Kaiserslautern, Germany

GUSTAVO PETRI, ARM, United Kingdom

CARLA FERREIRA, Universidade NOVA de Lisboa, Portugal

MARC SHAPIRO, Sorbonne-Université (LIP6) and Inria, France

Conceptually, a database backend is just a memory for storing and retrieving data. However, the quest for performance and reliability adds significant complexity. To ensure correctness, we derive our implementation from a formal specification. We start from the formalisation of a basic transactional API. We implement this specification in a verbatim, unoptimised persistent *map-based store*. We also formalise a journal-based variant, prove equivalence with the map semantics, and implement an unoptimised crash-tolerant *journal-based store*. Finally, we study the correct *composition* of stores. We leverage compositionality to create variants of increasing complexity: Features such as caching, checkpointing, and write-ahead logging are each represented as an instance of a map or journal store, restricted to a suitable domain, which is composed with another, larger store. This paper reports on the formal specification, the approach to correctness of specifications and implementation, and a preliminary experimental evaluation. This is work in progress.

## CONTENTS

Abstract	1
Contents	1
1 Introduction	3
2 Approach	4
3 System model	6
3.1 Store	6
3.1.1 Keys	6
3.1.2 Values	6
3.1.3 Effects	7
3.1.4 Timestamps and clocks	7
3.1.5 Mappings	7
3.2 Data types and concurrent effects	8
3.3 Transactions	8
3.4 Visibility	9
3.5 Fault models	9
4 Formal model of transactions	9
4.1 Composing effects	9
4.2 Semantics of transactions	11
4.2.1 Informal presentation	11
4.2.2 Parameters	11
4.2.3 Transaction begin	12
4.2.4 Reads and writes	12
4.2.5 Transaction termination	12
4.2.6 Correctness and performance challenges	13

4.3	Traces and histories	14
4.4	Consistency	14
4.5	**** TODO ****	14
5	▷(TODO:) Todo◁ General Lemmas	14
5.1	General results	14
6	Map-based store	15
6.1	Map store semantics	15
6.2	Map store implementation	17
6.3	Checkpoints	17
7	Journal store	17
7.1	▷(TODO:) TODO◁ Journal store proofs	18
7.2	Journal store implementation	18
8	Design challenges and implementation challenge	19
9	Composing stores	20
9.1	Domain of a store	21
9.2	Composition of stores	22
9.3	Top-level (composed) store	23
9.4	Changing composition	23
9.5	Garbage collection	23
9.6	**** lost text ****	23
9.7	Proof goals	23
9.8	Compositon with the Conductor	23
9.9	Implementation of a simple composed store (Saalik)	24
9.10	Bounded Stores (Gustavo, Annette, Carla)	25
9.11	Moving the bounds (Gustavo, Annette, Carla)	25
9.12	Checkpointing and truncation (Saalik)	25
10	Implementation and experimental results (Saalik)	26
11	Related work (everyone)	27
11.1	Compiling specifications to executable code	27
11.2	Verification of complex storage software	27
11.3	Formal specification of transactions and isolation models	27
12	Lessons learned and future work (Marc)	27
	Acknowledgments	27
	References	27

## STORY FOR PL PAPER

Title “Rigorous design by composition: a formally-verified transactional database backend design.”

- DB backends are inherently complex.
- Simpler model by composition of primitive backends. Clarify/justify common features/optimisations: caching, journaling, checkpointing, layered stores, etc.
- Operational semantics of transactions and backend components.
- Operational model of composition.
- Show correctness of widely-used optimisations.

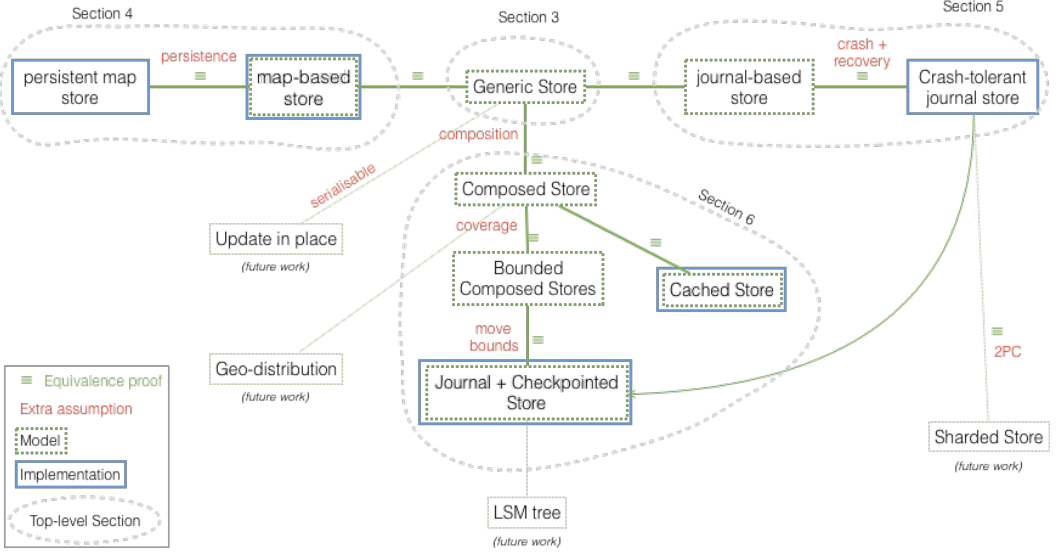


Fig. 1. Structure of the approach

## 1 INTRODUCTION

A database system manages a collection of digital data, to be organised, secured, updated, queried, and so on. An essential component is the *backend*, in charge of recording the lowest level of data into some memory or *store*. Although conceptually simple at a high level, actual backends are complex, due to the demands for fast response, high volume, limited footprint, concurrency, distribution, reliability, and so on. For instance, the open-source RocksDB has \*XXX\* LOC, Redis is \*XXX\* LOC, and the AntidoteDB backend is 1300 LOC [? ? ?]. Any such complex software has bugs; database backend bugs are critical, possibly violating data integrity or security. >(Marc:) Cite some episodes?<

To ensure correctness, formally-verified development is a powerful tool, but is often perceived as at odds with the performance requirements of backends. To our knowledge, formally-verified development of a serious backend has not been attempted before. One obstacle to the formal approach is the very complexity of a modern backend: fully specifying all the moving pieces is a daunting task.

This article reports an incremental approach to the rigorous and modular development of a full-featured backend. We first formalise the simplest specification of a store, a versioned key-value interface; we implement it directly, without any optimisation. As this specification makes versioning explicit, the state is a deterministic function; we show that the precise order of execution does not matter, and that multiple instances of store are mutually equivalent. We next consider a journal-based variant, again in its simplest form; we formally prove that it is equivalent with the base variant; and we implement it directly without optimisation. To support common optimisations such as caching, garbage collection, checkpointing and write-ahead logging, we proceed by *composition*, using the insight that each such optimisation can be represented as one of the basic stores, restricted to a suitable domain.

Our work combines a formal operational semantics (on paper), mechanised formal proofs (in Coq), and implementation (in Java).

In future work, we believe our same approach can justify sharding, geo-distribution, layered checkpoints in the style of LSM-Trees, and update-in-place under serialisability.

Our operational semantics specifies the rules for transactions under convergent causal consistency (TCC+), a variant of PSI [18]. In the style of MVCC (multi-value concurrency control), a transaction reads from a causally-consistent snapshot. It terminates by either committing atomically, or by aborting without effect. Our results generalise to stronger models such as SI or strict serialisability (which differ only by the conditions on transaction begin and commit).

The semantics specify the behaviour of a store’s specialised book-keeping operations (called `doUpdate`, `doCommit` and `lookup`), implicitly assuming infinite memory and no failures. To bound a store’s memory footprint, we restrict its domain, and show equivalence over the restricted domain. Furthermore we exhibit a (classical) crash-tolerant journal implementation, which satisfies the failure-free semantics by masking crash-recovery events.

Tools exist to compile a Coq specification to executable code [?]. We did not take that path for pragmatic reasons, because it is too far from the ordinary programmer’s experience. Instead, we manually transcribe the specification to Java, verbatim, resisting the temptation to optimise. We check through testing that it behaves like the specification, and that variants that were proved equivalent do have the same runtime behaviour. We report on the implementation throughout the paper, and an experimental evaluation in Section 10.

This paper focuses on safety properties, and does not consider security issues.

Our contributions can be summarised as follows:

- A formal model for backends; proof of equivalence between concurrent and sequential execution; correctness conditions for store composition.
- Specification, verification and implementation of variants supporting state-of-the-art features such as journaling, caching, crash-tolerant write-ahead logging and journal truncation.
- Experimental evaluation, showing that the rigorous approach does not preclude state-of-the-art performance, and justifying retrospectively the features.

▷(TODO:) Update me.◁ This paper progresses as illustrated in Figure 1. After this introduction, Section 3 formalises the system model and semantics. Section 2 explains our approach, and establishes some key lemmas. As a warm-up, we produce a correct map-based persistent store in Section ?? . Then, Section ?? addresses correct journal-based semantics and its implementation. In Section ?? we apply store composition to achieve caching, checkpointing, and truncating the journal. Our experimental evaluation is in Section 10. Section 11 discusses related work, and we discuss lessons learned in Section 12.

## 2 APPROACH

This work proceeds as follows. We first model two basic variants. The first variant is a classical map-based, versioned key-value store. As a transaction executes, it *eagerly* computes new key versions, which it copies into a map upon commit, labelled with the transaction’s commit timestamp. Reading a key at some timestamp searches the map for the most recent corresponding version. We mechanise the proof that the map-based model satisfies the specification, i.e., that for any given trace  $\Theta$ , a call to `read(k, t)` returns the expected value noted `safe(Θ, k, t)` for any key and timestamp pair. We translate the formal specification verbatim into Java, forsaking optimisations, and test empirically that this implementation satisfies the specification.

Our second variant uses a journal. As a transaction executes, it appends individual effects to the journal, tagged with the transaction identifier. Committing appends a commit record to the

$k$	$\in \text{Key}$	Primitive Key
$t$	$\in \text{TS}$	Primitive Timestamp
$v$	$\in \text{Value}$	Primitive data-type values
$\sigma$	$\in \Sigma = \text{Key} \times \text{TS} \rightarrow \mathcal{E}\mathcal{f}\mathcal{f}_\perp$	Store
$\mathcal{D}_\sigma$	$\subseteq \text{Key} \times \text{TS}$	Domain of store $\sigma$
$\mathcal{S}$	$\{\sigma_1, \sigma_2, \dots\}$	Composition of stores
$\delta$	$\in \mathcal{E}\mathcal{f}\mathcal{f} = \text{Value}_\perp \rightarrow \text{Value}$	Effect
$\perp$	$\notin \mathcal{E}\mathcal{f}\mathcal{f}$	Absence of effect
$\delta_k^t$	$\in \mathcal{T}\mathcal{E}\mathcal{f}\mathcal{f} = \mathcal{E}\mathcal{f}\mathcal{f} \times \text{Key} \times \text{TS}$	Effect to key $k$ with commit timestamp $t$
$\delta < \delta'$		Visibility
$\Theta$	$\in (\mathcal{T}\mathcal{E}\mathcal{f}\mathcal{f}, <)$	Trace
$\tau$	$\in \mathcal{T}_D$	Transaction identifier
$\text{dt}$	$\in \text{TS}$	Dependency timestamp (of transaction)
$\mathcal{E}$	$\in \mathcal{P}(\text{Key})$	Dirty set = set of keys modified in transaction
$\mathcal{R}$	$\in \text{Rbuf} = \text{Key} \rightarrow \text{Value}_\perp$	Read buffer (of transaction)
$\text{ct}$	$\in \text{TS}$	Commit timestamp (of transaction)
$(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, \text{ct})$	$\in \mathcal{T}_{\text{DESC}} = \mathcal{T}_D \times \text{TS} \times \mathcal{P}(\text{Key}) \times \text{Rbuf} \times \text{TS}$	Transaction descriptor
$\mathcal{T}_a$	$\subseteq \mathcal{T}_{\text{DESC}}$	Aborted transactions
$\mathcal{T}_c$	$\subseteq \mathcal{T}_{\text{DESC}}$	Committed transactions
$\mathcal{T}_r$	$\subseteq \mathcal{T}_{\text{DESC}}$	Running transactions
$\Pi_x \mathcal{T}$		Project set of tuples $\mathcal{T}$ along dimension $x$
$\sigma, \mathcal{T}_r, \dots$		Before transition
$\sigma', \mathcal{T}_r', \dots$		After transition

Table 1. Notation

journal, containing its commit timestamp.  $\text{read}(k, t)$  *lazily* applies the effects to  $k$ , recorded in the journal, that committed before  $t$ .

Again, we prove mechanically that the journal-based store satisfies the abstract specification. We transcribe the specification to a Java crash-tolerant journal store, and show empirically that it satisfies the specification.

The above implementations are impractical, because they are slow, consume a lot of memory, and are not fault tolerant. Features such as caching, sharding, replication, checkpointing, journal truncation, etc., are known to improve performance. Our key insight is the realisation that each such feature can be viewed as a composition of “little stores” that each implement the specification over some limited domain. For instance, we compose a truncated journal with a map-based checkpoint at the time of truncation.

We formulate rules for the composition of little stores. The domain of a composition is the union of the composed domains. Reading composes the results from all the little stores that are in domain, and writing updates all the stores in domain. Domains may overlap, as the little stores cannot contain conflicting information; furthermore, a useful composition must not have gaps, i.e., all composed reads must be safe.

A first application is caching. Consider some slow persistent store, either map- or journal-based, with a large domain  $K \times T$ . We compose it with a cache, an in-memory map-based store, whose domain is some small, arbitrary subset  $\{(k_1, t_1), (k_2, t_2), \dots\} \subset K \times T$ . For any key-version that is in the cache, both the cache and the larger store would read the same, safe, value, but the cache responds more quickly. Serving a cache miss grows the cache domain; eviction shrinks it

arbitrarily. We reuse the unoptimised Java map-based store, equipped with an associated domain. Unsurprisingly, our experiments confirm that caching improves response time; our contribution is achieving this through rigorous programming discipline.

Our second application is to split the store between an “old generation” over time domain  $[0, \max\_old)$  and a “new generation” over  $[\min\_new, \text{now}())$ , where  $\min\_new \leq \max\_old$ . As above, any key-version in the overlap will be the same in both stores. The bounds of generations may change, as long as the invariant  $\min\_new \leq \max\_old$  holds, for instance by first increasing  $\max\_old$  by some amount, before increasing  $\min\_new$  by the same amount. A practical use case is a persistent map-based store as the old generation and a journal-based store as the new one. Advancing  $\max\_old$  corresponds to checkpointing, and advancing  $\min\_new$  to truncating the journal. Again, we combine the corresponding unoptimised Java implementations, and test that it satisfies the specification. We leave it as an exercise for the reader to address multiple levels of checkpointing, and to compose the split store with a cache, for performance.

In future work, we believe that the same theory of composition can justify sharding, local replicas at multiple sites, update-in-place (under strict serialisability), and layering in the style of LSM-Trees.

### 3 SYSTEM MODEL

►(Marc 2023-01-07:) Give concrete examples◄

Table 1 summarises our notations, which are explained hereafter.

This study envisages different kinds of backend stores under a single system model. They follow a common API, described later, consisting of query method lookup and update methods `doUpdate` and `doCommit`.

A store runs under a general model of transactions, with Multi-Version Concurrency Control (MVCC). The transaction model guarantees Transactional Causal Plus Consistency (TCC+) [1, 19], also called Non-Monotonic Snapshot Isolation (NMSI) [3, 15]. Our results remain true for stronger consistency/isolation models, e.g., Parallel Snapshot Isolation, Snapshot Isolation, or Strict Serialisability; however this is mostly out of the scope of this paper.

#### 3.1 Store

We define an abstract model of a backend storage system (hereafter abbreviated to just *store*) with a versioned key-value interface. A store is denoted through the type  $\Sigma$ , and ranged over by the meta-variable  $\sigma$ , which is a function from key ( $k$ ) and timestamp ( $t$ ) pairs to (a means to compute) a value ( $v$ ). In the case where either the key or the timestamp is not defined in the store,  $\perp \notin \text{Value}$  represents the absence of a value.

Its API `lookup( $\sigma, k, t$ )` returns the value that store  $\sigma$  associates with key  $k$  at time  $t$ .

**3.1.1 Keys.** Keys are denoted by the set *Key*, for which we use meta-variable  $k$ . Keys are opaque, and are compared only for equality.

**3.1.2 Values.** The set of values is denoted by *Value* and ranged over with the meta-variable  $v$ . For the purpose of this paper, the type of values is a collection of data types, with some constraints on semantics that are detailed later.

**3.1.3 Effects.** Classically, a store records values. We take a more general approach, where the store records operations that compute a value, called *effects*. We use the meta-variable  $\delta$  to denote an effect, and the domain of effects  $\mathcal{Eff}$ .

We define an effect to be a function from value to value. An example is “inc\_10(arg),” which returns its argument plus 10. *Assignment* is modelled by the constant function; for example “assign\_27( )” is the effect that always returns 27. Abusing notation, we sometimes identify an assignment, e.g., and note simply 27 rather than “assign\_27( ).”

**3.1.4 Timestamps and clocks.** *Timestamps* serve to identify events and to distinguish the versions of a key. We assume the following properties about timestamps:

- *Partial Order*: the set of timestamps (TS) defines a Partial Order, denoted  $\leq$ .
- *Lattice*: for any two timestamps, there exists for them a unique least upper bound and a unique greatest lower bound.

In general, except when mentioned specifically, we do not assume that timestamps are totally ordered. Two timestamps that are not mutually ordered, i.e.,  $t_1 \not\leq t_2 \wedge t_2 \not\leq t_1$ , are said concurrent, noted  $t_1 \parallel t_2$ .

The type for timestamps is denoted by TS and generically instantiated by the meta-variable  $t$ . We use the following notation:

- $\max(t_1, t_2) \triangleq t_1 \sqcup t_2$  is the least upper bound of timestamps.
- We generalise  $\max$  to a set in the obvious way. As a shorthand,  $\max_{\mathcal{T}}(dt) \triangleq \max(\Pi_{dt}(\mathcal{T}))$  is the least-upper-bound of the projection of the set of tuples  $\mathcal{T}$  over the  $dt$  dimension; and similarly for the  $ct$  dimension.
- $t_1 < t_2 \triangleq t_1 \leq t_2 \wedge t_1 \neq t_2$ .
- $t_1 \geq t_2 \triangleq t_1 \not< t_2$ .  $\geq$  is read “greater, equal, or concurrent.”

The notations  $\min$ ,  $>$  and  $\leq$  are defined symmetrically.

A vector of integers is a popular concrete representation for timestamps [7, 13]. In this case, for indexes  $i, j$  ranging over such a vector:

- $t_1 \leq t_2 \iff \forall i : t_1[i] \leq t_2[i]$
- $t_1 < t_2 \iff t_1 \leq t_2 \wedge \exists i : t_1[i] \neq t_2[i]$
- $\forall i : \max(t_1, t_2)[i] = \max(t_1[i], t_2[i])$
- $t_1 \parallel t_2 \iff \exists i, j : t_1[i] < t_2[i] \wedge t_1[j] > t_2[j]$
- and so on.

A *clock* is an object that returns unique and monotonically-increasing timestamps. Concretely, any participant in the system may have its local clock, noted  $now()$ , such that only events timestamped  $\leq now()$  are visible to it; if it attempts an operation with a larger timestamp, the operation will not return until  $now()$  increases sufficiently. We do not assume a global clock.

**3.1.5 Mappings.** In the spirit of MVCC, to update a key’s mapping creates a new version, without overwriting an existing mapping; a version is write-once. A version is valid for the range of timestamps starting from its commit timestamp, until the next version, if any.

A store records mappings of a key and a timestamp to an effect. Updating a key  $k$  under timestamp  $t$  creates a new version tagged with  $t$ , without overwriting any existing versions.

For example, suppose store  $\sigma$  records effect 27 in version  $t_1$  of key  $k$ . This ensures that  $\text{lookup}(\sigma, k, t_1)$  will return 27. To increment  $k$  by 10 it might record a new version containing effect

“inc\_10” under the same key, with timestamp  $t_2 > t_1$ . If there are no versions between  $t_1$  and  $t_2$ , then,  $\text{lookup}(\sigma, k, t_2)$  will return 37.

Note that an assignment masks any earlier versions of the same key.

### 3.2 Data types and concurrent effects

We distinguish two kinds of value data types: classical vs. RDTs.

*Classical data types.* A classical data type is one whose only kind of effect is assignment, and does not allow assignments with concurrent timestamps. A standard sequential data type, such as integer or struct, is classical; so is a linearisable data type. A classical data type implies strong consistency, defined as totally-ordered timestamps. The merge operator discussed later is not defined for classical data types.

*Replicated Data Types (RDTs).* Replicated Data Types (RDTs) support more general and concurrent effects, by supporting a commutative, associative and idempotent merge operator [2, 17]. The current value of an RDT key results from applying the effects that follow the most recent assignment, and merging concurrent effects.

One example is a counter equipped with increment and decrement operations; as increment and decrement commute, merging simply composes them in any order. A classical data type augmented with a suitable merge becomes an RDT [10, 16]. A simple example is the last-writer-wins register, that merges two values by choosing the one with the larger timestamp, deterministically.

*Assumptions about data type semantics.* We make the following assumptions, which are true for both classical data types and RDTs:

- Effects are associative, i.e.,  $(\delta'' \circ \delta') \circ \delta = \delta'' \circ (\delta' \circ \delta)$ , where  $\circ$  denotes functional composition.
- Effects to different keys commute, i.e.,  $k \neq k' \implies \delta'_{k'} \circ \delta_k = \delta_k \circ \delta'_{k'}$ .
- Where  $\text{merge}(\delta, \delta')$  is defined, its arguments commute, i.e.,  $\text{merge}(\delta, \delta') = \text{merge}(\delta', \delta)$ .

We assume that all versions of a same key have the same data type.

### 3.3 Transactions

A transaction is a sequence of effects, with the following properties. All the transaction’s reads come from a same *snapshot*. A running transaction is *isolated*, i.e., its effects are not visible outside of it. A transaction terminates in an all-or-nothing manner: it either *aborts*, and its effects are discarded, making no changes to the store; or it *commits atomically*, and its effects become visible in the store all at once.

An effect  $\delta$  is said committed if it belongs to a committed transaction, and uncommitted otherwise. When a transaction commits, it is assigned a *commit timestamp*, noted  $ct$ . Notation  $\delta_k^{ct}$  refers to effect  $\delta$  updating key  $k$  and committed with timestamp  $ct$ .

When it starts, a transaction has a *snapshot timestamp* (a.k.a. dependency timestamp) noted  $dt$ . A transaction’s snapshot  $dt$  has visibility of an effect if and only if the latter’s commit timestamp  $ct$  is less or equal  $dt$ .

\*\*\*\* GIVE CONCRETE EXAMPLE



### 3.4 Visibility

Effects are ordered by the *visibility* relation  $\delta_1 < \delta_2$  (read “ $\delta_1$  is visible to  $\delta_2$ ”) defined as follows:

**(External visibility EXT)**  $\delta < \delta'$  if they belong to different transactions,  $\tau$  and  $\tau'$  respectively, where  $\tau'$  can read from  $\tau$ ; formally,  $\tau$  has committed, and  $\tau.ct < \tau'.dt$ .

**(Internal visibility INT)**  $\delta < \delta'$  if both execute in the same transaction, and  $\delta'$  executes before  $\delta$ .

►(TODO:) In Section ??, we will prove that visibility is also transitive.◄

Visibility is a partial order.<sup>1</sup> Committed effects that are not ordered are said *concurrent*, noted  $\parallel$ . Note that an uncommitted effect is *not* considered concurrent to another transaction, since it is not visible outside of its transaction.

### 3.5 Fault models

Failures are non-byzantine stop-restart. The system can stop due to a clean power-down or to a crash. Memory supports blocking reads and non-blocking writes. We consider three kinds of memory:

- Volatile memory (the default) models DRAM. Volatile memory is lost on restart.
- Persistent memory models secondary storage. Before a power-down, in-progress writes first terminate. If the system crashes, in-progress writes may take effect or not, non-deterministically.
- Crash-tolerant memory models an append-mode log file. Writes are sequentially ordered. If some write succeeds, all preceding writes will have succeeded. Furthermore, it supports a blocking *flush* operation. If *flush* returns successfully, or is followed by a power-down, all preceding writes will have succeeded. If a crash occurs before *flush* returns, it is guaranteed that some prefix of the preceding writes terminated successfully.

## 4 FORMAL MODEL OF TRANSACTIONS

The semantics of transactions is generic and independent of store variant. We first specify transactions formally, then study the implementation challenges.

### 4.1 Composing effects

The current state of a key is the result of applying the effects to the key in visibility order, and merging concurrent ones. We formalise this intuition with the *effect composition* rule in Figure 4. Here,  $\circ$  denotes functional composition:  $\forall v, (\delta_2 \circ \delta_1)(v) \triangleq \delta_2(\delta_1(v))$ . Recall that  $\circ$  binds right-to-left.

In detail: The first two lines state that the null effect  $\perp$  (corresponding for instance to a non-initialised key) can be ignored. The third line ensures that the same effect is not applied in duplicate. The fourth line states that keys are independent: effects to different keys commute (they can be applied in any order). The next two lines apply effects that are related by visibility in visibility order. The penultimate one ensures that concurrent effects are merged. Any other case would compose a non-committed effect with an effect that is not visible, and would therefore be an error; hence the final line.

The  $\sqcup$  operator is associative, commutative and idempotent. It defines a least upper bound between effects, in a join-semilattice whose least element is  $\perp$ .

<sup>1</sup>Note that even if timestamps form a total order (strong consistency), it is still possible to have concurrent effects. Indeed, a consistency models such as Snapshot Isolation allows the interval between snapshot and commit of committed transactions to overlap.

$$\begin{array}{c}
\text{BEGIN\_TXN} \\
\frac{\tau \notin \Pi_{\tau}(\mathcal{T}_a \cup \mathcal{T}_c \cup \mathcal{T}_r) \quad \forall t \in \Pi_{\text{ct}}(\mathcal{T}_r) : t \not\prec \text{dt} \quad \mathcal{T}_r' = \mathcal{T}_r \cup \{(\tau, \text{dt}, \emptyset, \emptyset, +\infty)\}}{(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r) \xrightarrow[\tau]{\text{begin}(\text{dt})} (\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r')} \\
\\
\text{INIT\_KEY} \\
\frac{k \notin \text{dom}(\mathcal{R}) \quad \text{lookup}(\sigma, k, \text{dt}) = \delta \quad \mathcal{R}' = \mathcal{R}[k \leftarrow \delta] \quad \mathcal{T}_r' = \mathcal{T}_r'' \cup \{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, +\infty)\}}{(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r) \xrightarrow[\tau]{} (\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r')} \\
\\
\text{READ} \\
\frac{\mathcal{T}_r = \mathcal{T}_r'' \cup \{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, +\infty)\} \quad k \in \text{dom}(\mathcal{R}) \quad \mathcal{R}(k) = v}{(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r) \xrightarrow[\tau]{\text{read}_{\tau}(k) \rightarrow \{v\}} (\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r)} \\
\\
\text{UPDATE} \\
\frac{k \in \mathcal{E} \quad \sigma' = \text{doUpdate}(\sigma, \tau, k, \delta) \quad \mathcal{T}_r' = \mathcal{T}_r'' \cup \{(\tau, \text{dt}, \mathcal{E} \cup k, \mathcal{R}[k \leftarrow \delta(\mathcal{R}(k))], +\infty)\}}{(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r) \xrightarrow[\tau]{\text{update}(\tau, k, \delta)} (\sigma', \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r')} \\
\\
\text{ABORT} \\
\frac{\mathcal{T}_r = \mathcal{T}_r' \cup \{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, +\infty)\} \quad \mathcal{T}_a' = \mathcal{T}_a \cup \{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, +\infty)\}}{(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r) \xrightarrow[\tau]{\text{abort}_{\tau}} (\sigma, \mathcal{T}_a', \mathcal{T}_c, \mathcal{T}_r')} \\
\\
\text{COMMIT} \\
\frac{\text{ct} \notin \Pi_{\text{ct}}(\mathcal{T}_c) \quad \text{dt} \leq \text{ct} \quad \mathcal{T}_r = \mathcal{T}_r' \cup \{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, +\infty)\} \quad \sigma' = \text{doCommit}(\sigma, \tau, \mathcal{E}, \mathcal{R}, \text{ct}) \quad \mathcal{T}_c' = \mathcal{T}_c \cup \{(\tau, \text{dt}, \mathcal{E}, \mathcal{R}, \text{ct})\}}{(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r) \xrightarrow[\tau]{\text{commit}_{\tau}(\text{ct})} (\sigma', \mathcal{T}_a, \mathcal{T}_c', \mathcal{T}_r')} \\
\\
\text{►(Marc 2023-02-08:) Somewhere, add the constraint that } \delta \text{ above must resolve to a Value in the top-level composed store.◄} \\
\\
\text{►(Marc 2023-02-02:) Updated to allow ct=dt◄} \\
\text{►(Marc 2023-02-13:) Moved visibility clause to BEGIN\_TXN◄}
\end{array}$$

Fig. 2. Operational Semantics of Transactions

$$\begin{array}{ll}
\text{lookup} & : \Sigma \times \text{Key} \times \text{TS} \rightarrow \mathcal{E}\mathcal{F}_{\perp} \\
\text{doUpdate} & : \Sigma \times \mathcal{T}_{ID} \times \text{Key} \times \mathcal{E}\mathcal{F} \rightarrow \Sigma \\
\text{doCommit} & : \Sigma \times \mathcal{T}_{ID} \times \mathcal{P}(\text{Key}) \times \text{Rbuf} \times \text{TS} \rightarrow \Sigma
\end{array}$$

Fig. 3. Store interface

The value expected of key  $k$ , at some timestamp  $t$ , results from applying effects to  $k$  up to  $t$ , in visibility order, while merging concurrent effects, i.e.,  $\text{safe}(\Theta, k, t) \triangleq \bigsqcup \{\delta_k' : \delta_k' \in \Theta \wedge t' < t\}$ , where  $\Theta$  is a legal trace of effects, and  $\delta_k'$  is an effect to key  $k$  committed with timestamp  $t'$ .<sup>2</sup>

Note that an assignment overwrites the previous history of the key, i.e.,  $\delta_k < \delta_k' \wedge \delta_k' \in \text{Value} \implies \delta_k \sqcup \delta_k' = \delta_k'$ . More generally, if a key's value is known at some point, earlier updates to

<sup>2</sup>A legal trace is one produced by the semantic rules specified in Section 4.2 hereafter.

$$\delta_k \sqcup \delta'_{k'} = \delta'_{k'} \sqcup \delta_k = \begin{cases} \delta_k & \text{if } \delta'_{k'} = \perp & (\text{null effect}) \\ \delta'_{k'} & \text{if } \delta_k = \perp & (\text{null effect}) \\ \delta_k & \text{if } \delta_k = \delta'_{k'} & (\text{idempotence}) \\ \delta'_{k'} \circ \delta_k & \text{if } k \neq k' & (\text{keys are independent}) \\ \delta'_{k'} \circ \delta_k & \text{if } \delta_k < \delta'_{k'} & (\text{apply in } < \text{ order}) \\ \delta_k \circ \delta'_{k'} & \text{if } \delta'_{k'} < \delta_k & (\text{apply in } < \text{ order}) \\ \text{merge}(\delta_k, \delta'_{k'}) & \text{if } \delta_k \parallel \delta'_{k'} & (\text{merge concurrent}) \\ \text{undefined} & \text{otherwise} & (\text{error, see text}) \end{cases}$$

Fig. 4. Effect composition

that key can be ignored (by associativity of  $\sqcup$ ). Say the value of  $k$  at timestamp  $t_0 < t$  is known, then:  $\text{safe}(\Theta, k, t) = \text{safe}(\Theta, k, t_0) \sqcup \sqcup \{\delta'_k : \delta'_k \in \Theta \wedge t_0 < t' < t\}$ . This justifies the common algorithm for implementing lookup: locate the most recently-known value of the key (e.g., a cached value, or a checkpoint), and apply later updates to it in visibility order. If we impose that the dependency timestamp of running transactions be greater than  $t_0$ , then the prefix of  $\Theta$  up to  $t_0$  can be replaced by the *checkpoint* of values computed at  $t_0$ . We return to checkpointing in Section 9.12.

## 4.2 Semantics of transactions

Figure 2 presents the semantics of transactions, in a standard formal syntax. Although intimidating, it can be read as pseudocode, as we explain now. In addition, it is fully formal and unambiguous and easily translated to Coq statements.

For reference, we summarise our notations in Table 1. They are explained hereafter.

**4.2.1 Informal presentation.** Each rule describes an atomic transition: if the conditions above the horizontal line are satisfied, the transition under the line can take place. Unprimed variables represent pre-transition state (pre-state), primed ones post-transition (post-state). A label on the transition arrow under the line represents a client API call.

As an example, consider rule `BEGIN_TXN`. Below the horizontal line is API `begin(dt)` to start a new transaction, transitioning from pre-state  $(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r)$  on the left of the arrow  $\xrightarrow[\tau]{\text{begin}(dt)}$ , to post-state  $(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r')$  on the right. In the post-state, only the set of running transactions changes (as indicated by fact that  $\mathcal{T}_r$  is primed); notably, the store  $\sigma$  is not modified.

The part above the horizontal line can be read as a kind of pseudocode. The first clause chooses a transaction identifier  $\tau$  that was not used by any transaction (it is not in the union of transaction descriptors,  $\mathcal{T}_a \cup \mathcal{T}_c \cup \mathcal{T}_r$ , projected by  $\Pi_\tau$  along the  $\tau$  dimension). The second clause (with the universal quantifier) contains a precondition (detailed later) which must be true. The next clause creates a transaction descriptor, a 5-tuple (i.e., a struct), whose role is described later. Note that the transition is labeled by  $\tau$ , indicating that multiple transactions can begin independently of one another; however the semantics is that each such transition is atomic.

**4.2.2 Parameters.** The rules describe a transaction system, which is a tuple  $(\sigma, \mathcal{T}_a, \mathcal{T}_c, \mathcal{T}_r)$  consisting of a store  $\sigma$ , and sets of transaction descriptors  $\mathcal{T}_a$ ,  $\mathcal{T}_c$ , and  $\mathcal{T}_r$ , which keep track of aborted, committed and running (ongoing) transactions respectively.

A transaction descriptor is a tuple whose elements are, respectively, the transaction identifier  $\tau$ , its *dependency timestamp*  $dt$ , its *dirty set*  $\mathcal{E}$ , its *read buffer*  $\mathcal{R}$ , and its *commit timestamp*  $ct$ . As

stated previously (Section 3.4), the two timestamps define visibility between transactions. Initially, after rule `BEGIN_TXN`, both buffers are empty and the commit timestamp is invalid. For each key that is accessed, rule `INIT_KEY` initialises the buffers, and rule `UPDATE` updates them. Computation of the actual commit timestamp is deferred to the `COMMIT` rule.

The rules are parameterised by commands `lookup`, `doUpdate` and `doCommit`, specified in Figure 3. These commands will be specialised for each specific store variant, as detailed later: the map-based variant in Section 6, the journal-based variant in Section 7, and composition variant in Section 9.

**4.2.3 Transaction begin.** The API `begin(dt)` is enabled by Rule `BEGIN_TXN`, to begin a new transaction with dependency timestamp  $dt$ . Its semantics was briefly explained above. The rule does not set any constraints on  $dt$ ; semantically, any value is legal.

The *visibility* clause  $\forall t \in \Pi_{ct}(\mathcal{T}_r) : t \not\prec dt$  states that the transaction may not depend on some running transaction (i.e., not yet terminated). Suppose this clause was not present: then this transaction might attempt to read a value that has not yet been written or will be aborted. We discuss the synchronisation required for enforcing it in Section ??.

The visibility clause can be re-written as  $\min_{\mathcal{T}}(ct) \not\prec dt$  or equivalently  $\min_{\mathcal{T}}(ct) \geq dt$ , where  $\geq$  is read “greater than, equal to, or concurrent with.”

**4.2.4 Reads and writes.** The rule `INIT_KEY` initialises the read buffer for some key  $k$ . As it does not have an API label, it can be called arbitrarily. It does not modify the store, but does modify the current transaction’s descriptor within the set of running transactions  $\mathcal{T}_r$ . The first clause above the horizontal line simply separates the descriptor of the current transaction (identified by  $\tau$ ) from the others. The second clause checks that  $k$  is not already in the read buffer, ensuring that this rule is called no more than once per key. The third clause calls `lookup` (specific to a store variant) for the requested key and the transaction’s dependency timestamp. The next clause initialises the read buffer with the return value of `lookup`. The final clause above the line stores the updated buffers into the descriptor.

The API `read $_{\tau}(k)$` , which returns the content of the read buffer, is enabled by Rule `READ`. It does not modify the store. The first clause pulls out the transaction descriptor. The third one identifies the return value as the value of  $k$  stored in the read buffer. The second clause requires that the requested key  $k$  be already mapped in the read buffer; if not (this key has not been used in this transaction before), this can be interpreted as a cache miss, which is resolved by first executing transition `INIT_KEY`.

The API `update( $\tau, k, \delta$ )` applies effect  $\delta$  to key  $k$ ; it is enabled by Rule `UPDATE`. It updates both the store and the transaction descriptor. The first two clauses are similar to `READ`, and similarly cause a cache miss if the key has not been used before (avoiding blind writes). It updates the read buffer (to support “read-my-own-writes”), and marks the key as “dirty,” i.e., modified by this transaction. It calls the variant-specific command `doUpdate`, discussed later in the context of each variant.

**4.2.5 Transaction termination.** A transaction terminates, either by aborting with no effect, or by committing, which applies its effects to the store.

Rule `ABORT` moves the current transaction’s descriptor from  $\mathcal{T}_r$  to  $\mathcal{T}_a$ , marking it as aborted. It does not make any other change.

API `commit $_{\tau}(ct)$`  takes a commit timestamp argument. It is enabled by rule `COMMIT`, which modifies the store, the running set, and the committed set. Above the line, the first clause is as

usual. The client is free to choose any value for  $ct$  as long as it satisfies the constraints stated in the rule: it is unique (it does not appear in  $\mathcal{T}_c$ ); it is greater or equal to the dependency.

Command `doCommit` (discussed later, in the context of each store variant) provides the new state of the store; it should ensure that the effects of the committed transaction become visible in the store, labelled with the commit timestamp. Finally, the transaction descriptor, now containing the commit timestamp, is moved to the set of committed transactions.

**4.2.6 Correctness and performance challenges.** The specification is unambiguous and fits in a single page. In Section 5, we translate it into Coq for verification. Furthermore, we just showed how to read it as a kind of pseudocode. Translating it to sequential, imperative code is relatively straightforward; our Java implementation is called *transaction coordinator* and can be found [in the code repository](#).

Different coordinators can execute in parallel, and failures may occur. Examination of the specification in this light reveals a number of correctness challenges:

- (1) The store is accessed concurrently by different transaction coordinators, being read in rules `INIT_KEY` and written in rules `UPDATE` and `COMMIT`. Unsurprisingly, care is needed to avoid interference, in order to maintain the illusion that transitions are atomic.
- (2) The read buffer and dirty set, being private to a coordinator, pose no difficulty; the set  $\mathcal{T}_a$  of aborted transactions is not used anywhere and does not need to be implemented. The set of committed transactions  $\mathcal{T}_c$  requires a modest amount of synchronisation between concurrent commits.
- (3) The visibility clause of `BEGIN_TXN`  $\forall t \in \Pi_{ct}(\mathcal{T}_r) : t \not\prec dt$  reads from  $\mathcal{T}_r$ , which is modified by both `BEGIN_TXN` and `COMMIT`. Furthermore the choice of  $ct$  may be delayed until `COMMIT`. Therefore, enforcing this clause synchronisation between transaction begin and commit.
- (4) `COMMIT` is the rule that ensures that a transaction is durable. A major challenge is ensuring that its transition appears atomic despite failures.

The visibility clause of `BEGIN_TXN` forbids to read from a non-committed transaction.<sup>3</sup> To illustrate why it is necessary, consider the following example: (i) Transaction  $T1$  has commit timestamp 1; (ii) Transaction  $T2$  starts with dependency timestamp  $2 > 1$ ; thus  $T2$  must read the writes of  $T1$ ; (iii) however,  $T1$  is slow and its committed effects reach the store only after the first read by  $T2$ . Clearly, this would be incorrect:  $T2$  should not start as long as  $T1$  is still running and has not finalised its transition to committed.

The specification also reveals challenges related to performance:

- The store accumulates versions without bounds. An implementation will address this issue by restricting the use of timestamps, and garbage-collecting obsolete versions.<sup>4</sup>
- In practice, the specification can be implemented in many ways, with different performance characteristics. The classical ad-hoc approach introduces complexity and rapidly diverges from the specification. We argue that this can be addressed in a systematic way, by suitable composition of a small number of basic variants.

We address these challenges in the rest of this paper.

<sup>3</sup>Alternatively, we could write a precondition clause on `COMMIT` or on `INIT_KEY`, with similar effects.

<sup>4</sup>The sets  $\mathcal{T}_a$  and  $\mathcal{T}_c$  also appear to grow without bounds. However, the former is a convenience and need not be implemented. The latter is used by `BEGIN_TXN` and `COMMIT`; inspection of their preconditions (visibility and uniqueness clauses) shows that it suffices to record the lower and upper bounds of dependency and commit timestamps.

### 4.3 Traces and histories

Notation  $\delta_k^{\text{ct}}$ , called a *tagged effect*, refers to  $\delta$  updating key  $k$  at commit timestamp  $\text{ct}$ . We call (legal) *trace* a set of tagged effects (ordered by  $<$ ), produced by the rules of Figure 2; a trace is denoted by meta-variable  $\Theta$ . A *history* is a linearisation of a trace.

### 4.4 Consistency

Cerone et al. [3] define Transactional Causal Consistency by the following properties: (i) (EXT) The value of a key first read by a transaction is that computed by previous transactions; (ii) (INT) later reads include the transaction’s own updates; (iii) (TRANSITIVE VISIBILITY) Visibility is transitive. We claim without proof that the rules of Figure 2 enforce EXT and INT. This should be obvious by inspection.

▷(TODO:) We prove in Coq that the rules also ensure TRANSITIVE VISIBILITY. Interestingly, if the clause  $\forall t \in \Pi_{\text{dt}} (\mathcal{T}_c \cup \mathcal{T}_r) : \text{ct} \not\prec t$  is omitted from COMMIT (as was the case in our first attempt), transitive visibility can be violated.◁ (Marc to @Annette+@Carla): That clause is crucial to correctness. I suggest to try the proof with and without it! ◁

Furthermore, when restricted to sequential data types and totally-ordered timestamps, the rules satisfy their specification of Parallel Snapshot Isolation (a.k.a. Non-Monotonic Snapshot Isolation) of Cerone et al. [3].

▷(TODO:) Assuming total order and enforcing no-write-conflict, prove SI.◁

We conjecture that SI and SSER can be encoded by strengthening the preconditions in BEGIN\_TXN and COMMIT, as follows. In SI, snapshots follow the prefix property; under Strict Serialisability, visibility is a total order. More detail is out of scope of this paper.

▷(FUTURE WORK:) Consistency/isolation levels:◁

- weak/TCC. Geo-replicate across  $N$  sites. Timestamp is a vector size  $N$ . Partial order abstracts geo-replication. Each entry in a *vector* (clock or timestamp) counts the events in the *sequential* history of a *replica*. Events of a replica are totally ordered w.r.t. the same replica.
- SI. Total order of commits, no concurrent writes to same key. Timestamps totally ordered (scalar).
- SSER. No concurrent reads or writes to same key. Timestamps totally ordered (scalar).

### 4.5 \*\*\*\* TODO \*\*\*\*

\*\*\*\* For implementation convenience, commit timestamps are locally monotonically increasing, i.e., greater than any commit that the session/client has observed so far. Set to *now()*. Not imposed by the model.

\*\*\*\* To enable GC, we will assume a monotonically-increasing minimal dependence timestamp. Not imposed by the model. Any key-versions underneath it can be replaced by their composition and GC’ed away.

; the history of a store. Unfolds with time. Current time = *now()*. Snapshots/commits visible iff  $\leq \text{now()}$  (non-visible events don’t return). Prefix of ongoing transaction visible. Geo-distribution: each site has its own *now()*.

## 5 ▷(TODO:) TODO◁ GENERAL LEMMAS

### 5.1 General results

All fine-grain interleavings of effects in a history, while maintaining the relative order of transactions, are equivalent:

$$\begin{aligned}
\text{lookup}(\sigma, k, dt) &= \bigsqcup \max_{t < dt} \{\sigma(k, t)\} \\
\text{doUpdate}(\sigma, \_, \_, \_) &= \sigma \\
\text{doCommit}(\sigma, \_, \mathcal{E}, \mathcal{R}, ct)(k) &= \begin{cases} \sigma(k) \cup \{(\mathcal{R}(k), ct)\} & \text{if } k \in \mathcal{E} \\ \sigma(k) & \text{otherwise} \end{cases}
\end{aligned}$$

where  $\max_{t < dt}$  returns the element with the highest timestamp less than  $dt$ .

►(Marc 2023-01-15:) Premature optimisation? Let lookup return the composition of all values, and optimise it to max. ◀

Fig. 5. Operations of map store

$$\begin{aligned}
\text{lookup}(\sigma, k, dt) &= \bigsqcup (\text{committed\_before}(\sigma, k, dt)) \\
\text{doUpdate}(\sigma, \tau, k, \delta) &= \sigma \triangleright (\text{update}, \tau, k, \delta) \\
\text{doCommit}(\sigma, \tau, \_, \_, ct) &= \sigma \triangleright (\text{commit}, \tau, ct)
\end{aligned}$$

with the following notation:

- $\sigma$  is a journal-based store.
- $\triangleright$  represents concatenation.
- ‘update’ and ‘commit’ are tags to distinguish the type of a journal record.
- $\text{committed\_before}(\sigma, k, dt)$  denotes the subsequence of journal  $\sigma$ , of records tagged with key  $k$ , that have a commit timestamp less or equal to  $dt$ ; formally:

$$\text{committed\_before}(\sigma, k, dt) \triangleq \{(\delta, t) \mid \sigma = \sigma_0 \triangleright (\text{update}, \tau, k, \delta) \triangleright \sigma_1 \triangleright (\text{commit}, \tau, t) \triangleright \sigma_2 \text{ and } t < dt\}$$

Fig. 6. Operations of Journal store

LEMMA 5.1 (REDUCIBILITY TO SEQUENTIAL TRACE). *Concurrent trace reducible to equivalent sequential trace, under MVCC/GSI model.*

Other lemmas:

- For a given  $k$  and  $t$ ,  $\text{lookup}(k, t)$  is unique. The store is a function. (Note: reads block until  $t \leq \text{now}().$ )
- Any two variants of the generic store are equivalent (same reads, same writes).
- It follows that *composition*, as defined in Section 9, is OK.

Invariants:

- Commit timestamps assigned by a site are monotonically increasing. Anything stronger, e.g., the LUB of site timestamps is monotonically increasing?
- Any other useful invariants?

\*\*\*\* ►(TODO:) \*\*\*\* Lemma (prove it!): For all  $\sigma, \sigma'$  built from the same trace  $\Theta$ ,  $\sigma(k, t) = \sigma'(k, t)$ . ◀

\*\*\*\* ►(TODO:) \*\*\*\* A more recent value hides any previous effects, i.e.,  $\bigsqcup \max_{t < dt} \{\sigma(k, t)\} \in \text{Value} = \bigsqcup \{\sigma(k, t) \mid t < dt\}$  ◀

\*\*\*\* ►(TODO:) \*\*\*\* Therefore: Effect summaries: cut off the history at some lower bound. Saves space, against the promise that all current or future transaction's  $dt$  is at least the lower bound. ◀

## 6 MAP-BASED STORE

### 6.1 Map store semantics

Our first variant is the *map-based store*, which stores assignments. When a transaction commits,  $\text{doCommit}$  eagerly assigns new versions of the keys marked dirty. When reading,  $\text{lookup}$  returns the most recent version in the store, since it replaces any older one.

		$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$		
(a)	$\sigma$	0	1	2	3	4	5	6	7
history			$x = 0$	$y = 0$	$x + 4; y = 0$	$x + 1$	$x - 3$	$x + 2$	
(b)	$x$		$x.1 = 0$		$x.3 + 4$	$x.4 + 1$	$x.5 - 3$	$x.6 + 2$	
journal	$y$			$y.2 = 0$	$y.3 = 0$				
		J:0(0,3]				J:3(3,7]			
(c)	$x$		$x.1 = 0$		$x.3 = 4$	$x.4 = 5$	$x.5 = 2$	$x.6 = 4$	
map	$y$			$y.2 = 0$	$y.3 = 0$				
		M:0(0,4]				M:4(4,7]			
(d)	$x$		$x.1 = 0$	$\Rightarrow$				$x.6 = 4$	
cache	$y$			$y.2 = 0$		$y.4 = 0$	$\Rightarrow$	$\Rightarrow$	
		M:{x.[1,2], x.6, y.2, y.[4,6]}							
(e)	$x$			$x.1 = 0$		$x.4 = 5$			$x.6 = 4$
check-point	$y$			$y.2 = 0$		$y.3 = 3$			
		M:0(1,2]			M:2(3,4]		M:4(6,7]		

Fig. 7. Store variants and store composition. Subfigure (a) represents a history; the following ones different variants and ranges. To simplify the figure, timestamps are assumed integer (scalar), and the history is represented as a sequence of transactions with strictly monotonic dependency and commit timestamp. Notation:  $x, y$ : keys;  $=0$ : assign key with value 0;  $+2, -3$ : increment key's value by 2, decrement by 3;  $y.4 + 1$ : increment  $y$  by 1 with commit timestamp 4.  $\Rightarrow$  extends validity of a cache entry; in  $J:1(2, 3]$ :  $J$  = journal/ $M$  = map,  $low\_history=1$ , domain ( $low\_lookup = 2, high\_lookup = 3$ ).

A map store is illustrated informally in Figure 7(b).<sup>5</sup> A map store can be viewed as a matrix of keys (rows) and timestamps (columns), where a cell contains either a Assignment or  $\perp$ . Each successive value associated with a key is a new, write-once version, called *key-version*. This map store grows without bound; we consider bounding its size later (Section 9.10).

The cell at index  $(k, t)$  is populated iff some transaction commits an update to key  $k$  at with timestamp  $t$ . Conversely, if key-timestamp pair  $(k, t)$  exists in the domain of map store  $\sigma$ , noted  $(k, t) \in \text{dom}(\sigma)$  or equivalently  $\sigma(k, t) \neq \perp$ , that means that a transaction with commit timestamp  $t$  assigned  $k$  with value  $\sigma(k, t)$ .

This invariant is enforced by the map store semantics of Figure 5, starting from an initially empty map store, as follows:

- Command  $\text{lookup}(\sigma, k, dt)$  (called in the context of rule `INIT_KEY`), finds in store  $\sigma$  the most recent version of key  $k$  that is visible from the current transaction's dependency timestamp  $dt$ . If there are multiple concurrent most-recent versions, it merges them.
- Command `doUpdate`, invoked by rule `UPDATE`, is a no-op for this variant.

<sup>5</sup>Ignore for now the dividing line between 4 and 5, and the bound indications underneath.



- In the context of COMMIT, command  $\text{doCommit}(\sigma, \tau, \mathcal{E}, \mathcal{R}, \text{ct})$  copies, for each key updated by transaction  $\tau$ , as indicated by dirty set  $\mathcal{E}$ , into a new version, labelled by the transaction's commit timestamp  $\text{ct}$ .<sup>6</sup>

## 6.2 Map store implementation

The map-based store (MStore) is implemented in Java and use a hash map to store the values. The MStore is a crash-free store, which means that it does not have a recovery mechanism. This means that the MStore assumes that the system will always gracefully shutdown.

The challenge talked about in *section 4.2.6* are addressed in the MStore implementation in the following way. The lookup in the MStore is simple as it queries the map with the pair  $k$  and  $\text{adept}$  and returns a value with a lower *committimestamp*. To serve multiple concurrent transaction the MStore uses MVCC (Multi-Version Concurrency Control) to ensure that reads and writes do not interfere with each other. For each  $k$  the map stores  $v$  with a timestamp that represent its commit time. When a read operation is performed, the MStore selects the most recent version of the key-value pair that is older than the timestamp of the read operation. This ensures that the read operation sees a consistent view of the  $k$ . This means the  $\text{doUpdate}$  is a no-op as any update to the map would be visible but updates the  $\mathcal{R}$  and the  $\mathcal{E}$  of the transaction object instead. An advantage of  $\text{doUpdate}$  is that it does not require keeping track of aborted transactions as an abort simply deletes the transaction object. This also means that each client can work on its local copy of an object without concurrency issues. The  $\text{doCommit}$  takes a commit timestamp and updates the commit timestamp of every  $v$  that is present in the Write Set and then add them to the map. Once finished, the transaction is added to the  $\mathcal{T}_c$ , and the timestamp server is notified that the commit has finished advancing *max\_allowed\_timestamp*.

## 6.3 Checkpoints

»(TODO:) A checkpoint is a specific map store. A checkpoint is equivalent to { the full store restricted to timestamps above the checkpoint's lower bound }.« (Marc:) Move into composition section?«

## 7 JOURNAL STORE

An alternative variant is the journal-based store. The journal stores computations, and materialises values lazily on lookup (this is in contrast to the map-based store, which eagerly records Assignments in the store).<sup>7</sup>

Figure 7(b) provides an informal illustration.<sup>8</sup>

Figure 6 gives the formal semantics of a journal store. Let us start with  $\text{doUpdate}$ , called in the context of the Rule UPDATE from Figure 2. It simply appends an update record to the journal, containing arguments transaction identifier  $\tau$ , key  $k$ , and effect  $\delta$ . Similarly,  $\text{doCommit}$ , called in the context of Rule COMMIT, appends a commit record containing the transaction identifier and commit timestamp. The effect and read buffers are discarded.

The real action is in lookup. To find the value of  $k$ , it extracts from the log the effects of all transactions that the current transaction depends upon, i.e., those whose commit-timestamp is before  $\text{dt}$ , and composes these effects (in visibility order, merging concurrent effects, per Figure 4).

»(TODO:) Stop backwards search at Assignment.«

<sup>6</sup>Strictly speaking, copying only the dirty keys is not logically necessary. However, this is justified in practice, since the space of keys is essentially unbounded (e.g., if keys are arbitrary strings).

<sup>7</sup>Technically, this is a redo log. An undo log would store inverse effects.

<sup>8</sup>Again, ignore for now the dividing line between 3 and 4, and the bounds indicated underneath.

Note that, just like the map-based store, the memory footprint of the journal-based store is unbounded. Reading and writing the journal requires care, as it is a shared resource among concurrent transactions. As a journal is sequential, flushing the single commit record to secondary storage ensures crash-atomicity. Indeed, in the event of a crash, either the commit record and all preceding records were written, making the whole transaction durable, or it has not, and the transaction is considered aborted on recovery. Thanks to sequential writes, a persistent journal generally has better throughput than a persistent map store. Our experiments in Section 10 confirm these well-known differences.

### 7.1 ~~▶(TODO:) TODO◀~~ Journal store proofs

In the Coq specification we proved that there is a bisimulation between the journal and the store. The following theorems were defined and proved:

```
Theorem journal_to_store_trace :
  forall tr j trx,
    journal_trace (empty_journal, ([], [], None)) tr (j, trx) ->
    exists s, store_trace (empty_store, ([], [], None)) tr (s, trx).
```

```
Theorem store_to_journal_trace :
  forall tr s trx,
    store_trace (empty_store, ([], [], None)) tr (s, trx) ->
    exists j, journal_trace (empty_journal, ([], [], None)) tr (j, trx).
```

Note that the semantic model does not consider interleaving of transactions, i.e., transactions are sequentially executed.

Proof objective:

- Journal bi-simulates generic store for the same history.
- From Lemma “Concurrent trace reducible to equivalent sequential trace,” it follows that The absolute ordering in a journal stream is irrelevant, as long as it respects visibility: (i) the relative ordering within a transaction, and (ii) the timestamp ordering of commit records of different transactions.

### 7.2 Journal store implementation

The journal-based store (JStore) is implemented in java using the same API as the MStore. JStore is implemented as a lazy store, i.e., it does not materialize values until they are requested. It is built using an append-only structure in Java where all the transaction operations are recorded in the form of records. There are three main types of records in this implementation:

- Begin record, system record that marks the beginning of a transaction. It contains a  $\mathcal{T}_{ID}$  and a dependency timestamp.
- Effect records, contains updates performed inside the transaction and the transaction identifier it belongs to.
- End record, system record that marks the end of a transaction. It is either an abort record or a commit record. The former only contains the transaction identifier and the latter also contains the ct of the transaction.

During the life cycle of a transaction the *begin*, updates and end operations of a transaction are persisted in the journal in a write-ahead manner. This provides crash resilience as it is possible to implement a recovery mechanism that will replay the journal and recover the state of the JStore.

The Coordinator still maintains a *dirtyset*  $\mathcal{E}$  and a *readbuffer*  $\mathcal{R}$  but it is used only for read-your-own-writes.

Because of the lazy nature of the JStore, the lookup is implemented by reading the journal and applying the effects of the transactions that are visible to the current transaction. This is necessary because all the updates are directly written to the JStore through the `doUpdate` operation. This put the burden of checking the visibility on the lookup which checks if updates have a matching commit record in the journal.

The JStore is implemented using a single thread executor to ensure that the journal is written sequentially. Reading while a concurrent transaction is writing to the journal is not a problem as any  $\text{futur.ct} \text{not} < \mathcal{T}_r.\text{ct}$ . By storing the all the  $\delta$  in the journal, the JStore can materialize any  $v$  at any time solving any concurrent read problem.

Similarly to the MStore, the JStore ensures that `doCommit` is visible atomically by notifying the end of the commit to the timestamp server.

On the other hand  $\mathcal{T}_a$  leave a trace in the journal therefore the journal writes an abort record that invalidates the transaction update records.

The JStore being crash resistant the recovery ensures that for every transaction there is a matching begin and end record. If a commit record is missing, the transaction is aborted by writing an abort record. The recovery also initializes the timestamp server with the highest  $\text{ct}$  found in the journal.

Once completed the JStore can resume operations.

## 8 DESIGN CHALLENGES AND IMPLEMENTATION CHALLENGE

This section describes the general architecture of the system and the challenges we faced during the implementation around the store.

One design choice was to decouple the client handling from the different stores. This was done by implementing a coordinator role that handles the client requests and would be responsible for handling the transactions. This give us the flexibility to implement different stores and to easily switch between them. The coordinator maintains a transaction object, and execute the operations on the store when appropriate.

Another design choice was to implement a timestamp server that would be responsible for generating monotonically increasing unique timestamps and allow the coordinator to query atomically the current timestamp. The timestamp server was fairly simple in the beginning where it would maintain a single timestamps that it would modify and read atomically.

The  $\text{ct}$  of a transaction is selected at the beginning of the transition from running to committed. The coordinator ask for the timestamp server to generate a new timestamp atomically and assign it to the transaction. `doCommit` of the store is called with the  $\text{ct}$  of the transaction and can take some time to complete.

This meant the following situation could occur: a  $T1$  transaction start its commit phase, the timestamp server is queried and the `doCommit` of the store is called with  $T1.\text{ct}$ . Meanwhile another transaction starts and is assigned a higher timestamp than  $T1.\text{ct}$  without the  $T1$  transaction having committed. This would lead to a violation of the definition of visibility as  $T1$  would be visible partially to other transactions. This problem is solved by enforcing:  $dt \not\prec \min_{\mathcal{T}_R}(\text{ct})$

To avoid strong synchronisation we maintain two different atomic integers  $\text{max\_allowed\_dt}$  and  $\text{min\_allowed\_ct}$ . Both are monotonically increasing and  $\text{max\_allowed\_dt} \not\prec \text{min\_allowed\_ct}$ .

$$\begin{aligned}
\text{lookup}(\sigma, k, dt) &= \begin{cases} \text{as in Figure 5 or 6} & \text{if } (k, t) \in \mathcal{D}_\sigma \\ \perp & \text{otherwise} \end{cases} \\
\text{doUpdate}(\sigma, \tau, k, \delta) &= \begin{cases} \text{as in Figure 5 or 6} & \text{if } (k, t) \in \mathcal{D}_\sigma \\ \sigma & \text{otherwise} \end{cases} \\
\text{doCommit}(\sigma, \tau, \mathcal{E}, \mathcal{R}, ct) &= \begin{cases} \text{as in Figure 5 or 6} & \text{if } (k, t) \in \mathcal{D}_\sigma \\ \sigma & \text{otherwise} \end{cases}
\end{aligned}$$

Fig. 8. Store operations, restricted to domain

$$\begin{aligned}
\text{lookup}(\{\sigma_1, \sigma_2\}, k, dt) &= \text{lookup}(\sigma_1, k, dt) \sqcup \text{lookup}(\sigma_2, k, dt) \\
\text{doUpdate}(\{\sigma_1, \sigma_2\}, \tau, k, \delta) &= \text{doUpdate}(\sigma_1, \tau, k, \delta) \parallel \text{doUpdate}(\sigma_2, \tau, k, \delta) \\
\text{doCommit}(\{\sigma_1, \sigma_2\}, \tau, \mathcal{E}, \mathcal{R}, ct) &= \text{doCommit}(\sigma_1, \tau, \mathcal{E}, \mathcal{R}, ct) \parallel \text{doCommit}(\sigma_2, \tau, \mathcal{E}, \mathcal{R}, ct)
\end{aligned}$$

Fig. 9. Operations of composed store

ADD\_MINISTORE

$$\begin{array}{c}
S = \{\sigma_1, \sigma_2, \dots\} \quad \mathcal{D}_S = \mathcal{D}_{\sigma_1} \cup \mathcal{D}_{\sigma_2} \cup \dots \\
\forall (k, t) \in \mathcal{D}_{\sigma_0} \cap \mathcal{D}_S : \text{lookup}(\sigma_0, k, t) = \text{lookup}(S, k, t) \quad S' = S \cup \sigma_0 \quad \mathcal{D}_{S'} = \mathcal{D}_S \cup \mathcal{D}_{\sigma_0} \\
\hline
(S, \mathcal{D}) \xrightarrow{\text{addmini}(S, \mathcal{D}_S, \sigma_0, \mathcal{D}_{\sigma_0})} (S', \mathcal{D}_{S'})
\end{array}$$

REMOVE\_MINISTORE

$$\begin{array}{c}
S = \{\sigma_0, \sigma_1, \sigma_2, \dots\} \\
\mathcal{D}_S = \mathcal{D}_{\sigma_0} \cup \mathcal{D}_{\sigma_1} \cup \mathcal{D}_{\sigma_2} \cup \dots \quad S' = S \setminus \sigma_0 \quad \mathcal{D}_{S'} = \mathcal{D}_{\sigma_1} \cup \mathcal{D}_{\sigma_2} \cup \dots \\
\hline
(S, \mathcal{D}) \xrightarrow{\text{rmvmini}(S, \mathcal{D}_S, \sigma_0, \mathcal{D}_{\sigma_0})} (S', \mathcal{D}_{S'})
\end{array}$$

»(Marc:) Make  $\mathcal{D}$  explicit argument, or part of store?«Fig. 10. Operational semantics of store composition.  $S$  is a composition of ministores  $\sigma_1, \sigma_2, \dots$ , and  $\mathcal{D}_\sigma$  is the domain of  $\sigma$ .

The *min\_allowed\_ct* is updated by the timestamp server when a new timestamp is generated. And the *max\_allowed\_dt* is updated by the store when a transaction finishes committing. We also ensure that all the transaction are visible in the order of commit timestamp issues by the timestamp server. The timestamp server keeps a list of transaction that have committed and only advances the *max\_allowed\_dt* in the order of the list.

Our simple stores implementation present challenges related to performance that are difficult to solve. Both stores grow without bounds and both having every  $v$  ever written it makes a garbage collection necessary. One way to this problem is to partially forget the *history* of the store without violation the consistency requirement of the store.

## 9 COMPOSING STORES

»(Implementation Goals:)

- **journal + checkpoint + truncate. Maintains the fault tolerance properties.**
- **checkpointing**
- **cache**

◀

▷(TERMINOLOGY:) In this section, “domain” has a slightly different meaning from previous ones. Alternatives: realm, area, field, province?◀ ▷(NOTATION:) It may be more consistent if a segment is closed on its upper bound and open on its lower bound?◀

The basic store variants described so far (map and journal-based) are simple but have performance issues, e.g., they grow unbounded. We argue in this section that the complex design of modern stores can be explained as a *composition* of basic variants.

The composition rules are particularly simple. A composed store is simply a set of stores (called *ministores*), each restricted to a specific *domain*. The lookup, doUpdate and doCommit operations over the composition extend recursively to its component ministores.

For instance, a cache is an in-memory record of recently-used versions; ignoring the details of the caching policy, the cache is a map ministore, whose domain is an arbitrary subset of key-timestamp pairs. For instance, Figure 7(d) represents a cache containing versions  $\{(x, 1), (x, 6), (y, 2), (y, 4)\}$ , where  $(x, 1)$  and  $(y, 4)$  are known to remain valid until timestamps 2 and 6 respectively. A cache is fast, but must be composed with a (possibly slower) persistent store covering the full range of keys and timestamps; for instance either the journal in Figure 7(b) or the map in (c).

Another example is a write-ahead log (WAL). A WAL combines the high sequential throughput and fault-tolerance of a journal ministore with a fast-response random-access map ministore. When an update commits, it is written (in a crash-tolerant manner) to the journal; later, it is copied and persisted to the map. Thus, the journal’s time domain includes the most recent updates, whereas the map’s is slightly delayed. For instance, the first map in Figure 7(c) might be combined with the second journal in Subfigure (b). Furthermore, this WAL might be paired with the in-memory cache of Subfigure (d) to decrease average response time.

The LSM-Tree approach [14] decomposes a map into a series of layers. The top layer contains the most recent updates, which percolate downwards towards the bottom layer as they age. When reading, the system queries the layers from top to bottom in succession, returning the first value found. We will characterise this approach as a composition of map ministores that differ in the time domain, most recent at the top.

In future work, we plan to describe a geo-distributed store, such as Antidote [1], as a composition of mini-stores with concurrent time domains.

### 9.1 Domain of a store

We associate store  $\sigma$  with its *domain*, noted  $\mathcal{D}_\sigma \subseteq \text{Key} \times \text{TS}$ , and modify the store operations lookup, doUpdate and doCommit to be significant only within the associated domain; outside, they are no-ops. This is formalised in Figure 8.

We will call *top-level* store one whose domain is the full universe  $\text{Key} \times \text{TS}$ , and *ministore* one whose domain is a strict subset. A ministore may restrict the key domain, the time domain, or both. For instance, *sharding* composes ministores covering disjoint subsets of the key domain.

Caches aside, the rest of this paper focuses on stores whose time domain is a continuous *segment* of time (*low\_lookup*, *high\_lookup*). We also associate each such ministore timestamp *low\_history* that represents the beginning of the history used to compute this segment. We summarise this

information with notation  $X:\text{low\_history}(\text{low\_lookup}, \text{high\_lookup}]$ , where  $X$  is  $M$  (for a map) or  $J$  (for a journal).

For instance (assuming for simplicity that timestamps are integers), a ministore  $\sigma = J:10(20, 30]$  is a journal, representing the history between times 10 (exclusive) and 30 (inclusive), but for which  $\text{lookup}(\sigma, k, t)$  is significant only if  $20 < t \leq 30$ . Figure 7(b) represents two journal ministores  $J:0(0, 3]$  and  $J:3(3, 7]$ ; the second contains only the incremental effects in range  $(3, 7]$ , ignoring those at or before its  $\text{low\_history} = 3$ . Subfigure (c) represents two map ministores; the second one reflects only the incremental versions created starting at timestamp 5. Subfigure (e) shows three incremental checkpoints, recording only the last value in the domain, and only if updated between  $\text{low\_history}$  and  $\text{high\_lookup}$ .

Obviously,  $\text{low\_history} \leq \text{low\_lookup} < \text{high\_lookup}$ . Typically, either  $\text{low\_history} = \text{low\_lookup}$  (it records all versions in range) or  $\text{low\_lookup} + 1 = \text{high\_lookup}$  (a checkpoint summarising the updates between  $\text{low\_history}$  and  $\text{high\_lookup}$ ).

▷(Marc:) Needs improvement.◁

## 9.2 Composition of stores

A composed store is simply a set of ministores. Its domain is the union of components' domains. An operation on the composed store just calls itself recursively on the component ministores.

Formally, consider ministores  $\sigma_1$  and  $\sigma_2$ , with domains  $\mathcal{D}_1$  and  $\mathcal{D}_2$  respectively. Their composition  $\{\sigma_1, \sigma_2\}$  has domain  $\mathcal{D}_1 \cup \mathcal{D}_2$ . Its operations are defined by the equations in Figure 9, where  $\parallel$  denotes parallel composition.<sup>9</sup>

We showed earlier that stores are functional and that lookup on any map or journal returns the same result. It follows that, if a key-timestamp pair is in domain of two components, it is equivalent to perform lookup on one, on the other, or on both. Formally,  $(k, t) \in \mathcal{D}_1 \cap \mathcal{D}_2 \implies \text{lookup}(\sigma_1, k, t) = \text{lookup}(\sigma_2, k, t)$ . ▷(TODO:) Check this assertion! It might be incorrect for a journal whose history does not start at 0.◁ Therefore it is often most efficient to perform lookup on the most recent ministore first, and to stop as soon as a lookup returns an Assignment. Similarly, an operation can skip recursing into a component for which the arguments are not in domain.

By abuse of notation, we identify a composed store with a store, and note  $\sigma = \{\sigma_1, \sigma_2\}$ .

The composition of segment stores, whose domains are adjacent or overlap, behaves like single union segment. For instance, in Figure 7 the composition of mini-journals  $J:0(0, 3]$  and  $J:3(3, 7]$  behaves like  $J:0(0, 7]$ . Combining the second checkpoint  $M:2(3, 4]$  with mini-map  $M:4(4, 7]$  behaves like  $M:2(3, 7]$ . Similarly, mini-map  $M:0(0, 4]$  combined with mini-journal  $J:3(3, 7]$  behaves like  $M:0(0, 7]$ . It is not useful to compose  $M:0(1, 2]$  with  $J:3(3, 7]$ , as this would leave a gap.

Formally, composing  $M:\text{low\_history}_1(\text{low\_lookup}_1, \text{high\_lookup}_1]$  with  $M:\text{low\_history}_2(\text{low\_lookup}_2, \text{high\_lookup}_2]$  where there are no gaps between the two, i.e.,  $\text{high\_lookup}_1 \geq \text{low\_lookup}_2 \wedge \text{high\_lookup}_1 \geq \text{low\_history}_2$ , behaves like  $M:\text{low\_history}_1(\text{low\_lookup}_1, \text{high\_lookup}_2]$ . Similarly for J-segment ministores.

Another interesting case is when their histories are adjacent but not their domains. In this case, the history of the composition is the composition of histories, but the domain is just the higher one. For instance, in Figure 7, composing  $M:0(0, 4]$  with checkpoint  $M:4(6, 7]$  behaves like  $M:0(6, 7]$ .

Formally, if  $\text{high\_lookup}_1 \geq \text{low\_history}_2$ , then their composition behaves like  $M:\text{low\_history}_1(\text{low\_lookup}_2, \text{high\_lookup}_2]$ .

<sup>9</sup>We defer the discussion of how to ensure the composition remains atomic to Section ??.

»(Marc 2023-02-18 20:16:) \*\*\*\*\* stopped here«

### 9.3 Top-level (composed) store

Extra *no-gap* invariant on top-level: represents the whole history: time domain is an interval starting at 0.

Extra *no-broken-reads* invariant on top-level commit:  $low\_lookup < dt$ . »(Marc:) Does this duplicate the visibility clause of begin-txn?«

Exhibit new top-level commit rule.

### 9.4 Changing composition

Maintain lookup is functional.

Operations: addmini, rmvmini. How they affect interval domain.

Changing bounds: add new bound, remove old bound. Maintain the no-gaps invariant: add first.

### 9.5 Garbage collection

Advance  $low\_lookup$  monotonically. Frees memory. Loses information: snapshots that are in the freed zone would be in error.

Enforce no-broken-reads by either not advancing  $low\_lookup < \min_{\tau}(dt)$ , or by aborting running transactions that violate.

### 9.6 \*\*\*\*\* lost text \*\*\*\*\*

Therefore a composed store whose joint domain starts from time 0, without holes, is equivalent to the full store. Causal consistency implies that there are no holes!

### 9.7 Proof goals

»(Proof goals for composition:)

- Composing cache & either map or journal = equivalent to cacheless.
- Composing journal & checkpoint, where the lower bound of the journal overlaps the upper bound of the checkpoint, is equivalent to the full store.
- Move checkpoint + truncate. Can be done concurrently with transactions.
- Chunked (= incremental) checkpoints is equivalent to non-chunked (monolithic) checkpoint.
- Multiple layers are equivalent to a single layer

«

### 9.8 Compositon with the Conductor

In order to handle the compositon we introduce the Conductor. The Conductor sits between the client and the ComposedStore and maintains a topological view of the system.

The Conductor uses the same interface used by the JStore client and the MStore client. This allows a ComposedStore to replace a JStore or an MStore transparently in any system. The Conductor uses the topology of the composition to ensure safety between the different layers.

For every topology the Conductor maintains every read is started on the first layer and then transfered to next. But this mecanism is not suited for updates as they must be written in a layer that has persistency and if the target model requires it crash-resistance.

The Conductor chooses the first level accessible level in the topology that satisfies the memory model as the write point of the system.

On top of the topology the Conductor maintain a domain for each layer. These domains are used to ensure safety throughout the ComposedStore.

### 9.9 Implementation of a simple composed store (Saalik)

The initial composed store is a composition of two stores, a journal (*JStore*) and a map (*CheckpointStore*). As an implementation decision we use the journal as the level where we advance the *history* and the map as a checkpointing level. The higher bound of the *JStore* advances as transactions are committed. While we could have used the map as the first level of the store, we decided to use the journal because the journal gives us the ability to perform a recovery in case of a crash.

The implementation challenges described in the Section 4.2.6 are solved by using either the *JStore* and the *MStore*.

The first implementation challenge of the ComposedStore is to ensure that checkpointing is done properly. Here is the algorithm used to perform a checkpoint:

- (1) Choose the next checkpoint timestamp,  $t_{chk\_next}$  such as  $t_{chk\_next} > t_{min\_store}$  and  $t_{chk\_next} < min_{allowed\_ct}$  and  $t_{chk\_next} < \mathcal{T}_r.dt$
- (2) Find the commit record of the last committed transaction between  $t_{min\_store}$  and  $t_{chk\_next}$ .
- (3) Write the CheckpointBegin record with  $t_{chk\_prev}$  and  $t_{chk\_next}$ .
- (4) Lookup all the versions that are visible at  $t_{chk\_next}$  and write them to the Checkpoint Store.
- (5) Update the higher bound of the *CheckpointStore* to  $t_{chk\_next}$ .
- (6) Write the CheckpointEnd record with the  $record\_id$  that marks the new start of the *JStore*.
- (7) Advance the lower bound of the *JStore* to  $t_{chk\_next}$ .

While the journal had records for the begin, effect and end of a transaction, we encountered a problem where crashing during a checkpoint would result in a state where correctness was not guaranteed.

To solve the second challend we added two new records to the journal, the CheckpointBegin and CheckpointEnd records. The first record, CheckpointBegin, writes all the information about the next checkpoint in the journal. This includes the old checkpoint timestamp and the new checkpoint timestamp. Once the checkpoint is completed, the CheckpointEnd record is written to the journal, with the  $record\_id$  of the first begin record of a transaction that has not been checkpointed.

The two new records are not sufficient to ensure that the checkpointing process is safe in case of a crash. The higher bound of the *CheckpointStore* can be corrupted while updating so we use a two phase update. We use two higher bounds, an active and an inactive one. During the first phase we write the new higher bound to the inactive one and then we swap the two.

Recovery in our implementation must be idempotent, meaning that if the recovery process is interrupted, it can be restarted and it will not corrupt the state of the system. So using the the two previous mechanism we describe the recovery process as follows:

- (1) Journal recovers as describe in Section 7.2
- (2) If there is a CheckpointBegin record without a CheckpointEnd record it means that the checkpointing process was interrupted.
- (3) The checkpoint is restarted from the beginning, as writing a Key-Value pair that is already present in the *CheckpointStore* is idempotent.



- (4) Once all the key-versions are written to the map, the Conductor check if the inactive and the active  $t_{store\_max}$  are correct and if not are updated with the values found in CheckpointBegin record.
- (5) A CheckpointEnd record is written to the journal.
- (6) The Conductor then initialises its local states of the *JStore* and the *CheckpointStore* domain information.

### 9.10 Bounded Stores (Gustavo, Annette, Carla)

Consider two stores that satisfy the generic specification. The “old generation” has an upper bound, and the “new generation” has a lower bound.

Coverage requirement: the bounds overlap.

Updates go only to the new generation (i.e., no transactions underneath the lower bound).

▷(TODO:) Prove these lemmas◁

Lemma: Bounding does not violate the atomicity of transactions.

Garbage Collection Lemma: given some timestamp  $t$ , the set of key-versions committed at  $< t$  that are accessible by transactions is monotonically non-decreasing. Any key-version that is not accessible remains non-accessible; consequently, obsolete key-versions in the old generation can be GC'ed.

### 9.11 Moving the bounds (Gustavo, Annette, Carla)

Restrictions:

- New transactions go to the new generation.
- The bounds are monotonically non-decreasing.
- The old generation's bound grows before growing that of the new generation, ensuring that they continuously overlap. key-versions are moved from the new to the old, so that complementarity continues to hold.

▷(TODO:) Prove these lemmas.◁ Truncation Lemma: Moving the bounds does not violate the atomicity and durability requirements of transactions.

### 9.12 Checkpointing and truncation (Saalik)

As stated in Section 4.1, a history of updates up to some timestamp can be replaced by the values of keys at that timestamp. Define a *checkpoint*  $\sigma_{t_0}$  at to be a map store containing these values at time  $t_0$ :  $\sigma_{t_0} \triangleq \{t \mapsto \text{safe}(\Theta, k, t_0)\}$ . Then we can compute any later version of the key from the checkpoint and its later updates:  $\text{safe}(\Theta, k, t) = \text{lookup}(\sigma_{t_0} \sqcup \bigsqcup \{\delta'_k : \delta'_k \in \Theta \wedge t_0 < t' < t\})$ .

One version of a composed store is done with a journal containing every new transaction and a Map-Store containing old information that we will call Checkpoint Store.

Every new transaction writes in the journal until it is over. And periodically the journal is checkpointed in the Checkpoint Store. Every time a checkpoint is created we persist a checkpoint record in the Journal. The checkpoint start record contains information about range of records that will be checkpointed. Once the checkpoint is persisted to the checkpoint store an endCheckpoint record is written to disk with the updated lower bound of the journal.

To represent the active (Safe? Important?) portion of the journal, the system maintains two values. A Low-Watermark, that represents the oldest record visible to an ongoing transaction. It is

advanced every time a new checkpoint is written to the store and can be recovered using the last endCheckpoint written in the journal. The second one High-Watermark is the most recent record that can be visible to an on-going transaction. It is advanced every time a transaction commit and can be recovered using the last commit record written in the journal.

To ensure that the Low-Watermark is safely advanced we define some rules for the checkpointing. First, at any given time there is only one checkpoint happening in the journal (this is not mandatory but is easier here). Second in the range of records that are checkpointed there must be no ongoing transaction. If there is, we either wait for the transaction to finish or we abort them forcibly. Finally, once a checkpoint is started new transactions must have a dependency higher than the Low-Watermark.

XXXXXXXXXXXXXXXXXXXXX

Before moving bounds checkpoint records are written to the Journal. The first record is a beginCheckpoint record that contains the future lower bound of the journal. Ahead of writing the record the system ensures that there are no ongoing transaction between the active and the future lower bound. If there are, we either wait for the transaction to finish or we abort them forcibly. Then the record is written in the journal. All objects are then persisted to the checkpoint store. After a checkpoint is done, the Checkpoint Store serves the new generation.

Once a checkpoint is written the store can trigger the garbage collection of unattainable objects in the store as well as deleting all the record that precede the lower bound of the journal.

Old generation: checkpoint store. New generation: bounded-size journal.

Moving the bounds in practice: ensure no ongoing transactions between the old and new lower bound; either by waiting for them to terminate, or by terminating them forcibly. Stop transactions from having dependencies inside the bounds of an ongoing checkpointing

Preview of experimental results: GC'ing obsolete key-versions.

10 IMPLEMENTATION AND EXPERIMENTAL RESULTS (SAALIK)

»(Goals:)

- Compare single variants with composed variant. Show that cache improves performance. Trade-off between write throughput (journal) and read latency (map).
- Snapshot is performed in the background.
- Can recover from a crash (journal, composed store). Recovery time is linear in the size of the log.
- Performance does not degrade over time for versions with snapshots or caches. Steady-state performance same order of magnitude as RocksDB, etc.

◀

We plan on comparing all the different store between them.

First we establish a baseline with the journal store. We write a number of random transaction, with random writes to the store. Using this we having a baseline throughput of the Journal.

Once this is done we run the same set of transactions to the checkpoint store and and the composed store to show what are the performance benefits and hits of the different features we added.

Finally we run an experiment where we add caching to the composed store and show how caching improves performance.

Objective of experiments:

- Caching improves performance.
- Truncation bounds the size of the journal.
- GC'ing the checkpoint store is effective.

## 11 RELATED WORK (EVERYONE)

### 11.1 Compiling specifications to executable code

Costa [5], Hackett et al. [8].

### 11.2 Verification of complex storage software

Chajed et al. [4], Hance et al. [9], Malecha et al. [12].

### 11.3 Formal specification of transactions and isolation models

Cerone et al. [3], Crooks et al. [6], Lesani et al. [11]

## 12 LESSONS LEARNED AND FUTURE WORK (MARC)

XXXX

Explicit modeling of sites.

Sharded store.

LSM Tree.

Update in place.

## ACKNOWLEDGMENTS

This research was funded in part by Agence Nationale de la Recherche (ANR) under grant ANR-19-CE25-0007 (<https://www.irif.fr/~gio/adecods/index.xhtml>). Gustavo: acknowledgment to ARM?

## REFERENCES

- [1] Deepthi Devaki Akkooorath, Alejandro Z. Tomsic, Manuel Bravo, Zhongmiao Li, Tyler Crain, Annette Bieniusa, Nuno Preguiça, and Marc Shapiro. 2016. Cure: Strong semantics meets high availability and low latency. In *Int. Conf. on Distributed Comp. Sys. (ICDCS)*. Nara, Japan, 405–414. <https://doi.org/10.1109/ICDCS.2016.98>
- [2] Sebastian Burckhardt, Alexey Gotsman, Hongseok Yang, and Marek Zawirski. 2014. Replicated Data Types: Specification, Verification, Optimality. In *Symp. on Principles of Prog. Lang. (POPL)*. San Diego, CA, USA, 271–284. <https://doi.org/10.1145/2535838.2535848>
- [3] Andrea Cerone, Giovanni Bernardi, and Alexey Gotsman. 2015. A Framework for Transactional Consistency Models with Atomic Visibility. In *Int. Conf. on Concurrency Theory (CONCUR) (Leibniz Int. Proc. in Informatics (LIPIcs))*, Luca Aceto and David de Frutos Escrig (Eds.), Vol. 42. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, Dagstuhl, Germany, 58–71. <https://doi.org/10.4230/LIPIcs.CONCUR.2015.58>
- [4] Tej Chajed, Joseph Tassarotti, Mark Theng, M. Frans Kaashoek, and Nickolai Zeldovich. 2022. Verifying the DaisyNFS concurrent and crash-safe file system with sequential reasoning. In *Symp. on Op. Sys. Design and Implementation (OSDI)*. Usenix, Carlsbad, CA, USA, 447–463. <https://www.usenix.org/conference/osdi22/presentation/chajed>
- [5] Renato Mascarenhas Costa. 2019. *Compiling distributed system specifications into implementations*. Master’s thesis. U. of British Columbia, Vancouver, BC, Canada. <https://doi.org/10.14288/1.0378287>
- [6] Natacha Crooks, Youer Pu, Lorenzo Alvisi, and Allen Clement. 2017. Seeing is Believing: A Client-Centric Specification of Database Isolation. In *Symp. on Principles of Dist. Comp. (PODC) (PODC ’17)*. Assoc. for Computing Machinery, Washington, DC, USA, 73–82. <https://doi.org/10.1145/3087801.3087802>
- [7] C. J. Fidge. 1988. Timestamps in message-passing systems that preserve the partial ordering. In *11th Australian Computer Science Conference*. University of Queensland, Australia, 55–66.
- [8] Finn Hackett, Shayan Hosseini, Renato Costa, Matthew Do, and Ivan Beschastnikh. 2023. Compiling Distributed System Models with PGo. In *Int. Conf. on Archi. Support for Prog. Lang. and Systems (ASPLOS) (ASPLOS 2023)*. Assoc. for Computing Machinery, Vancouver, BC, Canada, 159–175. <https://doi.org/10.1145/3575693.3575695>
- [9] Travis Hance, Yi Zhou, Andrea Lattuada, Reto Achermann, Alex Conway, Ryan Stutsman, Gerd Zellweger, Chris Hawblitzel, Jon Howell, and Bryan Parno. 2023. Sharding the State Machine: Automated Modular Reasoning for Complex Concurrent Systems. In *Symp. on Op. Sys. Design and Implementation (OSDI)*. Usenix, Boston, MA, USA, TBD. <https://doi.org/TBD>
- [10] Paul R. Johnson and Robert H. Thomas. 1976. *The maintenance of duplicate databases*. Internet Request for Comments RFC 677. Information Sciences Institute. <http://www.rfc-editor.org/rfc.html>

- [11] Mohsen Lesani, Li-yao Xia, Anders Kaseorg, Christian J. Bell, Adam Chlipala, Benjamin C. Pierce, and Steve Zdancewic. 2022. C4: Verified Transactional Objects. *Proc. ACM Program. Lang.* 6, OOPSLA1 (apr 2022). <https://doi.org/10.1145/3527324>
- [12] Gregory Malecha, Greg Morrisett, Avraham Shinnar, and Ryan Wisnesky. 2010. Toward a Verified Relational Database Management System. In *Symp. on Principles of Prog. Lang. (POPL) (POPL '10)*. Assoc. for Computing Machinery, Madrid, Spain, 237–248. <https://doi.org/10.1145/1706299.1706329>
- [13] Friedmann Mattern. 1989. Virtual Time and Global States of Distributed Systems. In *Int. W. on Parallel and Distributed Algorithms*. Elsevier Science Publishers B.V. (North-Holland), 215–226.
- [14] Patrick O’Neil, Edward Cheng, Dieter Gawlick, and Elizabeth O’Neil. 1996. The log-structured merge-tree (LSM-tree). *Acta Informatica* 33, 4 (June 1996), 351–385. <https://doi.org/10.1007/s002360050048>
- [15] Masoud Saeida Ardekani, Pierre Sutra, and Marc Shapiro. 2013. Non-Monotonic Snapshot Isolation: scalable and strong consistency for geo-replicated transactional systems. In *Symp. on Reliable Dist. Sys. (SRDS)*. IEEE Comp. Society, Braga, Portugal, 163–172. <https://doi.org/10.1109/SRDS.2013.25>
- [16] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. 2011. *A comprehensive study of Convergent and Commutative Replicated Data Types*. Rapport de Recherche 7506. Institut National de la Recherche en Informatique et Automatique (Inria), Rocquencourt, France.
- [17] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. 2011. Conflict-free Replicated Data Types. In *Int. Symp. on Stabilization, Safety, and Security of Dist. Sys. (SSS) (Lecture Notes in Comp. Sc. (LNCS))*, Xavier Défago, Franck Petit, and V. Villain (Eds.), Vol. 6976. Springer-Verlag, Grenoble, France, 386–400. [https://doi.org/10.1007/978-3-642-24550-3\\_29](https://doi.org/10.1007/978-3-642-24550-3_29)
- [18] Yair Sovran, Russell Power, Marcos K. Aguilera, and Jinyang Li. 2011. Transactional storage for geo-replicated systems. In *Symp. on Op. Sys. Principles (SOSP)*. Assoc. for Computing Machinery, Cascais, Portugal, 385–400. <https://doi.org/10.1145/2043556.2043592>
- [19] Ilyas Toumlilt, Pierre Sutra, and Marc Shapiro. 2021. Highly-available and consistent group collaboration at the edge with Colony. In *Int. Conf. on Middleware (MIDDLEWARE)*. ACM/IFIP, Québec, Canada (online), 336–351. <https://doi.org/10.1145/3464298.3493405>