



Formalizing and Checking Multilevel Consistency

Ahmed Bouajjani¹, Constantin Enea¹, Madhavan Mukund^{2,3},
Gautham Shenoy R.², and S. P. Suresh^{2,3}(✉)

¹ Université Paris Diderot, Paris, France
{abou,cenea}@irif.fr

² Chennai Mathematical Institute, Chennai, India
{madhavan,gautshen,spsuresh}@cmi.ac.in

³ CNRS UMI 2000 ReLaX, Chennai, India

Abstract. Developers of distributed data-stores must trade consistency for performance and availability. Such systems may in fact implement weak consistency models, e.g., causal consistency or eventual consistency, corresponding to different costs and guarantees to the clients. We consider the case of distributed systems that offer not just one level of consistency but multiple levels of consistency to the clients. This corresponds to many practical situations. For instance, popular data-stores such as Amazon DynamoDB and Apache’s Cassandra allow applications to tag each query within the same session with a separate consistency level. In this paper, we provide a formal framework for the specification of multilevel consistency, and we address the problem of checking the conformance of a computation to such a specification. We provide a principled algorithmic approach to this problem and apply it to several instances of models with multilevel consistency.

1 Introduction

To achieve availability and scalability, modern data-stores (key-value stores) rely on optimistic replication, allowing multiple clients to issue operations on shared data on a number of replicas, which communicate changes to each other using message passing. One benefit of such architectures is that the replicas remain locally available to clients even when network connections fail. Unfortunately, the famous CAP theorem [15] shows that such high Availability and tolerance to network Partitions are incompatible with strong Consistency, i.e., the illusion of a single centralized replica handling all operations. For this reason, modern replicated data-stores often provide weaker forms of consistency such as eventual consistency [23] or causal consistency [19], which have been formalized only recently [6, 7, 10, 22].

Programming applications on top of weakly-consistent data-stores is difficult. Some form of synchronization is often unavoidable to preserve correctness.

Partially supported by CEFIPRA DST-Inria-CNRS Project 2014-1, AVeCSO.

© Springer Nature Switzerland AG 2020

D. Beyer and D. Zufferey (Eds.): VMCAI 2020, LNCS 11990, pp. 379–400, 2020.

https://doi.org/10.1007/978-3-030-39322-9_18

Therefore, popular data-stores such as Amazon DynamoDB and Apache's Cassandra provide different levels of consistencies, ranging from weaker forms to strong consistency. Applications can tag queries to the data-store with a suitable level of consistency depending on their needs.

Implementations of large-scale data-stores are difficult to build and test. For instance, they must account for partial failures, where some components or the network can fail and produce incomplete results. Ensuring fault-tolerance relies on intricate protocols which are difficult to design and reason about. The black-box testing framework Jepsen¹ found a remarkably large number of subtle problems in many production distributed data-stores.

Testing a data-store raises two issues: (1) deriving a suitable set of testing scenarios, e.g., faults to inject into the system and the set of operations to be executed, and (2) efficient algorithms for checking whether a given execution satisfies the considered consistency models. The Jepsen framework shows that the first issue can be solved using randomization, e.g., introducing faults at random and choosing the operations randomly. The effectiveness of this solution has been proved formally in recent work [21]. The second issue is dependent on a suitable formalization of the consistency models.

In this work, we consider the problem of specifying data-stores which provide multiple levels of consistency and derive algorithms to check whether a given execution adheres to such a multilevel consistency specification.

We build on the specification framework in [10] which formalizes consistency models using two auxiliary relations: (i) a *visibility* relation, which specifies the set of operations observed by each operation, and (ii) an *arbitration order*, which specifies the order in which concurrent operations should be viewed by all replicas. An execution is said to satisfy a consistency criterion if there exists a visibility relation and an arbitration order that obey an associated set of constraints. For the case of a data-store providing multiple levels of consistency, we consider multiple visibility relations and arbitration orders, one for each level of consistency. Then, we consider a set of formulas which specifies each consistency level in isolation, and also, how visibility relations and arbitration orders of different consistency levels are related.

Based on this formalization, we investigate the problem of checking whether a given execution satisfies a certain multilevel consistency specification. In general, this problem is known to be NP-COMPLETE [6]. However, we show that for executions where each value is written at most once to a key, this problem is polynomial time for many practically-interesting multilevel consistency specifications. Since practical data-store implementations are data-independent [24], i.e., their behaviour doesn't depend on the concrete values read or written in the transactions, it suffices to consider executions where each value is written at most once. This complexity result uses the idea of *bad patterns* introduced in [6] for the case of causal consistency. Intuitively, a bad pattern is a set of operations occurring in a particular order corresponding to a consistency violation. In this paper, we provide a *systematic methodology* for deriving bad patterns characterizing a wide range of consistency models and combinations thereof.

¹ Available at <http://jepsen.io>.

Our contributions form an effective algorithmic framework for the verification of modern data-stores providing multiple levels of consistency. To the best of our knowledge, we are the first to investigate the asymptotic complexity for such a wide class of consistency models and their combinations, despite their prevalence in practice.

The paper is organized as follows. We begin with some real-life examples of multilevel consistency. In Sect. 3, we present a formal model for specifying and reasoning about multilevel consistency. Section 4 describes algorithms for verifying multilevel consistency. We conclude with a discussion of related work. Detailed proofs can be found in the technical report [8].

2 Multilevel Consistency in the Wild

In this section we present some real-world instances of multilevel consistency. We restrict our attention to distributed read-write key-value data-stores (henceforth referred to as read-write stores), consisting of unique memory locations addressed by *keys* or *variables*. We use *keys* and *variables* interchangeably in this work. The contents of these memory locations come from a domain, called *values*.

The read-write data-store provides two APIs to access and modify the contents of a particular memory location. The API to read the content of a particular memory location is typically named *Read* or *Get*, and the API to store a value into a particular memory location is typically named *Write* or *Put*. In this paper, we refer to these two methods as **Read** and **Write** respectively. The **Read** method does not update the state of the data-store and only reveals part of the state to the application session which invokes the method. The **Write** method on the other hand modifies the state of the data-store.

Typically, an application reads a location of the data-store, performs some local computation and writes a value back to the data-store. A sequence of related read and write operations performed by an application is called a *session*.

Applications expect some sort of consistency guarantee from the data-store in terms of how *fresh* or *stale* the data value is that they read from the data-store. They also seek some guarantees pertaining to monotonicity of the results that are presented to them. These guarantees provided by the data-store to the applications are called *consistency criterion*. Some of the popular consistency criteria include:

- **Read-Your-Writes:** The effects of prior operations in the session will be visible to later operations in the same session.
- **Monotonic Reads:** Once the effect of an operation becomes visible within a session, it remains visible to all subsequent operations in that session.
- **Monotonic Writes:** If the effect of a remote operation is visible in a session, then the effects of all prior operations in the session of the remote operation are also visible.
- **Causal consistency:** Effects of prior operations in a session are always visible to later operations. Further, if the effect of an operation is visible to

another operation, then every operation that has seen the effects of the latter would have seen the effects of the former.

- **Sequential Consistency:** Effects of the operations can be explained from a single sequential execution obtained by interleaving the reads and writes performed at individual sessions.

Most of the existing literature on testing the behaviour of read-write stores focuses on testing the correctness with respect to specific consistency criteria [6, 7, 14]. However, there are cases where data-stores such as DynamoDB and Cassandra offer to applications the choice of specifying the consistency level per read-operation [11]. There are distributed data-store libraries that allow consistency rationing [18] and also allow incremental consistency guarantees for the read operations [17]. In each of these cases we need to reason about the correctness of the behaviour of the data-store with respect to more than one consistency criterion.

We now look at some examples of multilevel consistency in the real world. In this work, we assume that the Read and the Write APIs are as follows.

Definition 1 (Read and Write APIs). *Let x be a key/variable, val denote a value read-from/written-to the data-store and $level$ denote the consistency level.*

- **Write(x, val)** : *Updates the content of the memory location addressed by the key/variable x with the value val .*
- **Read($x, val, level$)** : *The content of the memory location whose key is x is val with respect to the consistency level $level$.*

Read-Write Stores with Strong and Weak Reads

Consider the case of the data-store Cassandra, which allows the application a more fine grained choice of consistency levels, such as **ANY**, **ONE**, **QUORUM**, **ALL**. It achieves this by ensuring that when the Read is executed with **ANY**, the return value is provided by consulting any available replica of the data store. Similarly, if the Read operation is submitted with **ONE**, the return value is provided by consulting a replica that is known to contain at least one value for that key. On the other hand, if the Read is executed with **QUORUM**, the data-store returns the value after consulting a majority of the replicas. Finally, if Read is executed with **ALL**, then all the replicas are consulted before returning the response. Clearly, **ANY** is the weakest consistency criterion while **ALL** is the strongest consistency criterion. In general, a data-store offers responses pertaining to different consistency criteria by consulting the required subset of replicas to answer the query.

Typically a read operation under the stronger consistency criterion will take more time, since it might have to wait for all pending operations to become visible, or run a consensus protocol before returning the result. In certain cases, applications may be satisfied with Read operations that return values that are correct with respect to some weaker consistency criterion. Consider a web-application that displays the available seats in a movie theater. The application can choose to read the available seats based on a weaker consistency criterion, since:

- The number of users attempting to book seats is usually more than the seats available. Waiting for a consensus or a quorum can slow down the reads for everyone. So a quicker response is desirable.
- There is a lag between the time the user gets to see available seats and the time when the user decides to book particular seats. Since concurrent bookings are ongoing, the data displayed can become stale by the time the user books the seat.
- Users can change their minds before finally settling on a set of seats, and paying for them.

Thus, the web-application can opt for a read satisfying a weaker consistency criterion while allowing the user to pick a seat, and then perform a read satisfying a stronger consistency criterion only when the user pays for it.

Consider the example in Fig. 1 where all write requests are processed at the same replica. For each session, there is a (potentially different) designated replica from which the responses to the weak reads are returned.

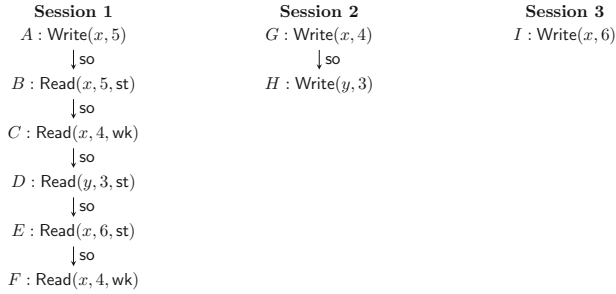


Fig. 1. An example of a read-write store behaviour with strong and weak reads. The *so* relation relates read and write operations from the same session in the order in which they happened in that session.

In this scenario, the strong reads (corresponding to the consistency level **ALL**) satisfy sequential consistency while the weak reads obey monotonic reads consistency. Hence, the fragment consisting of all the writes and the weak reads should be correct with respect to monotonic reads. Similarly, the fragment consisting of all the writes and the strong reads should be correct with respect to sequential consistency.

The weak fragment corresponding to the example in Fig. 1 can be seen in Fig. 2(a). This fragment is correct with respect to monotonic reads; once the write *G* is visible at session 1 to the read *C*, it remains visible throughout the session. The write *I* is not visible to any of the other sessions yet.

The strong fragment is represented in Fig. 2(b). This is correct with respect to sequential consistency, where the order of the operations obtained by consensus is $A \rightarrow B \rightarrow G \rightarrow H \rightarrow I \rightarrow D \rightarrow E$.

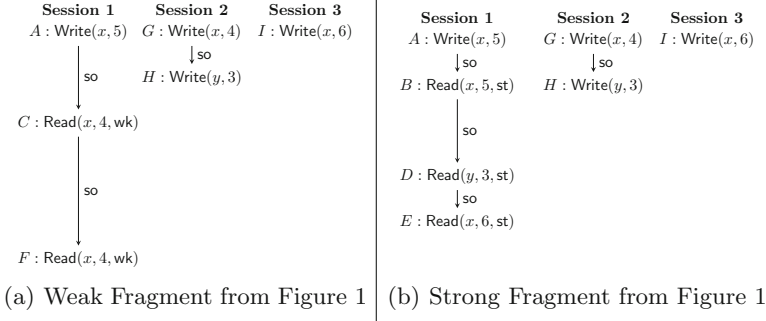


Fig. 2. Strong and Weak fragments of the hybrid behaviour

However, since the strong reads correspond to the level **ALL** where all the replicas have seen the prior writes and have agreed on the order of the concurrent writes, it behooves a weak read following a strong read to take into consideration the effects seen by the earlier strong read. Thus the data-store imposes an additional constraint. Once a write is visible to a strong read in a session, it is visible to all the subsequent weak reads in that session. This ensures that the weaker reads do incorporate the prior results seen by the session. Similarly, a write visible to a weak read is made from a replica which participates in the subsequent strong reads corresponding to the level **ALL**. Thus the effects visible to the prior weak reads in a session are also visible to the subsequent strong reads.

With these additional constraints, we can no longer explain the read operation F , since the effects of writes G and I are both visible at read F . The strong consistency criterion has already guaranteed that write I has happened after write G , thereby effectively overwriting the value 4 with the value 6. Hence this behaviour is incorrect in the multilevel setting.

Now consider the behaviour of Cassandra where writes are performed at one of the replicas (corresponds to the level **ONE**), weak reads are performed at one of the replicas (corresponds to the level **ONE**) and strong reads are performed at a quorum of replicas (corresponds to the level **QUORUM**). In this situation, it is not necessary that the effects of writes visible to prior weaker reads are visible at subsequent stronger reads, since the replica from which the weaker read is performed may be missing from the quorum of replicas from which the stronger read is made. Similarly, the effects of writes visible to prior stronger reads of a session need not be visible to the subsequent weaker reads in the session, as the writes from the quorum may not have reached the replica from which the weaker read is performed. Thus, the stronger and weaker reads can be independent of each other.

Finally consider the case of Amazon DynamoDB Accelerator (DAX) [1], which contains a write-through cache sitting between the application and the DynamoDB backend. Every write made by the application is first submitted to

the DynamoDB backend and also updated at the cache. By default, the reads are eventually consistent, i.e., the reads are performed from the cache. If the item does not exist in the cache, then it is fetched from the backend data-store and the cache is updated with the item before the value is returned to the application. However, the application can also request strongly consistent reads by invoking `ConsistentRead`. In this case, the value is read from the backend and returned to the application, without caching the results. Any subsequent eventually consistent reads made by the application may not reflect the value returned by the prior strongly consistent read. In the case of DAX, it can be observed that the effects of the writes visible to the weak eventually consistent reads are also visible to the subsequent strongly consistent reads as those writes are also present in the DynamoDB backend. However, it is not necessary that the effects of writes visible to the strongly consistent reads are visible to the subsequent weak eventually consistent reads.

From these examples of multilevel consistency, we can see that the presence of another consistency criterion can impose additional constraints on the choice of the visibility and arbitration relations chosen to explain the correctness of the history. In the next section, we provide a formal framework for modelling behaviours of read-write data-stores with multiple consistency levels.

3 Formalizing Multilevel Consistency

We extend the formal framework provided in [9] for modelling the behaviours of read-write stores. Each operation submitted to the data-store by the application is either a `Read` or a `Write` operation whose signature is given in Definition 1.

We denote the set of all variables in the read-write store by $Vars$ and assume that each value written to the read-write store is a natural number $val \in \mathbb{N}$. We assume that all variables are initially undefined, with value \perp .

For simplicity, we assume only two consistency levels, weak and strong, denoted by `wk` and `st`, respectively, where the consistency criterion corresponding to `wk`-level is strictly weaker than then the consistency criterion corresponding to the `st`-level. Comparison between consistency criteria is formally defined in Definition 7.

The behaviour of the read-write data-store as observed by an application is the sequence of reads and writes that it performs on the stores. The sequence of related read and write operations is termed a *session*. Thus the behaviour of the read-write store seen by each session is a total order of read/write operations performed in that session.

The behaviour of the read-write store is the collection of behaviours seen by all the sessions. In Fig. 1 we saw the behaviour of the data-store as observed by the three sessions accessing the data-store. We call such a behaviour a *hybrid history*, formally defined as follows:

Definition 2 (Hybrid History). *A hybrid history of a read-write store is a pair $H = (\mathcal{O}, so)$ where \mathcal{O} is the set of read-write operations and so is a collection of total orders called session orders.*

For a history H , we define the following subsets of \mathcal{O} .

- $\mathcal{O}_{\text{Read}}$ is the set of read operations occurring in H .
- $\mathcal{O}_{\text{Write}}$ is the set of write operations occurring in H .
- $\mathcal{O}_{\text{wk}} = \mathcal{O}_{\text{Write}} \cup \{\text{Read}(x, \text{val}, \text{level}) \in \mathcal{O}_{\text{Read}} \mid \text{level} = \text{wk}\}$ (the set of weak operations occurring in H).
- $\mathcal{O}_{\text{st}} = \mathcal{O}_{\text{Write}} \cup \{\text{Read}(x, \text{val}, \text{level}) \in \mathcal{O}_{\text{Read}} \mid \text{level} = \text{st}\}$ (the set of strong operations occurring in H).

The weak fragment of the history H is denoted H_{wk} and defined to be $(\mathcal{O}_{\text{wk}}, \text{so} \cap (\mathcal{O}_{\text{wk}} \times \mathcal{O}_{\text{wk}}))$. Similarly the strong fragment of the history H is denoted H_{st} and is defined to be $(\mathcal{O}_{\text{st}}, \text{so} \cap (\mathcal{O}_{\text{st}} \times \mathcal{O}_{\text{st}}))$. Note that we take the write operations to be part of both the strong and weak fragments.

- For $X \subseteq \mathcal{O} \times \mathcal{O}$ and $\ell \in \{\text{Read}, \text{Write}, \text{wk}, \text{st}\}$, $X \upharpoonright_{\ell} = X \cap (\mathcal{O}_{\ell} \times \mathcal{O}_{\ell})$.
- For $X, Y \subseteq \mathcal{O} \times \mathcal{O}$, $X; Y = \{(x, y) \mid \exists z : (x, z) \in X \text{ and } (z, y) \in Y\}$ is the composition of X and Y .
- For $X \subseteq \mathcal{O} \times \mathcal{O}$, $\text{total}(X)$ is used to mean that X is a total order.

When a replica of the read-write store receives an operation from an application, it decides how the effects of the older operations known to the replica (either received from applications, or from other replicas of the data-store) should be made visible to the new operation. A visibility relation over a history specifies the set of operations visible to an operation.

Definition 3 (Visibility Relation). A visibility relation vis over a history $H = (\mathcal{O}, \text{so})$ is an acyclic relation over \mathcal{O} . For $o, o' \in \mathcal{O}$, we write $o \xrightarrow{\text{vis}} o'$ to indicate that the effects of the operation o are visible to the operation o' .

If a pair of operations o, o' are not related by vis , we term them concurrent operations, denoted by $o \parallel_{\text{vis}} o'$.

We define the view of an operation o with respect to a visibility relation vis , denoted $\text{View}_{\text{vis}}(o)$ to be the set of all the Write operations visible to it.

For the history in Fig. 1, we can define a visibility relation to be

$$\{A \xrightarrow{\text{vis}} B, G \xrightarrow{\text{vis}} C, G \xrightarrow{\text{vis}} D, H \xrightarrow{\text{vis}} D, G \xrightarrow{\text{vis}} E, H \xrightarrow{\text{vis}} E, I \xrightarrow{\text{vis}} E, G \xrightarrow{\text{vis}} F\}$$

When the replicas communicate with each other, they need to reconcile the effects of concurrent write operations in order to converge to the same state eventually. In case of convergent data-stores this is done using a rule such as *Last Writer Wins* which totally orders all write operations. This is abstracted by an arbitration relation, which is a total order over all write operations in the history. We will denote the arbitration relation by arb . We assume that the arbitration relation is consistent with the visibility relation, in the sense that for a pair of writes o and o' , if o is visible to o' then o is before o' in arb .

Definition 4 (Arbitration Relation). An arbitration relation arb over a hybrid history $H = (\mathcal{O}, \text{so})$ is a total order over $\mathcal{O}_{\text{Write}}$. For $o_i, o_j \in \mathcal{O}$, we say $o_i \xrightarrow{\text{arb}} o_j$ to indicate that operation o_i has been ordered before the operation operation o_j .

For the history in Fig. 1 the arbitration relation can be the total order

$$A \xrightarrow{\text{arb}} G \xrightarrow{\text{arb}} H \xrightarrow{\text{arb}} I$$

We define the correctness of a hybrid history in terms of the *functional specification* of read-write stores.

Let H be a hybrid history. Let vis and arb be visibility and arbitration relations over H .

We say that a write operation o' is a *related-write* of a read operation o iff o' is in the view of o and both o and o' operate on the same variable. The set of all related writes of o , denoted $\text{RelWrites}_{\text{vis}}(o)$, is defined to be the set $\{o' \in \text{View}_{\text{vis}}(o) \mid o \text{ and } o' \text{ operate on the same variable}\}$.

$\text{MaxRelWrites}_{\text{vis}}(o)$, the set of maximal elements among these related writes with respect to vis , is defined to be

$$\{o' \in \text{RelWrites}_{\text{vis}}(o) \mid \forall o'' \in \text{RelWrites}_{\text{vis}}(o) : o'' \xrightarrow{\text{vis}} o' \vee o'' \parallel_{\text{vis}} o'\}$$

The effective write of a read-operation o , denoted by $\text{EffWrite}_{\text{vis}}^{\text{arb}}(o)$ is defined to be the maximum write operation from the set of maximal related writes of o as per the arbitration relation.

$$\text{EffWrite}_{\text{vis}}^{\text{arb}}(o) = \begin{cases} \max(\text{arb} \upharpoonright_{\text{MaxRelWrites}_{\text{vis}}(o)}) & \text{if } \text{MaxRelWrites}_{\text{vis}}(o) \neq \emptyset \\ \perp & \text{otherwise} \end{cases}$$

Definition 5 (Functional Correctness for Read-Write Stores). Let $H = (\mathcal{O}, \text{so})$ be a hybrid history of a read-write data store with visibility relation vis and arbitration relation arb . We say that $(H, \text{vis}, \text{arb})$ is functionally correct iff for every read operation $o = \text{Read}(x, \text{val}, \text{level})$, the following conditions hold.

- $\text{EffWrite}_{\text{vis}}^{\text{arb}}(o) = \perp$ iff $\text{val} = \perp$ (i.e., there was no write operation on x when o happened).
- If $o' = \text{EffWrite}_{\text{vis}}^{\text{arb}}(o)$ then o' wrote the value val .

Next, we formally define *consistency criteria* in terms of a set of formulas. Our definition is adapted from the definitions of constraints in [13].

Definition 6 (Consistency Criteria). A relation term τ is a composition of the form $r_1; \dots; r_k$ ($k \geq 1$), where each $r_i \in \{\text{so}, \text{vis}\}$. A consistency criterion is a subset of

$$\{\tau \subseteq \text{vis} \mid \tau \text{ is a relation term}\} \cup \{\text{total}(\text{vis})\}.$$

Thus a consistency criterion is a possibly empty collection of visibility constraints and an optional totality constraint. For simplicity of notation, we usually write a constraint as a conjunction.

Note that so and vis are variables which are usually interpreted as restrictions of the so and vis relations in a history. As we will see below, we always require an additional constraint that $\text{vis} \upharpoonright_{\text{Write}} \subseteq \text{arb}$ (and hence it is not explicitly included in the consistency criteria).

For a consistency criterion α , $RelTerms(\alpha)$ is the set of all relation terms occurring in α , and $VisBasic(\alpha)$ is the collection of all visibility constraints in α excluding the totality constraint $total(vis)$.

Definition 7 (Consistency Criterion in a history). Let $H = (\mathcal{O}, so)$ be a hybrid history, let vis and arb be a visibility and arbitration relation over H , and let α be a consistency criterion. We say that $H, vis \models \alpha$ iff:

1. for every $\tau \subseteq vis$ in α , $\tau[so := so, vis := vis] \subseteq vis$, and
2. if $total(vis) \in \alpha$, then $total(vis)$ holds.

Further we say that $H, vis, arb \models \alpha$ iff $H, vis \models \alpha$ and $vis \downarrow_{Write} \subseteq arb$.

Some well known consistency criteria are given in Table 1.

Table 1. Well known consistency criteria

Name	Description
<i>Basic Eventual Consistency (BEC)</i>	\top
<i>Read Your Writes (RYW)</i>	$so \subseteq vis$
<i>Monotonic Reads (MR)</i>	$vis; so \subseteq vis$
<i>Monotonic Writes (MW)</i>	$so; vis \subseteq vis$
<i>Strong Eventual Consistency (SEC)</i>	$so \subseteq vis \wedge vis; so \subseteq vis$
<i>FIFO Consistency (FIFO)</i>	$so \subseteq vis \wedge vis; so \subseteq vis \wedge so; vis \subseteq vis$
<i>Causal Consistency (CC)</i>	$so \subseteq vis \wedge vis; vis \subseteq vis$
<i>Sequential Consistency (SEQ)</i>	$so \subseteq vis \wedge vis; vis \subseteq vis \wedge total(vis)$

We say that a consistency criterion α is at least as *strong* as another consistency criterion α' if for every history H , visibility relation vis , and arbitration relation arb over H , if $H, vis, arb \models \alpha$ then $H, vis, arb \models \alpha'$.

Suppose $H = (\mathcal{O}, so)$ is a hybrid history. Let α_w and α_s respectively be the *wk* and *st* consistency criteria. We then want to choose *wk* and *st* visibility relations vis_{wk} , vis_{st} , respectively, and an arbitration relations arb such that $H_{wk}, vis_{wk}, arb \models \alpha_w$ and $H_{st}, vis_{st}, arb \models \alpha_s$.

As we had noted in the previous section, in a multilevel setting, it is not sufficient to separately satisfy the constraints corresponding to the *wk* and *st* consistency criteria. We now proceed to modelling multilevel consistency constraints.

Modelling Multilevel Consistency

Taking inspiration from DAX [1] and the cache-hierarchy in modern processors, we can model multilevel consistency as a series of data-stores arranged in increasing order of the consistency they guarantee, such that the data-store offering the

weakest level of consistency is closest to the application, and the data-store offering the strongest level of consistency is farthest away from the application. We shall further assume that these data-stores use the same arbitration strategy to order concurrent write operations and every weaker data-store has the capability to update its state to match that of a stronger data-store.

For the purpose of this paper, since we are restricting ourselves to only two levels, namely *wk* and *st*, this will reduce to having just two data-stores, where the data-store corresponding to the weaker consistency criterion sits as a cache between the application and the data-store corresponding to the stronger consistency criterion.

All the *wk*-reads are performed from the *wk* data-store.

There are two possible ways in which the writes can be performed.

1. **Write-Through:** The write is first performed at the *st*-data-store and eventually will be propagated to the *wk*-data-store.
2. **Write-Back:** The write is first performed at the *wk*-data-store and eventually will be propagated to the *st*-data-store.

There are two possible ways in which *st*-reads can be performed.

- (a) **Read-Through:** The result of the *st*-read performed at the *st*-data-store is directly sent to the application bypassing the *wk*-data-store.
- (b) **Read-Back:** The result of the *st*-read is updated at the *wk*-data-store before it is propagated to the application.

Thus, the system picks one of two ways to perform the write, and one of the two ways to perform the *st*-read.

Note that a system which picks the **Write-Through** strategy for performing the write will ensure that any write visible at the *wk* data-store will also be visible to the *st* data-store, as all the writes are first performed at the *st* data-store before they are propagated to the *wk* one. Hence, the effects of write operations visible to a *wk*-read operation are also visible to the subsequent *st*-operations in the session.

Similarly a system which picks the **Read-Back** strategy for performing the *st*-reads will ensure that any write that is visible to a *strong*-read will also be visible at a subsequent *wk*-read in the session as before returning the result of the *st*-read to the application, the result is merged into the *wk* data-store.

However, the **Write-Back** and **Read-Through** strategies do not provide any guarantees between the effects of writes visible to *wk* (resp. *st*) reads in relation to the subsequent *st* (resp. *wk*) reads in that session.

We now define the guarantees provided by each of these four strategies in the form of a constraint.

Definition 8 (Multilevel Constraints). *We define the following formulas:*

- $\psi_{thru}^{write} := (vis^{wk}; so) \upharpoonright_{st} \subseteq vis^{st}$
- $\psi_{back}^{write} := \top$
- $\psi_{thru}^{read} := \top$

$$- \psi_{back}^{read} := (vis^{st}, so) \downarrow_{wk} \subseteq vis^{wk}$$

A multilevel constraint φ is a conjunction of the form $\psi^{read} \wedge \psi^{write}$, where $\psi^{read} \in \{\psi_{thru}^{read}, \psi_{back}^{read}\}$ and $\psi^{write} \in \{\psi_{thru}^{write}, \psi_{back}^{write}\}$.

Suppose $H = (\mathcal{O}, so)$ is a history, and vis_{wk} and vis_{st} are two visibility relations respectively over \mathcal{O}_{wk} and \mathcal{O}_{st} . Let φ be a multilevel constraint. We say that $H, vis_{wk}, vis_{st} \models \varphi$ iff $\varphi[so := so, vis^{wk} := vis_{wk}, vis^{st} := vis_{st}]$ is true.

The formula ψ_{thru}^{write} imposes the constraint that the strong operations see the effects seen by the prior weak operations in the session. Similarly, the formula ψ_{back}^{read} imposes the constraint that the weak operations see the effects seen by the prior strong operations in the session. These two guarantee that the effect seen by reads of one consistency level remain monotonically visible to the subsequent reads of another consistency level.

Consider Cassandra's multilevel consistency with writes performed at level ONE, weak-reads at level ONE and strong-reads at level ALL which ensure that weaker reads see the effects visible to prior stronger reads and vice-versa. This can be modelled using $\psi_{thru}^{write} \wedge \psi_{back}^{read}$.

On the other hand, Cassandra's multilevel consistency with writes performed at level ONE, weak-reads at level ONE and strong-reads at level QUORUM neither ensures that weaker reads see the effects visible to prior stronger reads nor the converse. This can be modelled using $\psi_{back}^{write} \wedge \psi_{thru}^{read}$.

The DynamoDB's DAX case can be modelled using $\psi_{thru}^{write} \wedge \psi_{thru}^{read}$ which only allows for the effects of prior weak reads to be visible to subsequent stronger reads, but not the converse.

We now formally define when a hybrid history is correct.

Definition 9 (Multilevel Correctness of a Hybrid History). A hybrid history $H = (\mathcal{O}, so)$ of a read-write store is said to be multilevel correct with respect to a wk-consistency criterion α_w , st-consistency criterion α_s and multilevel consistency constraint φ , iff there exists visibility relations vis_{wk} and vis_{st} over H_{wk} and H_{st} respectively and arbitration relation arb such that

- (H_{wk}, vis_{wk}, arb) and (H_{st}, vis_{st}, arb) are functionally correct,
- $H_{wk}, vis_{wk}, arb \models \alpha_w$,
- $H_{st}, vis_{st}, arb \models \alpha_s$, and
- $H, vis_{wk}, vis_{st} \models \varphi$.

4 Testing Multilevel Correctness of a Hybrid History

Given a read-write hybrid history $H = (\mathcal{O}, so)$, we want to test it for multi-level correctness with respect to weak and strong consistency criteria α_w and α_s and multilevel constraints given by φ .

We note that for the history to be correct, for every read operation that returns a value that is not \perp , there should exist a write operation writing the same value to the variable that was read. The *reads-from* relation associates a write operation to the read that reads its effect. Our strategy for testing the

multilevel correctness of H is to enumerate all such *reads-from* relations rf , for each rf we find visibility relations vis_{wk} and vis_{st} , respectively, containing rf_{wk} and rf_{st} , such that they satisfy the visibility constraints imposed by the individual consistency criteria, as well as the multilevel constraints, i.e., $H_{\text{wk}}, \text{vis}_{\text{wk}} \models \alpha_w$, $H_{\text{st}}, \text{vis}_{\text{st}} \models \alpha_s$ and $H, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}} \models \varphi$. We then check for the presence of a finite number of *bad-patterns* in these visibility relations. The presence of a bad-pattern implies that for every arbitration relation arb , there is some level $\ell \in \{\text{wk}, \text{st}\}$ such that either the arbitration constraint $\text{vis}_{\ell} \upharpoonright_{\text{Write}} \subseteq \text{arb}$ is not satisfied, or the history $(H_{\ell}, \text{vis}_{\ell}, \text{arb})$ is not functionally correct.

If the history is multi-level correct, then we will find a witness consisting of a *reads-from* relation rf and visibility relations vis_{wk} and vis_{st} extending rf_{wk} and rf_{st} such that all the constraints are satisfied and there are no bad-patterns. If the history is not multi-level correct, then for every pair of weak and strong visibility relation extending every *reads-from* relation, either some constraint is not satisfied or there exists a bad-pattern.

We present the bad-pattern characterization for multilevel correctness of a hybrid history in the next subsection. In the following subsection, we provide a procedure for computing the minimal visibility relations vis_{wk} and vis_{st} for a given *reads-from* relation rf that satisfies α_w , α_s and φ .

4.1 Bad Pattern Characterization for Multilevel Correctness

We now characterize the correctness of hybrid histories based on the non-existence of certain bad patterns. This is a generalization of the bad-pattern characterization for causal consistency in [6].

Given a hybrid history, we can associate each Read with a unique write operation from the history whose effect the Read operation reads from. We call this the *reads-from* relation.

Definition 10 (Reads-From). A reads-from relation rf over a history $H = (\mathcal{O}, \text{so})$ is a binary relation such that

1. $(o_i, o_j) \in \text{rf} \implies o_i$ is a Write, o_j is a Read, both on the same variable, such that the value returned by o_j is the value written by o_i .
2. $(o_i, o_j) \in \text{rf} \wedge (o_k, o_j) \in \text{rf} \implies o_i = o_k$.
3. For all $o_j = \text{Read}(x, \text{val}, \text{level}) \in \mathcal{O}_{\text{Read}}$

$$[\exists o \in \mathcal{O}_{\text{Write}} \text{ which writes val to } x \implies \exists o_i \in \mathcal{O}_{\text{Write}} : (o_i, o_j) \in \text{rf}.]$$

Condition 1 associates a read operation with a write operation only if they operate on the same variable and that the return value of the read operation matches the argument of the write operation.

Condition 2 ensures that a read operation is associated with at most one write operation.

Finally, Condition 3 insists that if a Read is not related to any Write via rf , it is only because there is no matching Write in the hybrid history (i.e. a write of the same value to the same variable).

Let rf be a *reads-from* relation on a hybrid history $H = (\mathcal{O}, \text{so})$. For a Read operation $o \in \mathcal{O}$, if there exists a Write operation o' such that $(o', o) \in \text{rf}$, then we say that $\text{rf}^{-1}(o) = o'$. If no such o' exists, we set $\text{rf}^{-1}(o) = \perp$.

Further, we denote by rf_{wk} and rf_{st} the *reads-from* relation restricted to H_{wk} and H_{st} respectively.

Suppose rf_ℓ is a *reads-from* relation over H_ℓ . We say that a visibility relation vis_ℓ over H_ℓ *extends* rf_ℓ iff $\text{rf}_\ell \subseteq \text{vis}_\ell$. Suppose arb is an arbitration relation over H_ℓ . Then, we say that $(\text{vis}_\ell, \text{arb})$ *realize* rf_ℓ iff for all read operations $o \in \mathcal{O}_\ell$, $\text{rf}_\ell^{-1}(o) = \text{EffWrite}_{\text{vis}_\ell}^{\text{arb}}(o)$.

Given a *reads-from* relation rf_ℓ and a visibility relation vis_ℓ that extends it, we can define a conflict relation that orders all the remaining maximal related writes in $\text{MaxRelWrites}_{\text{vis}_\ell}(o)$ of a read-operation o before the write-operation $\text{rf}_\ell^{-1}(o)$. The conflict relation captures the essence of the arbitration relation for a given *reads-from* relation and a visibility relation extending it.

Definition 11 (Conflict Relation). Let $H_\ell = (\mathcal{O}_\ell, \text{so}_\ell)$ be a history. Let rf_ℓ be a *reads-from* relation over H_ℓ . Let $\text{vis}_\ell \supseteq \text{rf}_\ell$ be a visibility relation over H_ℓ . We define the conflict relation for rf_ℓ and vis_ℓ , denoted $\text{CF}(\text{rf}_\ell, \text{vis}_\ell)$, as the set

$$\{(o'', o') \mid \exists o \in \mathcal{O}_\ell \upharpoonright_{\text{Read}}: o'', o' \in \text{MaxRelWrites}_{\text{vis}_\ell}(o) \wedge o' = \text{rf}_\ell^{-1}(o)\}.$$

We now define the bad patterns that characterize the correctness of the hybrid history.

Definition 12 (Bad Patterns for a hybrid history). Let $H = (\mathcal{O}, \text{so})$ be a hybrid history with weak and strong consistency criteria α_w and α_s respectively and multilevel constraints φ . Let rf be a *reads-from* relation over H . For $\ell \in \{\text{wk}, \text{st}\}$, let vis_ℓ be a relation over \mathcal{O}_ℓ with $\text{vis}_\ell \supseteq \text{rf}_\ell$ such that $H_{\text{wk}}, \text{vis}_{\text{wk}} \models \alpha_w$, $H_{\text{st}}, \text{vis}_{\text{st}} \models \alpha_s$ and $H, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}} \models \varphi$. We define the following bad patterns for $(H, \text{rf}, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}})$. For some $\ell \in \{\text{wk}, \text{st}\}$:

- **BADVISIBILITY:** $\text{Cyclic}(\text{vis}_\ell)$
- **THINAIR:** $\exists o \in \mathcal{O}_{\text{Read}} \upharpoonright_\ell: o \text{ returns a value that is not } \perp, \text{ but } \text{rf}_\ell^{-1}(o) = \perp$
- **BADINITREAD:** $\exists o \in \mathcal{O}_{\text{Read}} \upharpoonright_\ell: o \text{ returns } \perp \text{ but } \text{RelWrites}_{\text{vis}_\ell}(o) \neq \emptyset$
- **BADREAD:** $\exists o \in \mathcal{O}_{\text{Read}} \upharpoonright_\ell: \text{rf}_\ell^{-1}(o) \notin \text{MaxRelWrites}_{\text{vis}_\ell}(o)$
- **BADARB:** $\text{Cyclic}(\bigcup_{\ell \in \{\text{wk}, \text{st}\}} (\text{CF}(\text{rf}_\ell, \text{vis}_\ell) \cup (\text{vis}_\ell)_{\text{Write}}))$

BADVISIBILITY says that one of the visibility relations has a cycle.

THINAIR says that there exists a read in the history which reads a non-initial value which is not written by any write operation in the hybrid history.

BADINITREAD says that there is a read operation on a variable which reads the initial value despite having a non-initial write to that variable in its view.

BADREAD says that the write operation from which the read-operation reads is not a maximal write, and there are other writes in the view of the read operation that would have overwritten the value written by that write.

BADARB says that the union of the conflict relations along visibility relation restricted to only the Write operations has a cycle indicating that there exists no total-order arb over $\mathcal{O}_{\text{Write}}$, such that $(\text{vis}_\ell, \text{arb})$ realizes rf_ℓ .

Multi-level correctness of a hybrid history can be characterized in terms of non-existence of these bad patterns. The proof can be found in the full version of this paper [8].

Theorem 1 (Bad patterns characterization). *A hybrid history $H = (\mathcal{O}, \text{so})$ is said to be multilevel correct with respect to weak and strong consistency criteria α_w , α_s and multilevel constraint φ iff there exists a reads-from relation rf , and relations $\text{vis}_{\text{wk}} \supseteq \text{rf}_{\text{wk}}$ and $\text{vis}_{\text{st}} \supseteq \text{rf}_{\text{st}}$ respectively over \mathcal{O}_{wk} and \mathcal{O}_{st} such that $H_{\text{wk}}, \text{vis}_{\text{wk}} \models \alpha_w$, $H_{\text{st}}, \text{vis}_{\text{st}} \models \alpha_s$ and $H, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}} \models \varphi$ and no bad pattern exists in $(H, \text{rf}, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}})$.*

4.2 Constructing Minimal Visibility Relations

Suppose $H = (\mathcal{O}, \text{so})$ is a hybrid history. Let α_w and α_s be the formulas defining the weak and strong consistency criteria, and let φ be the formula defining the multilevel constraints. Let $\alpha'_w = \text{VisBasic}(\alpha_w)$ and $\alpha'_s = \text{VisBasic}(\alpha_s)$.

We provide a procedure that iterates over all the possible *reads-from* relations and constructs a minimal visibility relation extending the *reads-from* relation such that it satisfies α_w , α_s and φ . The pseudo-code for the procedure is presented in Algorithm 1 and 2.

Algorithm 1 Constructing minimal visibility relations

<pre> 1 MinVisOne($\mathcal{O}_\ell, \text{so}_\ell, \text{vis}_\ell, \alpha_\ell$): 2 Let $\text{vis}_o := \text{vis}_\ell$; 3 4 while (True): 5 Let $\text{vis}_p := \text{vis}_o$; 6 for $\tau \in \text{RelTerms}(\alpha_\ell)$: 7 $\text{vis}_n := \text{vis}_p \cup \tau[\text{so}_\ell, \text{vis}_p]$; 8 $\text{vis}_p := \text{vis}_n$; 9 if ($\text{vis}_n == \text{vis}_o$) 10 return vis_n 11 $\text{vis}_o := \text{vis}_n$ 12 13 ComputeVisSet($\mathcal{O}_\ell, \text{so}_\ell, \text{vis}_\ell, \alpha_\ell$) 14 if total($\text{vis}$) is a subformula in α_ℓ: 15 $\text{visSet}_\ell := \{\text{totvis} \mid \text{totvis}$ 16 is a 17 total order over 18 \mathcal{O}_ℓ such that 19 $\text{vis}_\ell \subseteq \text{totvis}\}$ 20 else : 21 $\text{visSet}_\ell := \{\text{vis}_\ell\}$ 22 return visSet_ℓ </pre>	<pre> 21 MinVisMulti($\mathcal{O}, \text{so}, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}}, \psi$) 22 if $\psi \in \{\psi_{\text{back}}^{\text{write}}, \psi_{\text{thru}}^{\text{read}}\}$: 23 return ($\text{vis}_{\text{wk}}, \text{vis}_{\text{st}}$) 24 else if $\psi \in \{\psi_{\text{thru}}^{\text{write}}\}$: 25 Let $\ell = \text{st}, \ell' = \text{wk}$; 26 else if $\psi \in \{\psi_{\text{back}}^{\text{read}}\}$: 27 Let $\ell = \text{wk}, \ell' = \text{st}$; 28 29 Let $\text{vis}_\ell^o := \text{vis}_\ell$; 30 Let $\text{vis}_{\ell'}^o := \text{vis}_{\ell'}$; 31 32 if $\psi \in \{\psi_{\text{thru}}^{\text{write}}, \psi_{\text{back}}^{\text{read}}\}$: 33 Let $\text{vis}_\ell^n := \text{vis}_\ell^o \cup (\text{vis}_{\ell'}^o; \text{so}) \upharpoonright_\ell$; 34 Let $\text{vis}_{\ell'}^n := \text{vis}_{\ell'}^o$; 35 36 if $\psi \in \{\psi_{\text{back}}^{\text{read}}\}$: 37 return ($\text{vis}_\ell^n, \text{vis}_{\ell'}^n$) 38 else if $\psi \in \{\psi_{\text{thru}}^{\text{write}}\}$: 39 return ($\text{vis}_{\ell'}^n, \text{vis}_\ell^n$) 40 41 42 </pre>
--	---

In Lines 1–12 we have a method `MinVisOne` that takes as input a visibility relation vis_ℓ for the history $(\mathcal{O}_\ell, \text{so}_\ell)$ and constructs an extension vis_n that

satisfies the formula $VisBasic(\alpha_\ell)$. We achieve this by iterating over the $RelTerms$ appearing in $RelTerms(\alpha_\ell)$ (Line 6) and extending the previous visibility relation vis_p with the evaluation of the term (Line 7). We do this until we obtain a relation vis_n which we can no longer extend (Line 9). This final visibility relation vis_n extends vis_ℓ and satisfies the formula $VisBasic(\alpha_\ell)$.

In Lines 21–40, we have the procedure `MinVisMulti` which takes as inputs the hybrid history (\mathcal{O}, so) , visibility relations vis_{wk} and vis_{st} and an individual conjunct ψ appearing in the multilevel constraint φ . Since every visibility relation trivially satisfies ψ_{back}^{write} or ψ_{thru}^{read} , for these multilevel constraint, we simply return without modifying vis_{wk} or vis_{st} (Lines 22–23). In the remaining cases, when the multi-level constraint is either ψ_{back}^{write} or ψ_{thru}^{read} , for $\ell, \ell' \in \{wk, st\}$, the multilevel constraints relates the write operations visible to the operations of level ℓ in terms of the writes seen by operations of level ℓ' that have occurred previously in the session. Depending on the conjunct ψ , we set ℓ and ℓ' appropriately (Lines 24–27). We then extend the visibility relation for level ℓ by relating each ℓ -operation to the Writes that have been seen by any of the ℓ' -operations prior to the ℓ -operation in its session (Line 33). The visibility relation for level ℓ' remains unchanged in this case (Line 34).

We return these extended visibility relations as a pair, where the wk visibility extension is followed by st visibility extension (Lines 36–39).

Algorithm 2 Testing multilevel correctness of a hybrid history

<pre> 43 ComputeStableExt($\mathcal{O}, so, vis_{wk}, vis_{st}, \alpha_w, \alpha_s, \varphi$): 44 Let $vis_{wk}^o := vis_{wk}$, $vis_{st}^o := vis_{st}$ 45 46 while (True): 47 Let $vis_{wk}^p := vis_{wk}^o$, $vis_{st}^p := vis_{st}^o$ 48 49 Let $vis_{wk}^n :=$ MinVisOne($\mathcal{O}_{wk}, so_{wk}, vis_{wk}^p, \alpha_w$); 50 51 Let $vis_{st}^n :=$ MinVisOne($\mathcal{O}_{st}, so_{st}, vis_{st}^p, \alpha_s$); 52 53 for each subformula ψ_i 54 in the conjunction φ: 55 $vis_{wk}^p := vis_{wk}^n, vis_{st}^p := vis_{st}^n$ 56 57 (vis_{wk}^n, vis_{st}^n) = MinVisMulti($\mathcal{O}, so, vis_{wk}^p, vis_{st}^p, \psi_i$) 58 59 if $vis_{wk}^n = vis_{wk}^o$ and $vis_{st}^n = vis_{st}^o$: 60 return (vis_{wk}^n, vis_{st}^n) 61 62 $vis_{wk}^o := vis_{wk}^n, vis_{st}^o := vis_{st}^n$ </pre>	<pre> 63 TestMultiCorrect($\mathcal{O}, so, \alpha_w, \alpha_s, \varphi$): 64 Let rfSet := {rf rf is a reads-from relation over (\mathcal{O}, so)} 65 66 for rf \in rfSet: 67 Let $vis_{wk}^{min} :=$ MinVisOne($\mathcal{O}_{wk}, so_{wk}, rf_{wk}, \alpha_w$); 68 Let $vis_{st}^{min} :=$ ComputeVisSet($\mathcal{O}_{wk}, so_{wk}, vis_{wk}^{min}, \alpha_w$); 69 70 Let $vis_{st}^{min} :=$ MinVisOne($\mathcal{O}_{wk}, so_{st}, rf_{st}, \alpha_s$); 71 Let $vis_{st}^{min} :=$ ComputeVisSet($\mathcal{O}_{st}, so_{st}, vis_{st}^{min}, \alpha_s$); 72 73 for $vis_{wk} \in visSet_{wk}, vis_{st} \in visSet_{st}$: 74 Let ($vis_{wk}^{stb}, vis_{st}^{stb}$) := ComputeStableExt($\mathcal{O}, so, vis_{wk},$ $vis_{st}, \alpha_w, \alpha_s, \varphi$); 75 76 if NoBadPatterns($\mathcal{O}, so, rf,$ $vis_{wk}^{stb}, vis_{st}^{stb}$): 77 return (rf, $vis_{wk}^{stb}, vis_{st}^{stb}$) 78 79 return BadHistory </pre>
--	---

In Lines 43–61 we have the procedure `ComputeStableExt` which takes history (\mathcal{O}, so) a pair of visibility relations vis_{wk} and vis_{st} and extends it to vis_{wk}^n and vis_{st}^n such that they individually satisfy $VisBasic(\alpha_w)$ (Line 49) and $VisBasic(\alpha_s)$

(Line 51) respectively and jointly satisfy φ (Lines 53–56). We repeat this till we can extend these relations no longer, which implies that they have satisfied all the constraints (Lines 58–59).

The procedure **TestMultiCorrect** in Lines 42–57 takes as input a hybrid history $H = (\mathcal{O}, \text{so})$ whose multilevel correctness we want to check with respect to formulas α_w , α_s and φ .

We first enumerate the set of possible *reads-from* relations on the history (line 43). We then iterate through each of the *reads-from* relations rf to see whether it can be extended to construct a minimal visibility relation satisfying all the constraints and having no bad-patterns (Lines 44–55). For each rf , we construct minimal visibility relations $\text{vis}_{\text{wk}}^{\min}$ and $\text{vis}_{\text{st}}^{\min}$ extending rf_{wk} and rf_{st} respectively and satisfying the subformulas $\text{VisBasic}(\alpha_w)$ and $\text{VisBasic}(\alpha_s)$ respectively (Lines 45, 48).

If α_w (resp. α_s) contains the subformula $\text{total}(\text{vis})$, we enumerate the set of all the total orders extending $\text{vis}_{\text{wk}}^{\min}$ (resp. $\text{vis}_{\text{st}}^{\min}$) in the set $\text{visSet}_{\text{wk}}$ (resp. $\text{visSet}_{\text{st}}$) in Line 46 (resp. Line 49). If α_w (resp. α_s) does not contain the subformula $\text{total}(\text{vis})$, then, $\text{visSet}_{\text{wk}}$ (resp. $\text{visSet}_{\text{st}}$) will contain the only minimum visibility relation extending rf_{wk} (resp. rf_{st}), i.e., $\text{vis}_{\text{wk}}^{\min}$ (resp. $\text{vis}_{\text{st}}^{\min}$).

For each pair of visibility relations from $\text{visSet}_{\text{wk}}$ and $\text{visSet}_{\text{st}}$ we compute their stable extensions $\text{vis}_{\text{wk}}^{\text{stb}}$ and $\text{vis}_{\text{st}}^{\text{stb}}$ which individually satisfy α_w and α_s , respectively, and jointly satisfy φ (line 52). We then check if this computed extension has a bad pattern (Line 54). If no bad patterns are found, we return the $(\text{rf}, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}})$ as the witness.

If none of the rf can be extended to obtain the required visibility relation, we declare that the history is a bad history. We formally prove the correctness of **TestMultiCorrect** in the full version of this work [8].

Theorem 2 (Correctness of TestMultiCorrect procedure). *For a hybrid read-write history $H = (\mathcal{O}, \text{so})$ with weak and strong consistency criteria given by α_w and α_s , respectively, and multilevel constraints given by φ , the procedure **TestMultiCorrect** returns a witness $(\text{rf}, \text{vis}_{\text{wk}}, \text{vis}_{\text{st}})$ over H iff H is multi-level correct with respect to α_w , α_s and φ .*

4.3 Complexity

Suppose $H = (\mathcal{O}, \text{so})$ is history with $|\mathcal{O}| = N$.

We note that in the procedure **ComputeStableExt**, at the end of every iteration of the outer **while**-loop, the values of vis_{wk}^n and vis_{st}^n monotonically increase from the end of the previous iteration. Since they are binary relations over finite history $H = (\mathcal{O}, \text{so})$ their size is upper bounded by $O(N^2)$. The time taken to evaluate each term in $\text{RelTerms}(\alpha_\ell)$ is again polynomial in N . Hence, the time-complexity of **ComputeStableExt** is polynomial in N , say $f(N)$.

We can observe from the procedure **TestMultiCorrect** that the main part that adds to the complexity is iterating through all the *reads-from* relation, as well as the total orders if α_w or α_s contain the subformula $\text{total}(\text{vis})$. Suppose the number of read operations are k . Then the number of write operations is $N - k$,

and there are $O((N - k)^k)$ *reads-from* relations. Since $k = O(N)$, this can be bound by $O(2^{N \log N})$. Furthermore, for a given *rf*, if any of the levels $\ell \in \{\text{wk}, \text{st}\}$ require that the visibility relation be a total order, then we iterate over all the total-orders containing the minimal visibility relation extending *rf*. Iterating through this requires time bounded by $O(2^{N \log N})$. Thus the worst case time complexity of the procedure is $O(f(N) \cdot 2^{N \log N})$.

In general, the problem of testing the correctness of a hybrid history is in NP. We need to guess the *reads-from* relation, and then, extend it to obtain the minimal visibility relations satisfying the visibility constraints of the *wk* and the *st* consistency criteria. If the visibility relation is required to be a total order, we can guess the order. Extending this to derive fixed-point minimal visibility relations that satisfy all the visibility constraints via **ComputeStableExt** requires polynomial time. Subsequently checking for each of the bad-patterns requires polynomial time.

Note that we can reduce the testing of the correctness of a regular history (that contains only a single level of Read and Write operations) with respect to consistency criterion α to this procedure by defining the level of all the read operations to *st*. We set α_s to α , α_w to \top , and φ to $\psi_{back}^{write} \wedge \psi_{thru}^{read}$. For any *reads-from* relation *rf*, $\text{rf}_{\text{wk}} = \emptyset$. Thus $\text{vis}_{\text{wk}} = \emptyset$, trivially satisfying α_w as well as φ . Thus, the lower bound for testing the correctness of the hybrid history H is the complexity of testing the correctness of the H_{wk} and H_{st} with respect to their respective consistency criteria. It has been shown in [14] that testing the correctness of a read-write history with respect to sequential consistency is NP-COMPLETE. In [6], the authors use the same reduction to show that testing the correctness with respect to causal consistency is NP-COMPLETE. However, it can be shown that the reduction works for any consistency criterion stronger than FIFO consistency, and checking correctness with respect to such a consistency criterion is NP-COMPLETE. Thus, in general, though testing the multi-level correctness of a hybrid history is a hard problem, the hardness is not due to the multilevel constraints but due to the constraints of the individual consistency criteria and the read-write specification.

In [6], the authors identify the class of read-write data-stores called *data-independent* data-stores whose behaviour is not dependent on the exact values written to the keys. Thus, for such stores, if there is a bad history, there is an equivalent bad *differentiated history* where a particular value is written to a particular memory location at most once. Thus, we can restrict our testing to only the correctness of differentiated histories. The authors show that the problem of testing the correctness of differentiated-histories with respect to causal consistency is solvable in polynomial time.

Note that for differentiated histories, there is exactly one *reads-from* relation which associates every Read operation with at most one Write operation which has written that value to the memory location read by the Read operation. Thus, if neither of α_w or α_s contain the subformula $\text{total}(\text{vis})$, the procedure **TestMultiCorrect** terminates in polynomial time. Thus, our procedure generalizes the result from [6] to all the consistency criteria defined in terms of the

set of formulas involving only visibility, but not totality constraints. Our procedure checks the multi-level correctness of hybrid histories where the individual consistency levels do not require the visibility relation to be a total order, in polynomial time.

On the other hand, if one of α_w or α_s contains $\text{total}(\text{vis})$, then the worst case complexity remains $O(2^{N \log N})$. Once again, this does not come as a surprise, since the problem of testing the correctness of a differentiated history w.r.t. sequential consistency is not known to have a polynomial time solution.

5 Related Work

There is prior work that illustrates the need for multiple levels of consistency provided by the distributed data-stores to provide a trade off between consistency and availability/latency [2, 17, 18, 20]. The work by Kraska et al. [18] provides a transactional paradigm that allows applications to define the consistency level on data instead of transactions, and also allows the application to switch consistency guarantees at runtime. In the work by Guerraoui et al. [17], the authors provide a generic library that allows applications to request multiple responses to the same query, where the response that comes later in time is *more-correct* than the prior responses. Thus, later responses are supposed to have more knowledge of the state of the system compared to earlier responses. In our work, we have defined multilevel constraints, which can model the requirement of incremental consistency guarantees by requiring that subsequent strong responses see the effects observed by prior weak responses.

Burckhardt [9] provides a generic methodology for formalizing the specification of distributed data-stores in terms of histories, visibility and arbitration orders and provides an axiomatic characterization for consistency criteria. In our work, we have derived the specification for read-write stores based on this formalism. We have adapted this characterization to define consistency criteria as a conjunction of individual formulas. Our work extends [9] in terms of the definition of hybrid histories and provides a definition of multi-level correctness for read-write stores.

There is prior work on verifying the correctness of a behaviour with respect to individual consistency criteria. Examples include [7], which deals with verifying the correctness with respect to eventual consistency, [5], which investigates the feasibility of checking a concurrent implementation with respect to a consistency criterion that has a sequential specification, including sequential consistency, linearizability and conflict-serializability and [6], which focusses on correctness with respect to causal consistency. Our work provides a generic procedure for checking the correctness of read-write histories for all these individual consistency criteria. Further, [6] show that verification of correctness of a history with respect to causal consistency is NP-COMPLETE. However, for differentiated histories, the problem is solvable in polynomial time. In our work, we generalize the technique of computing the minimal visibility relation and checking for the absence of bad patterns for all the consistency criteria defined using our syntax. In [12], the authors model quiescent consistency using Mazurkiewicz Trace

Theory. They show that the testing problem (which they call the membership problem) for a history is **NP-COMPLETE**. We cannot model quiescent consistency in our framework since we cannot model a quiescent point. In [14], the authors present a detailed complexity analysis of the problem of testing the correctness of a history with respect to various consistency criteria. Our findings are consistent with the results from [14] with respect to hardness of testing consistency criteria that require the visibility relation to be a total order. In a recent work [13], the authors provide a technique for testing the correctness of a history of a data-store with respect to a weak consistency criterion. That work also characterizes correctness in terms of minimal visibility relation extending the session order (called program-order there) and the happened-before relation (called returns-before relation in [9]). Our work applies this concept to read-write stores, where we observe that correctness with respect to visibility constraints can be satisfied by constructing a minimal visibility relation while the correctness with respect to read-write specifications and arbitration constraints can be reduced to checking for absence of certain bad patterns. In particular, our characterization of the arbitration relation in terms of the conflict relation saves the step of searching through all possible arbitration relations which is used in [13].

[16] deals with verification of *red-blue* consistency where, in a history, a subset of operations are labelled *red* while the remaining are labelled *blue*. The *blue* operations are expected to satisfy a weaker consistency criterion, while the *red* operations are supposed to satisfy a stronger consistency criterion. The effects of the strong operations and weak operations are visible to each other. We can model this by setting $\varphi = \psi_{thru}^{write} \wedge \psi_{back}^{read}$.

Our work should also be contrasted with [3], which addresses the problem of checking the consistency of CRDTs against their specifications, and covers a wide range of CRDTs including replicated sets, flags, counters, registers, etc. The relevant data structure in our case is registers, where the results are comparable (checking w.r.t. the weaker consistency criterion is tractable). However, we also consider registers with multiple consistency criteria in this paper, which is not considered there.

Another related work is [4], which uses the reads-from relation (called the *write-read* relation there) to show that testing the correctness of an execution (containing transactions) with respect to various consistency criteria like Read Committed (RC), Read Atomic (RA), Causal Consistency (CC), Prefix Consistency, and Snapshot Isolation. The key difference in the current work is that we consider histories having multiple consistency levels simultaneously while [4] considers executions consisting of transactions, under a single consistency criterion.

References

1. Amazon DynamoDB Developer Guide (API Version 2012-08-10): DAX and DynamoDB consistency Models (2018). <https://docs.aws.amazon.com/amazondynamodb/latest/developerguide/DAX.consistency.html>. Accessed 26 Sep 2019

2. Bailis, P., Ghodsi, A., Hellerstein, J.M., Stoica, I.: Bolt-on causal consistency. In: Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data, SIGMOD 2013, pp. 761–772. ACM, New York (2013). <https://doi.org/10.1145/2463676.2465279>
3. Biswas, R., Emmi, M., Enea, C.: On the complexity of checking consistency for replicated data types. In: Dillig, I., Tasiran, S. (eds.) CAV 2019. LNCS, vol. 11562, pp. 324–343. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-25543-5_19
4. Biswas, R., Enea, C.: On the complexity of checking transactional consistency. Proc. ACM Program. Lang. **3**(OOPSLA), 165:1–165:28 (2019). <https://doi.org/10.1145/3360591>
5. Bouajjani, A., Emmi, M.: Analysis of recursively parallel programs. ACM Trans. Program. Lang. Syst. **35**(3), 10:1–10:49 (2013). <https://doi.org/10.1145/2518188>
6. Bouajjani, A., Enea, C., Guerraoui, R., Hamza, J.: On verifying causal consistency. In: Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, pp. 626–638. ACM, New York (2017). <https://doi.org/10.1145/3009837.3009888>
7. Bouajjani, A., Enea, C., Hamza, J.: Verifying eventual consistency of optimistic replication systems. In: The 41st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2014, San Diego, CA, USA, January 20–21, 2014, pp. 285–296 (2014). <https://doi.org/10.1145/2535838.2535877>
8. Bouajjani, A., Enea, C., Mukund, M., Shenoy, R.G., Suresh, S.P.: Formalizing and checking multilevel consistency. Tech. rep., Chennai Mathematical Institute (2019). <http://www.cmi.ac.in/~spsuresh/pdfs/vmcai2020-tr.pdf>
9. Burckhardt, S.: Principles of eventual consistency. Found. Trends Program. Lang. **1**(1–2), 1–150 (2014). <https://doi.org/10.1561/25000000011>
10. Burckhardt, S., Gotsman, A., Yang, H., Zawirski, M.: Replicated data types: specification, verification, optimality. In: The 41st Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2014, San Diego, CA, USA, January 20–21, 2014, pp. 271–284 (2014)
11. Damien: DynamoDB vs Cassandra (2017). <https://www.beyondthelines.net/databases/dynamodb-vs-cassandra/>. Accessed 16 Nov 2018
12. Dongol, B., Hierons, R.M.: Decidability and complexity for quiescent consistency. In: Proceedings of the 31st Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2016, pp. 116–125. ACM, New York (2016). <https://doi.org/10.1145/2933575.2933576>
13. Emmi, M., Enea, C.: Monitoring weak consistency. In: Chockler, H., Weissenbacher, G. (eds.) CAV 2018. LNCS, vol. 10981, pp. 487–506. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-96145-3_26
14. Furbach, F., Meyer, R., Schneider, K., Senftleben, M.: Memory model-aware testing - a unified complexity analysis. In: 14th International Conference on Application of Concurrency to System Design, ACSD 2014, Tunis La Marsa, Tunisia, June 23–27, 2014, pp. 92–101 (2014). <https://doi.org/10.1109/ACSD.2014.27>
15. Gilbert, S., Lynch, N.A.: Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. SIGACT News **33**(2), 51–59 (2002)
16. Gotsman, A., Yang, H., Ferreira, C., Najafzadeh, M., Shapiro, M.: ‘cause i’m strong enough: reasoning about consistency choices in distributed systems. In: Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2016, St. Petersburg, FL, USA, January 20–22, 2016, pp. 371–384 (2016). <https://doi.org/10.1145/2837614.2837625>

17. Guerraoui, R., Pavlovic, M., Seredinschi, D.A.: Incremental consistency guarantees for replicated objects. In: Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation, OSDI 2016, pp. 169–184. USENIX Association, Berkeley (2016). <http://dl.acm.org/citation.cfm?id=3026877.3026891>
18. Kraska, T., Hentschel, M., Alonso, G., Kossmann, D.: Consistency rationing in the cloud: pay only when it matters. *PVLDB* **2**(1), 253–264 (2009). <http://www.vldb.org/pvldb/2/vldb09-759.pdf>
19. Lamport, L.: How to make a multiprocessor computer that correctly executes multiprocess programs. *IEEE Trans. Comput.* **28**(9), 690–691 (1979)
20. Li, C., Porto, D., Clement, A., Gehrke, J., Preguiça, N., Rodrigues, R.: Making geo-replicated systems fast as possible, consistent when necessary. In: Proceedings of the 10th USENIX Conference on Operating Systems Design and Implementation, OSDI 2012, pp. 265–278. USENIX Association, Berkeley (2012). <http://dl.acm.org/citation.cfm?id=2387880.2387906>
21. Ozkan, B.K., Majumdar, R., Niksic, F., Befrouei, M.T., Weissenbacher, G.: Randomized testing of distributed systems with probabilistic guarantees. *PACMPL* **2**(OOPSLA), 160:1–160:28 (2018). <https://doi.org/10.1145/3276530>
22. Perrin, M., Mostefaoui, A., Jard, C.: Causal consistency: beyond memory. In: Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP 2016, pp. 26:1–26:12. ACM, New York (2016)
23. Terry, D.B., Theimer, M., Petersen, K., Demers, A.J., Spreitzer, M., Hauser, C.: Managing update conflicts in bayou, a weakly connected replicated storage system. In: Proceedings of the Fifteenth ACM Symposium on Operating System Principles, SOSP 1995, Copper Mountain Resort, Colorado, USA, December 3–6, 1995, pp. 172–183 (1995). <https://doi.org/10.1145/224056.224070>
24. Wolper, P.: Expressing interesting properties of programs in propositional temporal logic. In: Conference Record of the Thirteenth Annual ACM Symposium on Principles of Programming Languages, St. Petersburg Beach, Florida, USA, January 1986, pp. 184–193 (1986). <https://doi.org/10.1145/512644.512661>