

Testing Weakly Consistent Geo-Distributed Data-Storages and Applications

Crash in the Cloud

by

Ranadeep Biswas

Submitted at UNIVERSITÉ DE PARIS, IRIF, CNRS, F-75013 PARIS, FRANCE
on Jan, 2021 in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

ABSTRACT

Recent distributed systems have introduced various abstract concepts like *conflict-free replicated data types* (CRDTs), transaction, weak consistency and isolation level.

CRDTs like counters, registers, flags, and sets provide high availability and partition tolerance. They utilize nontrivial mechanisms to resolve the effects of concurrent updates to replicated data. Naturally these objects weaken their consistency guarantees to achieve availability and partition-tolerance, and various notions of *weak consistency* capture those guarantees.

Transactions simplify concurrent programming by enabling computations on shared data that are isolated from other concurrent computations and resilient to failures. Modern databases provide different consistency models or isolation level for transactions corresponding to different tradeoffs between consistency and availability.

This thesis studies the tractability of CRDT-consistency checking and transactional database checking for an isolation level. In both cases, to capture guarantees precisely, and facilitate symbolic reasoning, we propose novel logical characterizations. By developing novel reductions from propositional satisfiability problems, and novel consistency-checking algorithms, we discover both positive and negative results. We also demonstrate that tractability for hard cases can be redeemed if some parameter is bounded, usually the number of session or replicas in the network.

Lastly, this thesis presents MonkeyDB, a mock storage system for testing storage-backed applications. MonkeyDB supports a Key-Value interface as well as SQL queries under multiple isolation levels. It uses a logical specification of the isolation level to compute, on a read operation, the set of all possible return values. MonkeyDB then returns a value randomly from this set. We show that MonkeyDB provides good coverage of weak behaviors, which is complete in the limit. We test a variety of applications for assertions that fail only under weak isolation. MonkeyDB is able to break each of those assertions in a small number of attempts.

ACKNOWLEDGMENTS

CONTENTS

1	INTRODUCTION	1
1.1	Conflict Free Replicated Data Types	2
1.2	Transactional Systems	2
1.3	Applications implemented on Distributed Datastores	4
2	CONSISTENCY OF CONFLICT-FREE REPLICATED DATA TYPES	7
2.1	Introduction	7
2.2	A Logical Characterization of Replicated Data Types	9
2.2.1	Replicated Sets and Flags	10
2.2.2	Replicated Registers	11
2.2.3	Replicated Counters	12
2.2.4	Replicated Growable Array	12
2.3	Intractability for Registers, Sets, Flags, and Counters	13
2.4	Polynomial-Time Algorithms for Registers and Arrays	17
2.5	Polynomial-Time Algorithms for Replicated Counters	19
2.6	Polynomial-Time Algorithms for Sets and Flags	20
2.7	Related Work	21
2.8	Conclusion	22
3	TRANSACTIONAL CONSISTENCY	23
3.1	Introduction	23
3.2	Overview	25
3.3	Consistency Criteria	27
3.3.1	Histories	27
3.3.2	Axiomatic Framework	28
3.4	Checking Consistency Criteria	34
3.5	Checking Consistency of Bounded-Width Histories	38
3.5.1	Checking Serializability	38
3.5.2	Reducing Prefix Consistency to Serializability	41
3.5.3	Reducing Snapshot Isolation to Serializability	45
3.6	Communication graphs	46
3.7	Experimental Evaluation	49
3.8	Related Work	52
4	APPLICATIONS IMPLEMENTED ON DISTRIBUTED DATASTORES	55
4.1	Introduction	55
4.2	Programming Language	58

Contents

4.3	Operational Semantics for \mathcal{P}_{KV}	59
4.3.1	Definition of the Operational Semantics	60
4.3.2	Correctness of the Operational Semantics	62
4.4	Compiling SQL to Key-Value API	64
4.5	Implementation	66
4.6	Evaluation: Microbenchmarks	66
4.6.1	Applications	67
4.6.2	Assertion Checking	68
4.6.3	Coverage	68
4.7	Evaluation: OLTP Workloads	69
4.8	Testing SQL Database	72
4.9	Conclusion	73

1 INTRODUCTION

As *internet* grows to be cheaper and faster, distributed systems and applications are also becoming popular over past decades. Today they are the backbone of almost any online application like banking, shopping, ticketing, social networking etc.. As the popularity of these softwares increases, it is very important that these systems are reliable and secure.

Distributed systems are usually deployed over multiple locations connected through *internet*. There are multiple good reasons why an online application would need that. The users of the application may require data availability and fast response. Or, if some location goes offline, the other locations can still serve the users. Usually each distributed *nodes* maintains a replica of the global data. A user connects to the node *nearest* to them and that node serves their requests. This way the distributed system reduces responses time and distributes the workload to multiple nodes.

Distributed systems are the solution to a fast, concurrent, scalable online applications. But it actually offers a trade off. As we allow concurrent operations at multiple nodes, they still need to synchronize among them and maintain a meaningful or *consistent* global data. A global lock defeats the whole purpose of the distributed system, because taking a lock over a multiple distributed nodes means more communication over internet and slow response time.

Over the past few years, the community came up with many implementations *weakly-consistent* distributed system. They throw away strong consistent guarantees of the global data in favour of a scalable system. This is almost always alright. Because in most cases, the distributed applications do not really require strong guarantees over the data.

But these implementations are often not correctly reasoned. Many cases the guaranteed weak consistencies are explained in words and not formalized, which makes it difficult to reason about them. Also, often these systems promise a weak consistency on paper, but implements something else in practice. So it is very crucial, these systems are tested thoroughly against the consistency needs of an application.

This dissertation is an effort to give solutions to these problems. We first study *Conflict Free Replicated Data Type* or *CRDT*. CRDTs are popular data types like sets, flags, registers, but replicated and maintained at multiple nodes. Then we move on to *Transactional Systems* which allows to perform a set of operations isolated from other sets of operations. In both cases we reason about their consistency guarantees and provide algorithms to test their histories against the promised weak-consistency. Lastly, we study the correctness of applications implemented over distributed systems guaranteeing weak consistency level. We implement an artificial data storage, MonkeyDB, based on our previous formalization of consistency levels to simulate the possible abstraction of weak scenarios such as network partition, packet drops etc.. These applications can use MonkeyDB as a drop-in replacement to their backend databases to inject faults without any external control over network.

1.1 CONFLICT FREE REPLICATED DATA TYPES

Replication is used in any distributed system. But to maintain a reliable replicated data, one has to maintain some level of consistency through clever synchronization techniques. The standard *strong consistency* approach is to find a linearization of the set of all global operations. But as CAP theorem shows, achieving strong consistency is not possible while maintaining availability and network partition.

When network delays are big and network partition is an issue, *eventual consistency* promises better performance and availability. An update gets executed at some replica but it may be sent to other replicas later as a delayed synchronization. Although updates eventually reaches to all replicas, these updates may be seen concurrently, even in different orders, by different replicas. If these concurrent updates are *conflicting*, they need to be resolved. This is not trivial because the resolved outcome need to be the *same* in all replicas such that the replicas do not diverge.

Conflict-free replicated data types (CRDTs) [DBLP:conf/sss/ShapiroPBZ11] efficiently resolve the effects of concurrent updates to replicated data. In figure ??, we demonstrate a CRDT, specifically an *OR-Set*. OR-Set implements a distributed Set with natural `add(_)`, `remove(_)`, `contain(_)` operations. `add(_)` and `remove(_)` are the only update operations. Two updates are in conflict if they are trying to insert or remove a same element. As its CRDT based conflict-resolve technique, an `add(_)` operation will always win among multiple conflicting updates. In figure ??, the first replica receives the update from third replica. Even though that update was in conflict, there was no `add(x)`. So `contain(x)` returns false. But in the third replica receives all the concurrent updates. Since there is one concurrent `add(x)`, they all get resolved to `add(x)` and the later `contain(x)` returns true. Similarly, *And-Set* is a CRDT in which `remove(_)` always wins among multiple conflicting updates.

Recent distributed systems have introduced variations of familiar abstract data types (ADTs) like counters, registers, flags, and sets, that provide high availability and partition tolerance. Naturally they weaken consistency guarantees to achieve availability and partition-tolerance, and various notions of *weak consistency* capture such guarantees [DBLP:conf/pdis/TerryDPSTW94, DBLP:conf/sosp/TerryTPDSH95, DBLP:conf/popl/MansonPA05, DBLP:journals/ftpl/Burckhardt14, DBLP:conf/popl/BurckhardtGYZ14].

1.1.1 STATE OF THE ART

Many have considered consistency models applicable to CRDTs, including causal consistency [DBLP:journals/cacm/Lam78], sequential consistency [DBLP:journals/tc/Lamport79], linearizability [DBLP:journals/toplas/HerlihyW90], session consistency [DBLP:conf/pdis/TerryDPSTW94], eventual consistency [DBLP:conf/sosp/TerryTPDSH95], and happens-before consistency [DBLP:conf/popl/MansonPA05]. Burckhardt et al. [DBLP:journals/ftpl/Burckhardt14, DBLP:conf/popl/BurckhardtGYZ14] propose a unifying framework to formalize these models. Many have also studied the complexity of verifying data-type agnostic notions of consistency, including serializability, sequential consistency and linearizability [DBLP:journals/jacm/Papadimitriou79b, DBLP:journals/siamcomp/GibbonsK97, DBLP:journals/iandc/AlurMP00, DBLP:conf/spaa/BinghamCH03, DBLP:conf/cav/FarzanM08, DBLP:conf/esop/BouajjaniEEH13, DBLP:conf/netys/Hamza15], as well as causal consistency [DBLP:conf/popl/BouajjaniEGH17]. Bouajjani et al. [DBLP:journals/iandc/BouajjaniEEH18, DBLP:journals/pacmpl/EmmiE18] consider the complexity for individual linearizable collection types. Others have developed effective consistency checking algorithms for sequential consistency.

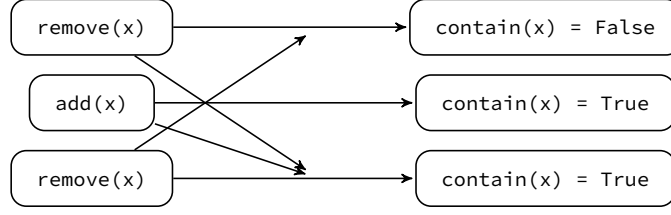


Figure 1.1: OR-Set

tency [DBLP:conf/cav/HenzingerQR99a, DBLP:journals/tpds/Qadeer03, DBLP:conf/cav/BinghamCHQZ04, DBLP:conf/pldi/BurckhardtAM07], serializability [DBLP:conf/fmcad/0002OPTZ07, DBLP:conf/cav/FarzanM08, DBLP:conf/pldi/GuerraouiHJS08, DBLP:conf/pldi/EmmiMM10], linearizability [DBLP:journals/jpdc/WingG93, DBLP:conf/pldi/BurckhardtDMT10, DBLP:conf/pldi/EmmiEH15, DBLP:journals/concurrency/Lowe17], and even weaker notions like eventual consistency [DBLP:conf/popl/BouajjaniEH14] and sequential happens-before consistency [DBLP:conf/cav/EmmiE18, DBLP:journals/pacmpl/EmmiE19].

1.1.2 CONTRIBUTION

In this work we study the tractability of CRDT consistency checking. In particular, we consider *runtime verification*: deciding whether a given execution of a CRDT is consistent with its ADT specification. This problem is particularly relevant as distributed-system testing tools like Jepsen [MISC:Jepsen] are appearing; without efficient, general consistency-checking algorithms, such tools could be limited to specialized classes of errors like node crashes.

Our study proceeds in three parts. First, to precisely characterize the consistency of various CRDTs, and facilitate symbolic reasoning, we develop novel logical characterizations to capture their guarantees. Second, we demonstrate the intractability of several replicated data types by reduction from propositional satisfiability (SAT) problems. Third, we develop tractable consistency-checking algorithms for individual data types and special cases: replicated growing arrays; multi-value and last-writer-wins registers, when each value is written only once; counters, when replicas are bounded; and sets and flags, when their sizes are also bounded.

1.2 TRANSACTIONAL SYSTEMS

Transactions simplify concurrent programming by enabling multiple computations on shared data that are isolated from other concurrent computations and resilient to failures. Modern databases provide transactions in various forms corresponding to different tradeoffs between consistency and availability. The strongest level of consistency is achieved with *serializable* transactions [DBLP:journals/jacm/Papadimitri78] whose outcome in concurrent executions is the same as if the transactions were executed atomically in some order. Unfortunately, serializability carries a significant penalty on the availability of the system assuming, for instance, that the database is accessed over a network that can suffer from partitions or failures. For this reason, modern databases often provide weaker guarantees about transactions, formalized by weak consistency models, e.g., causal consistency [DBLP:journals/cacm/Lamport78] and snapshot isolation [DBLP:conf/sigmod/BerensonBGM0095].

Implementations of large-scale databases providing transactions are difficult to build and test. For instance, distributed (replicated) databases must account for partial failures, where some components or the network can fail and produce incomplete results. Ensuring fault-tolerance relies on intricate protocols that are difficult to design and reason about. The black-box testing framework Jepsen [jepsen] found a remarkably large number of subtle problems in many production distributed databases.

1.2.1 STATE OF THE ART

Testing a transactional database raises two issues. First, deriving a suitable set of testing scenarios, e.g., faults to inject into the system and the set of transactions to be executed, and second, deriving efficient algorithms for checking whether a given execution satisfies the considered consistency model. The Jepsen framework aims to address the first issue by using randomization, e.g., introducing faults at random and choosing the operations in a transaction randomly. The effectiveness of this approach has been proved formally in recent work [DBLP:journals/pacmpl/OzkanMNBW18]. The second issue is, however, largely unexplored. Jepsen checks consistency in a rather ad-hoc way, focusing on specific classes of violations to a given consistency model, e.g., dirty reads (reading values from aborted transactions). This problem is challenging because the consistency specifications are non-trivial and they cannot be checked using, for instance, standard local assertions added to the client's code.

Besides serializability, the complexity of checking correctness of an execution w.r.t. some consistency model is unknown. Checking serializability has been shown to be NP-complete [DBLP:journals/jacm/Papadimitriou79] and checking causal consistency in a *non-transactional* context is known to be polynomial time [DBLP:conf/popl/Bouajjaj08].

1.2.2 CONTRIBUTION

In this work, we try to fill this gap by investigating the complexity of this problem w.r.t. several consistency models and, in the case of NP-completeness, devising algorithms that are polynomial time assuming fixed bounds for certain parameters of the input executions, e.g., the number of sessions. We consider several consistency models that are the most prevalent in practice. *Read Committed* (RC) [DBLP:conf/sigmod/BerensonBGM0095], *Read Atomic* (RA) [DBLP:conf/concur/Cerone0615], *Causal Consistency* (CC) [DBLP:journals/cacm/Lamport78], *Prefix Consistency* (PC) [DBLP:conf/ecoop/Burckhardt10], *Snapshot Isolation* (SI) [DBLP:conf/sigmod/BerensonBGM0095] and *Serializability* (SER).

We establish that checking whether an execution satisfies RC, RA, or CC is polynomial time, while the same problem is NP-complete for PC and SI. Moreover, in the case of the NP-complete consistency models (PC, SI, SER), we show that their verification problem becomes polynomial time provided that, roughly speaking, the number of sessions in the input executions is considered to be fixed (i.e., not counted for in the input size). In more detail, we establish that checking SER reduces to a search problem in a space that has polynomial size when the number of sessions is fixed. (This algorithm applies to arbitrary executions, but its complexity would be exponential in the number of sessions in general.) Then, we show that checking PC or SI can be reduced in polynomial time to checking SER using a transformation of executions that, roughly speaking, splits each transaction in two parts: one part containing all the reads, and one part containing all the writes (SI further requires adding some additional variables in order to deal with transactions writing on a common variable).

We provide an experimental evaluation of our algorithms on executions of CockroachDB [**cockroach**], which claims to implement serializability [**cockroach-claim**] acknowledging however the possibility of anomalies, Galera [**galera**], whose documentation contains contradicting claims about whether it implements snapshot isolation [**galera-claim**, **galera-notclaim**], and AntidoteDB [**antidote**], which claims to implement causal consistency [**antidote-claim**]. Our implementation reports violations of these criteria in all cases. The consistency violations we found for AntidoteDB are novel and have been confirmed by its developers. We show that our algorithms are efficient and scalable. In particular, we show that, although the asymptotic complexity of our algorithms is exponential in general (w.r.t. the number of sessions), the worst-case behavior is not exercised in practice.

1.3 APPLICATIONS IMPLEMENTED ON DISTRIBUTED DATASTORES

Modern applications require not just data reliability, but also high-throughput concurrent accesses. Applications concerning supply chains, banking, etc. use traditional relational databases for storing and processing data, whereas applications such as social networking software and e-commerce platforms use cloud-based storage systems (such as Azure CosmosDb [**cosmosdb**], Amazon DynamoDb [**DBLP:conf/sosp/DeCandiaHJKLPSVV07**], Facebook TAO [**DBLP:conf/usenix/BronsonACDDDFGKLM** etc.]).

Providing high-throughput processing, unfortunately, comes at an unavoidable cost of weakening the guarantees offered to users. Concurrently-connected clients may end up observing different views of the same data. These “anomalies” can be prevented by using a strong *isolation level* such as *serializability*, which essentially offers a single view of the data. However, serializability requires expensive synchronization and incurs a high performance cost. As a consequence, most storage systems use weaker isolation levels, such as *Causal Consistency* [**DBLP:journals/cacm/Lamport78**, **DBLP:conf/sosp/LloydFKA11**, **antidote-white-paper**], *Snapshot Isolation* [**DBLP:conf/sigmod/BerensonBGM0095**], *Read Committed* [**DBLP:conf/sigmod/BerensonBGM0095**], etc. for better performance. In a recent survey of database administrators [**DBLP:conf/sigmod/Pavlo17**], 86% of the participants responded that most or all of the transactions in their databases execute at read committed isolation level.

A weaker isolation level allows for more possible behaviors than stronger isolation levels. It is up to the developers then to ensure that their application can tolerate this larger set of behaviors. Unfortunately, weak isolation levels are hard to understand or reason about [**DBLP:conf/popl/BrutschyD0V17**, **adya-thesis**] and resulting application bugs can cause loss of business [**DBLP:conf/sigmod/WarszawskiB17**].

1.3.1 STATE OF THE ART

This work addresses the problem of *testing* code for correctness against weak behaviors: a developer should be able to write a test that runs their application and then asserts for correct behavior. The main difficulty today is getting coverage of weak behaviors during the test. If one runs the test against the actual production storage system, it is very likely to only result in serializable behaviors because of their optimized implementation. For instance, only 0.0004% of all reads performed on Facebook’s TAO storage system were not serializable [**DBLP:conf/sosp/LuVAHSTKL15**]. Emulators, offered by cloud providers for local development, on the other hand, do not support weaker isolation levels at all [**cosmosdb-local**]. Another option, possible when the storage system is avail-

able open-source, is to set it up with a tool like Jepsen [jepsen] to inject noise (bring down replicas or delay packets on the network). This approach is unable to provide good coverage at the level of client operations [DBLP:journals/pacmpl/RahmaniNDJ19] (§4.7). Another line of work has focussed on finding anomalies by identifying non-serializable behavior (§4.8). Anomalies, however, do not always correspond to bugs [DBLP:conf/pldi/BrutschyD0V18, DBLP:journals/pvldb/GanRRB020]; they may either not be important (e.g., gather statistics) or may already be handled in the application (e.g., checking and deleting duplicate items).

1.3.2 CONTRIBUTION

We present MonkeyDB, a mock in-memory storage system meant for testing correctness of storage-backed applications. MonkeyDB supports common APIs for accessing data (key-value updates, as well as SQL queries), making it an easy substitute for an actual storage system. MonkeyDB can be configured with one of several isolation levels.

MonkeyDB does not rely on stress generation, fault injection, or data replication. Rather, it works directly with a formalization of the given isolation level in order to compute allowed return values. The theory behind MonkeyDB builds on the axiomatic definitions of isolation levels introduced in chapter ?? [DBLP:journals/pacmpl/BiswasE19]. These definitions use logical constraints (called *axioms*) to characterize the set of executions of a key-value store that conform to a particular isolation level (we discuss SQL queries later).

We consider a set of micro-benchmarks Twitter [twissandra], Shopping Cart [DBLP:conf/pldi/Sivaramakrishnan15], Courseware [DBLP:conf/esop/NairP020], Treiber Stack [DBLP:conf/cav/NagarMJ20] inspired from real-world applications (§4.6.1) and evaluate the number of test iterations required to fail an invalid assertion (§4.6.2). We also measure the *coverage* of weak behaviors provided by MonkeyDB (§4.6.3). Each of these applications were implemented based on their specifications described in prior work; they all use MonkeyDB as a library, via its KV interface.

We also extend our semantics to cover SQL queries as well, by compiling SQL queries down to transactions with multiple key-value reads/writes. Naturally, we considered OLTPBench [DBLP:journals/pvldb/Difall] a popular benchmark suite of representative OLTP workloads for relational databases. We picked a subset of OLTPBench, including TPC-C, for which we had reasonable assertions. We showed using our technique, MonkeyDB was able to violate these assertion under weak isolation levels.

2

CONSISTENCY OF CONFLICT-FREE REPLICATED DATA TYPES

2.1 INTRODUCTION

Recent distributed systems have introduced variations of familiar abstract data types (ADTs) like counters, registers, flags, and sets, that provide high availability and partition tolerance. These *conflict-free replicated data types* (CRDTs) [DBLP:conf/sss/ShapiroPBZ11] efficiently resolve the effects of concurrent updates to replicated data. Naturally they weaken consistency guarantees to achieve availability and partition-tolerance, and various notions of *weak consistency* capture such guarantees [DBLP:conf/pdis/TerryDPSTW94, DBLP:conf/sosp/TerryTPDSH95, DBLP:conf/popl/MansonPA05, DBLP:journals/ftpl/Burckhardt14, DBLP:conf/popl/BurckhardtGYZ14].

In this work we study the tractability of CRDT consistency checking; Figure 2.1 summarizes our results. In particular, we consider *runtime verification*: deciding whether a given execution of a CRDT is consistent with its ADT specification. This problem is particularly relevant as distributed-system testing tools like Jepsen [MISC:Jepsen] are appearing; without efficient, general consistency-checking algorithms, such tools could be limited to specialized classes of errors like node crashes.

Our setting captures executions across a set of replicas as per-replica sequences of operations called *histories*. Roughly speaking, a history is *consistent* so long as each operation’s return value can be justified according to the operations that its replica has observed so far. In the setting of CRDTs, the determination of a replica’s observations is essentially an implementation choice: replicas are only obliged to observe their own operations, and the predecessors of those it has already observed. This relatively-weak constraint on replicas’ observations makes the CRDT consistency checking problem unique.

Our study proceeds in three parts. First, to precisely characterize the consistency of various CRDTs, and facilitate symbolic reasoning, we develop novel logical characterizations to capture their guarantees. Our logical models are built on a notion of *abstract execution*, which relates the operations of a given history with three separate relations: a *read-from* relation, governing the observations from which a given operation constitutes its own return value; a *happens-before* relation, capturing the causal relationships among operations; and a *linearization* relation, capturing any necessary arbitration among non-commutative effects which are executed concurrently, e.g., following a *last-writer-wins* policy. Accordingly, we capture data type specifications with logical axioms interpreted over the read-from, happens-before, and linearization relations of abstract executions, reducing the consistency problem to: does there exist an abstract execution over the given history which satisfies the axioms of the given data type?

Data Types	Complexity
Add-Wins Set, Remove-Wins Set	NP-complete
Enable-Wins Flag, Disable-Wins Flag	NP-complete
Sets & Flags — with bounded domains	P _{TIME}
Last-Writer-Wins Register (LWW)	NP-complete
Multi-Value Register (MVR)	NP-complete
Registers – with unique values	P _{TIME}
Replicated Counters	NP-complete
Counters – with bounded replicas	P _{TIME}
Replicated Growable Array (RGA)	P _{TIME}

Figure 2.1: The complexity of consistency checking for various replicated data types. We demonstrate intractability and tractability results in Sections 2.3 and 2.4, respectively.

Second, we demonstrate the intractability of several replicated data types by reduction from propositional satisfiability (SAT) problems. In particular, we consider the 1-in-3 SAT problem [DBLP:books/fm/Garey] which asks for a truth assignment to the variables of a given set of clauses such that exactly one literal per clause is assigned true. Our reductions essentially simulate the existential choice of a truth assignment with the existential choice of the read-from and happens-before relations of an abstract execution. For a given 1-in-3 SAT instance, we construct a history of replicas obeying carefully-tailored synchronization protocols, which is consistent exactly when the corresponding SAT instance is positive.

Third, we develop tractable consistency-checking algorithms for individual data types and special cases: replicated growing arrays; multi-value and last-writer-wins registers, when each value is written only once; counters, when replicas are bounded; and sets and flags, when their sizes are also bounded. While the algorithms for each case are tailored to the algebraic properties of the data types they handle, they essentially all function by constructing abstract executions incrementally, processing replicas' operations in prefix order.

The remainder of this article is organized around our three key contributions:

1. We develop novel logical characterizations of consistency for the replicated register, flag, set, counter, and array data types (§2.2);
2. We develop novel reductions from propositional satisfiability problems to consistency checking to demonstrate intractability for replicated flags, sets, counters, and registers (§2.3); and
3. We develop tractable consistency-checking algorithms for replicated growable arrays, registers, when written values are unique, counters, when replicas are bounded, and sets and flags, when their sizes are also bounded (§2.4–2.6).

Section 4.8 overviews related work, and Section 2.8 concludes.

2.2 A LOGICAL CHARACTERIZATION OF REPLICATED DATA TYPES

In this section we describe an axiomatic framework for defining the semantics of replicated data types. We consider a set of method names \mathbb{M} , and that each method $m \in \mathbb{M}$ has a number of arguments and a return value sampled from a data domain \mathbb{D} . We will use operation labels of the form $m(a) \xrightarrow{i} b$ to represent the call of a method $m \in \mathbb{M}$, with argument $a \in \mathbb{D}$, and resulting in the value $b \in \mathbb{D}$. Since there might be multiple calls to the same method with the same arguments and result, labels are tagged with a unique identifier i . We will ignore identifiers when unambiguous.

The interaction between a data type implementation and a client is represented by a *history* $h = \langle \text{Op}, \text{ro} \rangle$ which consists of a set of operation labels Op and a partial *replica order* ro ordering operations issued by the client on the same replica. Usually, ro is a union of sequences, each sequence representing the operations issued on the same replica, and the *width* of ro , i.e., the maximum number of mutually-unordered operations, gives the number of replicas in a given history.

To characterize the set of histories $h = \langle \text{Op}, \text{ro} \rangle$ admitted by a certain replicated data type, we use *abstract executions* $e = \langle \text{rf}, \text{hb}, \text{lin} \rangle$, which include:

- a *read-from* binary relation rf over operations in Op , which identifies the set of updates needed to “explain” a certain return value, e.g., a write operation explaining the return value of a read,
- a strict partial *happens-before* order hb , which includes ro and rf , representing the causality constraints in an execution, and
- a strict total *linearization* order lin , which includes hb , used to model conflict resolution policies based on timestamps.

In this work, we consider replicated data types which satisfy *causal consistency* [DBLP:journals/cacm/Lamport78], i.e., updates which are related by cause and effect relations are observed by all replicas in the same order. This follows from the fact that the happens-before order is constrained to be a partial order, and thus transitive (other forms of weak consistency don’t pose this constraint). Some of the replicated data types we consider in this paper do *not* consider resolution policies based on timestamps and in those cases, the linearization order can be ignored.

A *replicated data type* is defined by a set of first-order axioms Φ characterizing the relations in an abstract execution. A history h is *admitted* by a data type when there exists an abstract execution e such that $\langle h, e \rangle \models \Phi$. The satisfaction relation \models is defined as usual in first order logic. The *admissibility problem* is the problem of checking whether a history h is admitted by a given data type.

In the following, we define the replicated data types with respect to which we study the complexity of the admissibility problem. The axioms used to define them are listed in Figure 2.2 and Figure 2.3. These axioms use the function symbols *meth-od*, *arg-ument*, and *ret-urn* interpreted over operation labels, whose semantics is self-explanatory.

2 Consistency of Conflict-Free Replicated Data Types

<u>READFROM(R)</u>	<u>RETVALSET(X, v, Y)</u>
$\forall o_1, o_2. \text{rf}(o_1, o_2) \Rightarrow R(o_1, o_2)$	$\forall o_1. \text{meth}(o_1) = X \wedge \text{ret}(o_1) = v$
	$\Leftrightarrow \exists o_2. \text{rf}(o_2, o_1) \wedge \text{meth}(o_2) = Y$
	$\wedge \text{arg}(o_1) = \text{arg}(o_2)$
<u>READFROMMAXIMAL(R)</u>	<u>RETVALCOUNTER</u>
$\forall o_1, o_2, o_3. \text{rf}(o_1, o_2) \wedge R(o_3, o_2) \Rightarrow$	$\forall o_1. \text{meth}(o_1) = \text{read}$
$\neg \text{hb}(o_1, o_3) \vee \neg \text{hb}(o_3, o_2)$	$\Rightarrow \text{ret}(o_1) = \{o_2 : \text{meth}(o_2) = \text{inc} \wedge \text{rf}(o_2, o_1)\} $
	$- \{o_2 : \text{meth}(o_2) = \text{dec} \wedge \text{rf}(o_2, o_1)\} $
<u>READALLMAXIMALS(R)</u>	<u>LINLWW</u>
$\forall o_1, o_2. \text{hb}(o_1, o_2) \wedge R(o_1, o_2)$	$\forall o_1, o_2, o_3. \text{rf}(o_1, o_2) \wedge \text{meth}(o_3) = \text{write}$
$\Rightarrow \exists o_3. \text{hb}^*(o_1, o_3) \wedge \text{rf}(o_3, o_2)$	$\wedge \text{arg}_1(o_3) = \text{arg}(o_2) \wedge \text{hb}(o_3, o_2) \Rightarrow \text{lin}(o_3, o_1)$
<u>CLOSEDRF(R)</u>	
$\forall o_1, o_2, o_3. R(o_1, o_2) \wedge \text{hb}(o_1, o_3)$	
$\wedge \text{rf}(o_3, o_2) \Rightarrow \text{rf}(o_1, o_2)$	
<u>RETVALREG</u>	
$\forall o_1, v. \text{meth}(o_1) = \text{read} \wedge v \in \text{ret}(o_1) \Rightarrow \exists ! o_2. \text{rf}(o_2, o_1) \wedge \text{meth}(o_2) = \text{write} \wedge \text{arg}_2(o_2) = v$	

Figure 2.2: The axiomatic semantics of replicated data types. Quantified variables are implicitly distinct, and $\exists ! o$ denotes the existence of a unique operation o .

2.2.1 REPLICATED SETS AND FLAGS

The Add-Wins Set and Remove-Wins Set [DBLP:journals/eatcs/ShapiroPBZ11] are two implementations of a replicated set with operations $\text{add}(x)$, $\text{remove}(x)$, and $\text{contains}(x)$ for adding, removing, and checking membership of an element x . Although the meaning of these methods is self-evident from their names, the result of conflicting concurrent operations is not evident. When concurrent $\text{add}(x)$ and $\text{remove}(x)$ operations are delivered to a certain replica, the Add-Wins Set chooses to keep the element x in the set, so every subsequent invocation of $\text{contains}(x)$ on this replica returns *true*, while the Remove-Wins Set makes the dual choice of removing x from the set.

The formal definition of their semantics uses abstract executions where the read-from relation associates sets of $\text{add}(x)$ and $\text{remove}(x)$ operations to $\text{contains}(x)$ operations. Therefore, the predicate $\text{ReadOk}(o_1, o_2)$ is defined by

$$\text{meth}(o_1) \in \{\text{add}, \text{remove}\} \wedge \text{meth}(o_2) = \text{contains} \wedge \text{arg}(o_1) = \text{arg}(o_2)$$

and the Add-Wins Set is defined by the following set of axioms:

$$\begin{aligned} & \text{READFROM}(\text{ReadOk}) \wedge \text{READFROMMAXIMAL}(\text{ReadOk}) \wedge \\ & \text{READALLMAXIMALS}(\text{ReadOk}) \wedge \text{RETVALSET}(\text{contains}, \text{true}, \text{add}) \end{aligned}$$

READFROMMAXIMAL says that every operation read by a $\text{contains}(x)$ is maximal among its hb-predecessors that add or remove x while READALLMAXIMALS says that all such maximal hb-

predecessors are read. The RETVALSET instantiation ensures that a $\text{contains}(x)$ returns *true* iff it reads-from at least one $\text{add}(x)$.

The definition of the Remove-Wins Set is similar, except for the parameters of RETVALSET, which become $\text{RETVALSET}(\text{contains}, \text{false}, \text{remove})$, i.e., a $\text{contains}(x)$ returns *false* iff it reads-from at least one $\text{remove}(x)$.

The Enable-Wins Flag and Disable-Wins Flag are implementations of a set of flags with operations: $\text{enable}(x)$, $\text{disable}(x)$, and $\text{read}(x)$, where $\text{enable}(x)$ turns the flag x to true, $\text{disable}(x)$ turns x to false, while $\text{read}(x)$ returns the state of the flag x . Their semantics is similar to the Add-Wins Set and Remove-Wins Set, respectively, where $\text{enable}(x)$, $\text{disable}(x)$, and $\text{read}(x)$ play the role of $\text{add}(x)$, $\text{remove}(x)$, and $\text{contains}(x)$, respectively. Their axioms are defined as above.

2.2.2 REPLICATED REGISTERS

We consider two variations of replicated registers called Multi-Value Register (MVR) and Last-Writer-Wins Register (LWW) [DBLP:journals/eatcs/ShapiroPBZ11] which maintain a set of registers and provide $\text{write}(x, v)$ operations for writing a value v on a register x and $\text{read}(x)$ operations for reading the content of a register x (the domain of values is kept unspecified since it is irrelevant). While a $\text{read}(x)$ operation of MVR returns *all* the values written by concurrent writes which are maximal among its happens-before predecessors, therefore, leaving the responsibility for solving conflicts between concurrent writes to the client, a $\text{read}(x)$ operation of LWW returns a single value chosen using a conflict-resolution policy based on timestamps. Each written value is associated to a timestamp, and a read operation returns the most recent value w.r.t. the timestamps. This order between timestamps is modeled using the linearization order of an abstract execution.

Therefore, the predicate $\text{ReadOk}(o_1, o_2)$ is defined by

$$\text{meth}(o_1) = \text{write} \wedge \text{meth}(o_2) = \text{read} \wedge \text{arg}_1(o_1) = \text{arg}_2(o_2) \wedge \text{arg}_2(o_1) \in \text{ret}(o_2)$$

(we use $\text{arg}_1(o_1)$ to denote the first argument of a write operation, i.e., the register name, and $\text{arg}_2(o_1)$ to denote its second argument, i.e., the written value) and the MVR is defined by the following set of axioms:

$$\begin{aligned} & \text{READFROM}(\text{ReadOk}) \wedge \text{READFROMMAXIMAL}(\text{ReadOk}) \wedge \\ & \text{READALLMAXIMALS}(\text{ReadOk}) \wedge \text{RETVALREG} \end{aligned}$$

where RETVALREG ensures that a $\text{read}(x)$ operation reads from a $\text{write}(x, v)$ operation, for each value v in the set of returned values¹.

LWW is obtained from the definition of MVR by replacing READALLMAXIMALS with the axiom LINLWW which ensures that every $\text{write}(x, _)$ operation which happens-before a $\text{read}(x)$ operation is linearized before the $\text{write}(x, _)$ operation from where the $\text{read}(x)$ takes its value (when these two write operations are different). This definition of LWW is inspired by the “bad-

¹For simplicity, we assume that every history contains a set of write operations writing the initial values of variables, which precede every other operation in replica order.

READFROMRGA

$$\begin{aligned} \forall o_2. \text{meth}(o_2) = \text{addAfter} &\Rightarrow \text{arg}_1(o_2) = \circ \vee \\ &\quad \exists o_1. \text{meth}(o_1) = \text{addAfter} \wedge \text{arg}_2(o_1) = \text{arg}_1(o_2) \wedge \text{rf}(o_1, o_2) \\ \wedge \text{meth}(o_2) = \text{remove} &\Rightarrow \exists o_1. \text{meth}(o_1) = \text{addAfter} \wedge \text{arg}_2(o_1) = \text{arg}(o_2) \wedge \text{rf}(o_1, o_2) \\ \wedge \text{meth}(o_2) = \text{read} &\Rightarrow \forall v \in \text{ret}(o_2) \exists o_1. \text{meth}(o_1) = \text{addAfter} \wedge \text{arg}_2(o_1) = v \wedge \text{rf}(o_1, o_2) \end{aligned}$$

RETVALRGA

$$\begin{aligned} \forall o_1, o_2. \text{meth}(o_1) = \text{read} \wedge \text{meth}(o_2) = \text{addAfter} \wedge \text{hb}(o_2, o_1) \wedge \text{arg}_2(o_2) \notin \text{ret}(o_1) \\ \Rightarrow \exists o_3. \text{meth}(o_3) = \text{remove} \wedge \text{arg}(o_3) = \text{arg}_2(o_2) \wedge \text{rf}(o_3, o_1) \end{aligned}$$

LINRGA

$$\begin{aligned} \forall o_1, o_2. (\text{meth}(o_1) = \text{meth}(o_2) = \text{addAfter} \wedge \text{arg}_1(o_1) = \text{arg}_1(o_2) \wedge \\ \exists o_3, o_4, o_5. \text{meth}(o_3) = \text{meth}(o_4) = \text{addAfter} \wedge \text{rf}_{\text{addAfter}}^*(o_1, o_3) \wedge \text{rf}_{\text{addAfter}}^*(o_2, o_4) \wedge \\ \text{meth}(o_5) = \text{read} \wedge \text{arg}_2(o_4) <_{o_5} \text{arg}_2(o_3)) \Rightarrow \text{lin}(o_1, o_2) \end{aligned}$$

Figure 2.3: Axioms used to define the semantics of RGA.

pattern” characterization in [DBLP:conf/popl/BouajjaniEGH17], corresponding to their causal convergence criterion.

2.2.3 REPLICATED COUNTERS

The replicated counter datatype [DBLP:journals/eatcs/ShapiroPBZ11] maintains a set of counters interpreted as integers (the counters can become negative). This datatype provides operations $\text{inc}(x)$ and $\text{dec}(x)$ for incrementing and decrementing a counter x , and $\text{read}(x)$ operations to read the value of the counter x . The semantics of the replicated counter is quite standard: a $\text{read}(x)$ operation returns the value computed as the difference between the number of $\text{inc}(x)$ operations and $\text{dec}(x)$ operations among its happens-before predecessors. The axioms defined below will enforce the fact that a $\text{read}(x)$ operation reads-from all its happens-before predecessors which are $\text{inc}(x)$ or $\text{dec}(x)$ operations.

Therefore, the predicate $\text{ReadOk}(o_1, o_2)$ is defined by

$$\text{meth}(o_1) \in \{\text{inc}, \text{dec}\} \wedge \text{meth}(o_2) = \text{read} \wedge \text{arg}(o_1) = \text{arg}(o_2)$$

and the replicated counter is defined by the following set of axioms:

$$\text{READFROM}(\text{ReadOk}) \wedge \text{CLOSEDRF}(\text{ReadOk}) \wedge \text{RETVALCOUNTER}.$$

2.2.4 REPLICATED GROWABLE ARRAY

The Replicated Growing Array (RGA) [DBLP:journals/jpdc/RohJKL11] is a replicated list used for text-editing applications. RGA supports three operations: $\text{addAfter}(a, b)$ which adds the character b immediately after the occurrence of the character a assumed to be present in the list, $\text{remove}(a)$ which removes a assumed to be present in the list, and $\text{read}()$ which returns the list contents. It is assumed that a character is added at most once². The conflicts between concurrent addAfter operations that add a character immediately after the same character is solved using

²In a practical context, this can be enforced by tagging characters with replica identifiers and sequence numbers.

timestamps (i.e., each added character is associated to a timestamp and the order between characters depends on the order between the corresponding timestamps), which in the axioms below are modeled by the linearization order.

Figure 2.3 lists the axioms defining RGA. `READFROMRGA` ensures that:

- every `addAfter(a,b)` operation reads-from the `addAfter(⋮,a)` adding the character a , except when $a = \circ$ which denotes the “root” element of the list³,
- every `remove(a)` operation reads-from the operation adding a , and
- every read operation returning a list containing a reads-from the operation `addAfter(⋮,a)` adding a .

Then, `RETVLRGA` ensures that a read operation o_1 happening-after an operation adding a character a reads-from a `remove(a)` operation when a doesn’t occur in the list returned by o_1 (the history must contain a `remove(a)` operation because otherwise, a should have occurred in the list returned by the read).

Finally, `LINRGA` models the conflict resolution policy by constraining the linearization order between `addAfter(a,⋮)` operations adding some character immediately after the same character a . As a particular case, `LINRGA` enforces that `addAfter(a,b)` is linearized before `addAfter(a,c)` when a read operation returns a list where c precedes b (`addAfter(a,b)` results in the list $a \cdot b$ and applying `addAfter(a,c)` on $a \cdot b$ results in the list $a \cdot c \cdot b$). However, this is not sufficient: assume that the history contains the two operations `addAfter(a,b)` and `addAfter(a,c)` along with two operations `remove(b)` and `addAfter(b,d)`. Then, a read operation returning the list $a \cdot c \cdot d$ must enforce that `addAfter(a,b)` is linearized before `addAfter(a,c)` because this is the only order between these two operations that can lead to the result $a \cdot c \cdot d$, i.e., executing `addAfter(a,b)`, `addAfter(b,d)`, `remove(b)`, `addAfter(a,c)` in this order. `LINRGA` deals with any scenario where arbitrarily-many characters can be removed from the list: $\text{rf}_{\text{addAfter}}^*$ is the reflexive and transitive closure of the projection of `rf` on `addAfter` operations and $<_{o_5}$ denotes the order between characters in the list returned by the read operation o_5 .

2.3 INTRACTABILITY FOR REGISTERS, SETS, FLAGS, AND COUNTERS

In this section we demonstrate that checking the consistency is intractable for many widely-used data types. While this is not completely unexpected, since some related consistency-checking problems like sequential consistency are also intractable [DBLP:journals/siamcomp/GibbonsK97], this contrasts recent tractability results for checking strong consistency (i.e., linearizability) of common non-replicated data types like sets, maps, and queues [DBLP:journals/pacmpl/EmmiE19]. In fact, in many cases we show that intractability even holds if the number of replicas is fixed.

Our proofs of intractability follow the general structure of Gibbons and Korach’s proofs for the intractability of checking sequential consistency (SC) for atomic registers with read and write operations [DBLP:journals/siamcomp/GibbonsK97]. In particular, we reduce a specialized type of NP-hard propositional satisfiability (SAT) problem to checking whether histories are admitted by a given data type. While our construction borrows from Gibbons and Korach’s, the adaptation

³This element is not returned by read operations.

2 Consistency of Conflict-Free Replicated Data Types

	Replica 0	Replica 1	Replica 2
Round 0	$\left\{ \begin{array}{l} \text{Enable}(x_1) \\ \dots \\ \text{Enable}(x_n) \end{array} \right.$	$\left\{ \begin{array}{l} \text{Disable}(x_1) \\ \dots \\ \text{Disable}(x_n) \end{array} \right.$	
Barrier 1	$\left\{ \begin{array}{l} \text{Enable}(y_0) \\ \text{Read}(y_1) = \text{true} \\ \text{Read}(y_2) = \text{true} \end{array} \right.$	$\left\{ \begin{array}{l} \text{Enable}(y_1) \\ \text{Read}(y_0) = \text{true} \\ \text{Read}(y_2) = \text{true} \end{array} \right.$	$\left\{ \begin{array}{l} \text{Enable}(y_2) \\ \text{Read}(y_0) = \text{true} \\ \text{Read}(y_1) = \text{true} \end{array} \right.$
Round 1	$\left\{ \begin{array}{l} \text{Read}(\alpha_1) = \text{true} \\ \text{Read}(\beta_1) = \text{false} \\ \text{Read}(\gamma_1) = \text{false} \\ \text{Disable}(\alpha_1) \\ \text{Enable}(\beta_1) \end{array} \right.$	$\left\{ \begin{array}{l} \text{Read}(\alpha_1) = \text{false} \\ \text{Read}(\beta_1) = \text{true} \\ \text{Read}(\gamma_1) = \text{false} \\ \text{Disable}(\beta_1) \\ \text{Enable}(\gamma_1) \end{array} \right.$	$\left\{ \begin{array}{l} \text{Read}(\alpha_1) = \text{false} \\ \text{Read}(\beta_1) = \text{false} \\ \text{Read}(\gamma_1) = \text{true} \\ \text{Disable}(\gamma_1) \\ \text{Enable}(\alpha_1) \end{array} \right.$
Barrier 2	$\left\{ \begin{array}{l} \text{Disable}(y_0) \\ \text{Read}(y_1) = \text{false} \\ \text{Read}(y_2) = \text{false} \end{array} \right.$	$\left\{ \begin{array}{l} \text{Disable}(y_1) \\ \text{Read}(y_0) = \text{false} \\ \text{Read}(y_2) = \text{false} \end{array} \right.$	$\left\{ \begin{array}{l} \text{Disable}(y_2) \\ \text{Read}(y_0) = \text{false} \\ \text{Read}(y_1) = \text{false} \end{array} \right.$
...
Round m	$\left\{ \begin{array}{l} \text{Read}(\alpha_m) = \text{true} \\ \text{Read}(\beta_m) = \text{false} \\ \text{Read}(\gamma_m) = \text{false} \\ \text{Disable}(\alpha_m) \\ \text{Enable}(\beta_m) \end{array} \right.$	$\left\{ \begin{array}{l} \text{Read}(\alpha_m) = \text{false} \\ \text{Read}(\beta_m) = \text{true} \\ \text{Read}(\gamma_m) = \text{false} \\ \text{Disable}(\beta_m) \\ \text{Enable}(\gamma_m) \end{array} \right.$	$\left\{ \begin{array}{l} \text{Read}(\alpha_m) = \text{false} \\ \text{Read}(\beta_m) = \text{false} \\ \text{Read}(\gamma_m) = \text{true} \\ \text{Disable}(\gamma_m) \\ \text{Enable}(\alpha_m) \end{array} \right.$

Figure 2.4: The encoding of a 1-in-3 SAT problem $\bigwedge_{i=1}^m (\alpha_i \vee \beta_i \vee \gamma_i)$ over variables x_1, \dots, x_n as a 3-replica history of a flag data type. Besides the flag variable x_j for each propositional variable x_j , the encoding adds per-replica variables y_j for synchronization barriers.

from SC to CRDT consistency requires a significant extension to handle the consistency relaxation represented by abstract executions: rather than a direct sequencing of threads' operations, CRDT consistency requires the construction of three separate relations: read-from, happens-before, and linearization.

Technically, our reductions start from the 1-in-3 SAT problem [DBLP:books/fm/GareyJ79]: given a propositional formula $\bigwedge_{i=1}^m (\alpha_i \vee \beta_i \vee \gamma_i)$ over variables x_1, \dots, x_n with only positive literals, i.e., $\alpha_i, \beta_i, \gamma_i \in \{x_1, \dots, x_n\}$, does there exist an assignment to the variables such that exactly one of $\alpha_i, \beta_i, \gamma_i$ per clause is assigned *true*? The proofs of Theorems 2.3.1 and 2.3.2 reduce 1-in-3 SAT to CRDT consistency checking.

Theorem 2.3.1. *The admissibility problem is NP-hard when the number of replicas is fixed for the following data types: Add-Wins Set, Remove-Wins Set, Enable-Wins Flag, Disable-Wins Flag, Multi-Value Register, and Last-Writer-Wins Register.*

Proof. We demonstrate a reduction from the 1-in-3 SAT problem. For a given problem $p = \bigwedge_{i=1}^m (\alpha_i \vee \beta_i \vee \gamma_i)$ over variables x_1, \dots, x_n , we construct a 3-replica history h_p of the flag data type — either enable- or disable-wins — as illustrated in Figure 2.4. The encoding includes a flag variable x_j for each propositional variable x_j , along with a per-replica flag variable y_j used to implement synchronization barriers. Intuitively, executions of h_p proceed in $m + 1$ rounds: the first round corresponds to the assignment of a truth valuation, while subsequent rounds check the validity of each clause given the assignment. The reductions to sets and registers are slight variations on this proof, in which the Read, Enable, and Disable operations are replaced with Contains, Add, and Remove, respectively, and Read and Writes of values 1 and 0, respectively.

It suffices to show that the constructed history h_p is admitted if and only if the given problem p is satisfiable. Since the flag data type does not constrain the linearization relation of its abstract executions, we regard only the read-from and happens-before components. It is straightforward to verify that the happens-before relations of h_p 's abstract executions necessarily order:

1. every pair of operations in distinct rounds — due to barriers; and
2. every operation in a given round, over all replicas, without interleaving the operations of distinct replicas within the same round — since a replica's reads in a given round are only consistent with the other replicas' after the re-enabling and -disabling of flag variables.

In other words, replicas appear to execute atomically per round, in a round-robin fashion. Furthermore, since all operations in a given round happen before the operations of subsequent rounds, the values of flag variables are consistent across rounds — i.e., as read by the first replica to execute in a given round — and determined in the initial round either by conflict resolution — i.e., enable-or-disable-wins — or by happens-before, in case conflict resolution would have been inconsistent with subsequent reads.

In the “if” direction, let $\vec{r} \in \{0, 1, 2\}^m$ be the positions of positively-assigned variables in each clause, e.g., $r_i = 0$ implies $\alpha_i = \text{true}$ and $\beta_i = \gamma_i = \text{false}$. We construct an abstract execution $e_{\vec{r}}$ in which the happens-before relation sequences the operations of replica r_i before those of $r_i + 1 \bmod 3$, and in turn before $r_i + 2 \bmod 3$. In other words, the replicas in round i appear to execute in left-to-right order from starting with the replica r_i , whose reads correspond to the satisfying assignment of $(\alpha_i \vee \beta_i \vee \gamma_i)$. The read-from relation of $e_{\vec{r}}$ relates each $\text{Read}(x_j) = \text{true}$ operation to the most recent $\text{Enable}(x_j)$ operation in happens-before order, which is unique since happens-before sequences the operations of all rounds; the case for $\text{Read}(x_j) = \text{false}$ and $\text{Disable}(x_j)$ is symmetric. It is then straightforward to verify that $e_{\vec{r}}$ satisfies the axioms of the enable- or disable-wins flag, and thus h_p is admitted.

In the “only if” direction, let e be an abstract execution of h_p , and let $\vec{r} \in \{0, 1, 2\}^m$ be the replicas first to execute in each round according to the happens-before order of e . It is straightforward to verify that the assignment in which a given variable is set to true iff the replica encoding its positive assignment in some clause executes first in its round, i.e.,

$$x_j = \begin{cases} \text{true} & \text{if } \exists i. (r_i = 0 \wedge \alpha_i = x_j) \vee (r_i = 1 \wedge \beta_i = x_j) \vee (r_i = 2 \wedge \gamma_i = x_j) \\ \text{false} & \text{otherwise,} \end{cases}$$

is a satisfying assignment to p . □ □

Theorem 2.3.1 establishes intractability of consistency for the aforementioned sets, flags, and registers, independently from the number of replicas. In contrast, our proof of Theorem 2.3.2 for counter data types depends on the number of replicas, since our encoding requires two replicas per propositional variable. Intuitively, since counter increments and decrements are commutative, the initial round in the previous encoding would have fixed all counter values to zero. Instead, the next encoding isolates initial increments and decrements to independent replicas. The weaker result is indeed tight since checking counter consistency with a fixed number of replicas is polynomial time, as Section 2.5 demonstrates.

Theorem 2.3.2. *The admissibility problem for the Counter data type is NP-hard.*

2 Consistency of Conflict-Free Replicated Data Types

	Replica 0	Replica $2j+1$	Replica $2j+2$
Round 0 {	Read(y) = n	Inc(y) Inc(x_j)	Inc(y) Dec(x_j)
Barrier 1 {	Inc(z) Read(z) = 3	Replica 1 Inc(z) Read(z) = 3	Replica 2 Inc(z) Read(z) = 3
Round 1 {	Read(α_1) = 1 Read(β_1) = -1 Read(γ_1) = -1 Dec(α_1); Dec(α_1) Inc(β_1); Inc(β_1)	Read(α_1) = -1 Read(β_1) = 1 Read(γ_1) = -1 Dec(β_1); Dec(β_1) Inc(γ_1); Inc(γ_1)	Read(α_1) = -1 Read(β_1) = -1 Read(γ_1) = 1 Dec(γ_1); Dec(γ_1) Inc(α_1); Inc(α_1)
Barrier 2 {	Dec(z) Read(z) = 0	Dec(z) Read(z) = 0	Dec(z) Read(z) = 0
...
Round m {	Read(α_m) = 1 Read(β_m) = -1 Read(γ_m) = -1 Dec(α_m); Dec(α_m) Inc(β_m); Inc(β_m)	Read(α_m) = -1 Read(β_m) = 1 Read(γ_m) = -1 Dec(β_m); Dec(β_m) Inc(γ_m); Inc(γ_m)	Read(α_m) = -1 Read(β_m) = -1 Read(γ_m) = 1 Dec(γ_m); Dec(γ_m) Inc(α_m); Inc(α_m)
Barrier $m+1$ {	Inc(z) or Dec(z) Read(z) = 3 or 0	Inc(z) or Dec(z) Read(z) = 3 or 0	Inc(z) or Dec(z) Read(z) = 3 or 0
Round $m+1$ {	Read(y) = n		

Figure 2.5: The encoding of a 1-in-3 SAT problem $\bigwedge_{i=1}^m (\alpha_i \vee \beta_i \vee \gamma_i)$ over variables x_1, \dots, x_n as the history of a counter over $2n + 3$ replicas. Besides the counter variables x_j encoding propositional variables x_j , the encoding adds a variable y encoding the number of initial increments and decrements, and a variable z to implement synchronization barriers.

Proof. We demonstrate a reduction from the 1-in-3 SAT problem. For a given problem $p = \bigwedge_{i=1}^m (\alpha_i \vee \beta_i \vee \gamma_i)$ over variables x_1, \dots, x_n , we construct a history h_p of the counter data type over $2n + 3$ replicas, as illustrated in Figure 2.5.

Besides the differences imposed due to the commutativity of counter increments and decrements, our reduction follows the same strategy as in the proof of Theorem 2.3.1: the happens-before relation of h_p 's abstract executions order every pair of operations in distinct rounds (of Replicas 0–2), and every operation in a given (non-initial) round. As before, Replicas 0–2 appear to execute atomically per round, in a round-robin fashion, and counter variables are consistent across rounds. The key difference is that here abstract executions' happens-before relations only relate the operations of either Replica $2j+1$ or $2j+2$, for each $j = 1, \dots, n$, to operations in subsequent rounds: the other's operations are never observed by other replicas. Our encoding ensures that exactly one of each is observed by ensuring that the counter y is incremented exactly n times — and relying on the fact that every variable appears in some clause, so that a read that observed neither or both would yield the value zero, which is inconsistent with h_p . Otherwise, our reasoning follows the proof of Theorem 2.3.1, in which the read-from relation selects all increments and decrements of the same counter variable in happens-before order. \square \square

Input: A differentiated history $h = \langle \text{Op}, \text{ro} \rangle$ and a datatype T .

Output: *true* iff h satisfies the axioms of T .

```

1 rf ← ComputeRF( $h, \text{READFROM}[T], \text{RETVAL}[T]$ );
2 if rf =  $\perp$  then return false;
3 hb ← ( $\text{ro} \cup \text{rf}$ )+;
4 if hb is cyclic or  $\langle h, \text{rf}, \text{hb} \rangle \not\models \text{READFROMMAXIMAL}[T] \wedge \text{READALLMAXIMALS}[T]$  then
5   | return false;
6 lin ← hb;
7 lin ← LinClosure(hb, LIN[ $T$ ]);
8 if lin is cyclic then return false;
9 return true;

```

Algorithm 1: Consistency checking for RGA, LWW, and MVR. $\text{RE}\dots[T]$ refers to an axiom of T , or *true* when T lacks such an axiom. The relation R^+ denotes the transitive closure of R .

2.4 POLYNOMIAL-TIME ALGORITHMS FOR REGISTERS AND ARRAYS

We show that the problem of checking consistency is polynomial time for RGA, and even for LWW and MVR under the assumption that each value is written at most once, i.e., for each value v , the input history contains at most one write operation $\text{write}(x, v)$. Histories satisfying this assumption are called *differentiated*. The latter is a restriction motivated by the fact that practical implementations of these datatypes are data-independent [DBLP:conf/popl/Wolper86], i.e., their behavior doesn't depend on the concrete values read or written and any potential buggy behavior can be exposed in executions where each value is written at most once. Also, in a testing environment, this restriction can be enforced by tagging each value with a replica identifier and a sequence number.

In all three cases, the feature that enables polynomial time consistency checking is the fact that the read-from relation becomes fixed for a given history, i.e., if the history is consistent, then there exists exactly one read-from relation rf that satisfies the READFROM_- and RETVAL_- axioms, and rf can be derived syntactically from the operation labels (using those axioms). Then, our axiomatic characterizations enable a consistency checking algorithm which roughly, consists in instantiating those axioms in order to compute an abstract execution.

The consistency checking algorithm for RGA, LWW, and MVR is listed in Algorithm 1. It computes the three relations rf , hb , and lin of an abstract execution using the datatype's axioms. The history is declared consistent iff there exist satisfying rf and hb relations, and the relations hb and lin computed this way are acyclic. The acyclicity requirement comes from the definition of abstract executions where hb and lin are required to be partial/total orders. While an abstract execution would require that lin is a total order, this algorithm computes a partial linearization order. However, any total order compatible with this partial linearization would satisfy the axioms of the datatype.

ComputeRF computes the read-from relation rf satisfying the READFROM_- and RETVAL_- axioms. In the case of LWW and MVR, it defines rf as the set of all pairs formed of $\text{write}(x, v)$ and $\text{read}(x)$ operations where v belongs to the return value of the read. By RETVAL_- , each $\text{read}(x)$ operation must be associated to at least one $\text{write}(x, _)$ operation. Also, the fact that each value is written at most once implies that this rf relation is uniquely defined, e.g., for LWW, it is not possible to find two write operations that could be rf related to the same read operation. In gen-

<p>Input: A history $h = \langle \text{Op}, \text{ro} \rangle$ of RGA.</p> <p>Output: An rf satisfying $\text{READFROMRGA} \wedge \text{RETVALRGA}$, if exists; \perp o/w</p> <pre> 1 rf ← { (o₁, o₂) : meth(o₁) = addAfter, meth(o₂) ∈ {addAfter, remove, read}, arg₂(o₁) = arg₁(o₂) ∨ arg₂(o₁) ∈ ret(o₂) }; 2 if ⟨h, rf⟩ ⊭ READFROMRGA then return ⊥ ; 3 while true do 4 rf₁ ← ∅; 5 foreach o₁, o₂ ∈ Op s.t. ⟨o₂, o₁⟩ ∈ (rf ∪ ro)⁺ and meth(o₁) = read and meth(o₂) = addAfter and arg₂(o₂) ∉ ret(o₁) do 6 if ∃ o₃ ∈ Op s.t. meth(o₃) = remove and arg(o₃) = arg₂(o₂) then 7 rf₁ ← rf₁ ∪ {⟨o₃, o₁⟩}; 8 else 9 return ⊥; 10 if rf₁ ⊆ rf then break; 11 else rf ← rf ∪ rf₁; 12 return rf; </pre>

Algorithm 2: The procedure ComputeRF for RGA.

eral, if there exists no rf relation satisfying these axioms, then ComputeRF returns a distinguished value \perp to signal a consistency violation. Note that the computation of the read-from for LWW and MVR is quadratic time⁴ since the constraints imposed by the axioms relate only to the operation labels, the methods they invoke or their arguments. The case of RGA is slightly more involved because the axiom RETVALRGA introduces more read-from constraints based on the happens-before order which includes ro and the rf itself. In this case, the computation of rf relies on a fixpoint computation, which converges in at most quadratic time (the maximal size of rf), described in Algorithm 2. Essentially, we use the axiom READFROMRGA to populate the read-from relation and then, apply the axiom RETVALRGA iteratively, using the read-from constraints added in previous steps, until the computation converges.

After computing the read-from relation, our algorithm defines the happens-before relation hb as the transitive closure of ro union rf. This is sound because none of the axioms of these datatypes enforce new happens-before constraints, which are not already captured by ro and rf. Then, it checks whether the hb defined this way is acyclic and satisfies the datatype's axioms that constrain hb, i.e., READFROMMAXIMAL and READALLMAXIMALS (when they are present).

Finally, in the case of LWW and RGA, the algorithm computes a (partial) linearization order that satisfies the corresponding LIN_ axioms. Starting from an initial linearization order which is exactly the happens-before, it computes new constraints by instantiating the universally quantified axioms LINLWW and LINRGA. Since these axioms are not “recursive”, i.e., they don't enforce linearization order constraints based on other linearization order constraints, a standard instantiation of these axioms is enough to compute a partial linearization order such that any extension to a total order satisfies the datatype's axioms.

Theorem 2.4.1. *Algorithm 1 returns true iff the input history is consistent.*

The following holds because Algorithm 1 runs in polynomial time — the rank depends on the number of quantifiers in the datatype's axioms. Indeed, Algorithm 1 represents a least fixpoint

⁴ Assuming constant time lookup/insert operations (e.g., using hashmaps), this complexity is linear time.

computation which converges in at most a quadratic number of iterations (the maximal size of rf).

Corollary 2.4.1. *The admissibility problem is polynomial time for RGA, and for LWW and MVR on differentiated histories.*

2.5 POLYNOMIAL-TIME ALGORITHMS FOR REPLICATED COUNTERS

In this section, we show that checking consistency for the replicated counter datatype becomes polynomial time assuming the number of replicas in the input history is fixed (i.e., the width of the replica order ro is fixed). We present an algorithm which constructs a valid happens-before order (note that the semantics of the replicated counter doesn't constrain the linearization order) incrementally, following the replica order. At any time, the happens-before order is uniquely determined by a *prefix mapping* that associates to each replica a *prefix* of the history, i.e., a set of operations which is downward-closed w.r.t. replica order (i.e., if it contains an operation it contains all its ro predecessors). This models the fact that the replica order is included in the happens-before and therefore, if an operation o_1 happens-before another operation o_2 , then all the ro predecessors of o_1 happen-before o_2 . The happens-before order can be extended in two ways: (1) adding an operation issued on the replica i to the prefix of replica i , or (2) “merging” the prefix of a replica j to the prefix of a replica i (this models the delivery of an operation issued on replica j and all its happens-before predecessors to the replica i). Verifying that an extension of the happens-before is valid, i.e., that the return values of newly-added read operations satisfy the RETVALCOUNTER axiom, doesn't depend on the happens-before order between the operations in the prefix associated to some replica (it is enough to count the inc and dec operations in that prefix). Therefore, the algorithm can be seen as a search in the space of prefix mappings. If the number of replicas in the input history is fixed, then the number of possible prefix mappings is polynomial in the size of the history, which implies that the search can be done in polynomial time.

Let $h = (Op, ro)$ be a history. To simplify the notations, we assume that the replica order is a union of sequences, each sequence representing the operations issued on the same replica. Therefore, each operation $o \in Op$ is associated with a replica identifier $rep(o) \in [1..n_h]$, where n_h is the number of replicas in h .

A *prefix* of h is a set of operation $Op' \subseteq Op$ such that all the ro predecessors of operations in Op' are also in Op' , i.e., $\forall o \in Op. ro^{-1}(o) \in Op'$. Note that the union of two prefixes of h is also a prefix of h . The *last operation* of replica i in a prefix Op' is the ro -maximal operation o with $rep(o) = i$ included in Op' . A prefix Op' is called *valid* if (Op', ro') , where ro' is the projection of ro on Op' , is admitted by the replicated counter.

A *prefix map* is a mapping m which associates a prefix of h to each replica $i \in [1..n_h]$. Intuitively, a prefix map defines for each replica i the set of operations which are “known” to i , i.e., happen-before the last operation of i in its prefix. Formally, a prefix map m is *included* in a happens-before relation hb , denoted by $m \subseteq hb$, if for each replica $i \in [1..n_h]$, $hb(o, o_i)$ for each operation in $o \in m(i) \setminus \{o_i\}$, where o_i is the last operation of i in $m(i)$. We call o_i the *last operation* of i in m , and denoted it by $last_i(m)$. A prefix map m is *valid* if it associates a valid prefix to each replica, and *complete* if it associates the whole history h to each replica i .

Algorithm 3 lists our algorithm for checking consistency of replicated counter histories. It is defined as a recursive procedure `checkCounter` that searches for a sequence of valid extensions

<p>Input: History $h = (\text{Op}, \text{ro})$, prefix map m, and set seen of invalid prefix maps</p> <p>Output: <i>true</i> iff there exists read-from and happens-before relations rf and hb such that $m \subseteq \text{hb}$, and $\langle h, \text{rf}, \text{hb} \rangle$ satisfies the counter axioms.</p> <pre> 1 if m is complete then return true; 2 foreach replica i do 3 foreach replica $j \neq i$ do 4 $m' \leftarrow m[i \leftarrow m(i) \cup m(j)]$; 5 if $m' \notin \text{seen}$ and checkCounter(h, m', seen) then 6 return true; 7 $\text{seen} \leftarrow \text{seen} \cup \{m'\}$; 8 if $\exists o_1. \text{ro}^1(\text{last}_i(m), o_1)$ then 9 if $\text{meth}(o_1) = \text{read}$ and 1 $\arg(o_1) = x \wedge \text{ret}(o_1) \neq \{o \in m[i] \mid o = \text{inc}(x)\} - \{o \in m[i] \mid o = \text{dec}(x)\}$ then 11 return false; 12 $m' \leftarrow m[i \leftarrow m(i) \cup \{o_1\}]$; 13 if $m' \notin \text{seen}$ and checkCounter(h, m', seen) then 14 return true; 15 $\text{seen} \leftarrow \text{seen} \cup \{m'\}$; 16 return false; </pre>

Algorithm 3: The procedure checkCounter, where ro^1 denotes immediate ro-successor, and $f[a \leftarrow b]$ updates function f with mapping $a \mapsto b$.

of a given prefix map (initially, this prefix map is empty) until it becomes complete. The axiom RETVALCOUNTER is enforced whenever extending the prefix map with a new read operation (when the last operation of a replica i is “advanced” to a read operation). The following theorem states of the correctness of the algorithm.

Theorem 2.5.1. *checkCounter(h, \emptyset, \emptyset) returns true iff the input history is consistent.*

When the number of replicas is fixed, the number of prefix maps becomes polynomial in the size of the history. This follows from the fact that prefixes are uniquely defined by their ro-maximal operations, whose number is fixed.

Corollary 2.5.1. *The admissibility problem for replicated counters is polynomial-time when the number of replicas is fixed.*

2.6 POLYNOMIAL-TIME ALGORITHMS FOR SETS AND FLAGS

While Theorem 2.3.1 shows that the admissibility problem is NP-complete for replicated sets and flags even if the number of replicas is fixed, we show that this problem becomes polynomial time when additionally, the number of values added to the set, or the number of flags, is also fixed. Note that this doesn’t limit the number of operations in the input history which can still be arbitrarily large. In the following, we focus on the Add-Wins Set, the other cases being very similar.

We propose an algorithm for checking consistency which is actually an extension of the one presented in Section 2.5 for replicated counters. The additional complexity in checking consistency for the Add-Wins Set comes from the validity of contains(x) return values which requires identifying the maximal predecessors in the happens-before relation that add or remove x (which are not necessarily the maximal hb-predecessors all-together). In the case of counters, it was enough just

to count happens-before predecessors. Therefore, we extend the algorithm for replicated counters such that along with the prefix map, we also keep track of the hb-maximal $\text{add}(x)$ and $\text{remove}(x)$ operations for each element x and each replica i . When extending a prefix map with a contains operation, these hb-maximal operations (which define a witness for the read-from relation) are enough to verify the RetValSet axiom. Extending the prefix of a replica with an add or remove operation (issued on the same replica), or by merging the prefix of another replica, may require an update of these hb-maximal predecessors.

When the number of replicas and elements are fixed, the number of read-from maps is polynomial in the size of the history — recall that the number of operations associated by a read-from map to a replica and set element is bounded by the number of replicas. Combined with the number of prefix maps being polynomial when the number of replicas is fixed, we obtain the following result.

Theorem 2.6.1. *Checking whether a history is admitted by the Add-Wins Set, Remove-Wins Set, Enable-Wins Flag, or the Disable-Wins Flag is polynomial time provided that the number of replicas and elements/flags is fixed.*

2.7 RELATED WORK

Many have considered consistency models applicable to CRDTs, including causal consistency [DBLP:journals/cacm/Lamport79], sequential consistency [DBLP:journals/tc/Lamport79], linearizability [DBLP:journals/toplas/HerlihyW90], session consistency [DBLP:conf/pdis/TerryDPSTW94], eventual consistency [DBLP:conf/sosp/TerryTPDSH95], and happens-before consistency [DBLP:conf/popl/MansonPA05]. Burckhardt et al. [DBLP:journals/ftpl/Burckhardt14, DBLP:conf/popl/BurckhardtGYZ14] propose a unifying framework to formalize these models. Many have also studied the complexity of verifying data-type agnostic notions of consistency, including serializability, sequential consistency and linearizability [DBLP:journals/jacm/Papadimitriou79b, DBLP:journals/siamcomp/GibbonsK97, DBLP:journals/iandc/AlurMP00, DBLP:conf/spaa/BinghamCH03, DBLP:conf/cav/FarzanM08, DBLP:conf/esop/BouajjaniEEH13, DBLP:conf/netys/Hamza15], as well as causal consistency [DBLP:conf/popl/BouajjaniEGH17]. Our definition of the replicated LWW register corresponds to the notion of causal convergence in [DBLP:conf/popl/BouajjaniEGH17]. This work studies the complexity of the admissibility problem for the replicated LWW register. It shows that this problem is NP-complete in general and polynomial time when each value is written only once. Our NP-completeness result is stronger since it assumes a fixed number of replicas, and our algorithm for the case of unique values is more general and can be applied uniformly to MVR and RGA. While Bouajjani et al. [DBLP:journals/iandc/BouajjaniEEH18, DBLP:journals/pacmpl/EmmiE18] consider the complexity for individual linearizable collection types, we are the first to establish (in)tractability of individual replicated data types. Others have developed effective consistency checking algorithms for sequential consistency [DBLP:conf/cav/HenzingerQR99a, DBLP:journals/tpds/Qadeer03, DBLP:conf/cav/BinghamCHQZ04, DBLP:conf/pldi/BurckhardtAM07], serializability [DBLP:conf/fmcad/0002OPTZ07, DBLP:conf/cav/FarzanM08, DBLP:conf/pldi/GuerraouiHJS08, DBLP:conf/pldi/EmmiMM10], linearizability [DBLP:journals/jpdc/WingG93, DBLP:conf/pldi/BurckhardtDMT10, DBLP:conf/pldi/EmmiEH15, DBLP:journals/concurrency/Lowe17], and even weaker notions like eventual consistency [DBLP:conf/popl/BouajjaniEH14] and sequential happens-before consistency [DBLP:conf/cav/EmmiE18].

DBLP:journals/pacmpl/EmmiE19]. In contrast, we are the first to establish precise polynomial-time algorithms for runtime verification of replicated data types.

2.8 CONCLUSION

By developing novel logical characterizations of replicated data types, reductions from propositional satisfiability checking, and tractable algorithms, we have established a frontier of tractability for checking consistency of replicated data types. As far as we are aware, our results are the first to characterize the asymptotic complexity consistency checking for CRDTs.

3 TRANSACTIONAL CONSISTENCY

3.1 INTRODUCTION

Transactions simplify concurrent programming by enabling computations on shared data that are isolated from other concurrent computations and resilient to failures. Modern databases provide transactions in various forms corresponding to different tradeoffs between consistency and availability. The strongest level of consistency is achieved with *serializable* transactions [DBLP:journals/jacm/Papadimitriou79b] whose outcome in concurrent executions is the same as if the transactions were executed atomically in some order. Unfortunately, serializability carries a significant penalty on the availability of the system assuming, for instance, that the database is accessed over a network that can suffer from partitions or failures. For this reason, modern databases often provide weaker guarantees about transactions, formalized by weak consistency models, e.g., causal consistency [DBLP:journals/cacm/Lamport78] and snapshot isolation [DBLP:conf/sigmod/BerensonBGM0095].

Implementations of large-scale databases providing transactions are difficult to build and test. For instance, distributed (replicated) databases must account for partial failures, where some components or the network can fail and produce incomplete results. Ensuring fault-tolerance relies on intricate protocols that are difficult to design and reason about. The black-box testing framework Jepsen [jepsen] found a remarkably large number of subtle problems in many production distributed databases.

Testing a transactional database raises two issues: (1) deriving a suitable set of testing scenarios, e.g., faults to inject into the system and the set of transactions to be executed, and (2) deriving efficient algorithms for checking whether a given execution satisfies the considered consistency model. The Jepsen framework aims to address the first issue by using randomization, e.g., introducing faults at random and choosing the operations in a transaction randomly. The effectiveness of this approach has been proved formally in recent work [DBLP:journals/pacmpl/OzkanMNBW18]. The second issue is, however, largely unexplored. Jepsen checks consistency in a rather ad-hoc way, focusing on specific classes of violations to a given consistency model, e.g., dirty reads (reading values from aborted transactions). This problem is challenging because the consistency specifications are non-trivial and they cannot be checked using, for instance, standard local assertions added to the client's code.

Besides serializability, the complexity of checking correctness of an execution w.r.t. some consistency model is unknown. Checking serializability has been shown to be NP-complete [DBLP:journals/jacm/Papadimitriou79] and checking causal consistency in a *non-transactional* context is known to be polynomial time [DBLP:conf/popl/BouajjaniEG]. In this work, we try to fill this gap by investigating the complexity of this problem w.r.t. several consistency models and, in the case of NP-completeness, devising algorithms that are polynomial time assuming fixed bounds for certain parameters of the input executions, e.g., the number of sessions.

We consider several consistency models that are the most prevalent in practice. The weakest of them, *Read Committed* (RC) [DBLP:conf/sigmod/BerensonBGM0095], requires that every value read in a transaction is written by a committed transaction. *Read Atomic* (RA) [DBLP:conf/concur/Cerone0G15] requires that successive reads of the same variable in a transaction return the same value (also known as Repeatable Reads [DBLP:conf/sigmod/BerensonBGM0095]), and that a transaction “sees” the values written by previous transactions in the same session. In general, we assume that transactions are organized in *sessions* [DBLP:conf/pdis/TerryDPSTW94], an abstraction of the sequence of transactions performed during the execution of an application. *Causal Consistency* (CC) [DBLP:journals/cacm/Lamport78] requires that if a transaction t_1 “affects” another transaction t_2 , e.g., t_1 is ordered before t_2 in the same session or t_2 reads a value written by t_1 , then these two transactions are observed by any other transaction in this order. *Prefix Consistency* (PC) [DBLP:conf/ecoop/BurckhardtLPF15] requires that there exists a total commit order between all the transactions such that each transaction observes a prefix of this sequence. *Snapshot Isolation* (SI) [DBLP:conf/sigmod/BerensonBGM0095] further requires that two different transactions observe different prefixes if they both write to a common variable. Finally, we also provide new results concerning the problem of checking serializability (SER) that complement the known result about its NP-completeness.

The algorithmic issues we explore in this paper have led to a new specification framework for these consistency models that relies on the fact that the *write-read* relation in an execution (also known as *read-from*), relating reads with the transactions that wrote their value, can be defined effectively. The write-read relation can be extracted easily from executions where each value is written at most once (a variable can be written an arbitrary number of times). This can be easily enforced by tagging values with unique identifiers (e.g., a local counter that is incremented with every new write coupled with a client/session identifier)¹. Since practical database implementations are data-independent [DBLP:conf/popl/Wolper86], i.e., their behavior doesn’t depend on the concrete values read or written in the transactions, any potential buggy behavior can be exposed in executions where each value is written at most once. Therefore, this assumption is without loss of generality.

Previous work [DBLP:conf/popl/BouajjaniEGH17, DBLP:conf/popl/BurckhardtGYZ14, DBLP:conf/concur/Ceron] has formalized such consistency models using two auxiliary relations: a *visibility* relation defining for each transaction the set of transactions it observes, and a *commit order* defining the order in which transactions are committed to the “global” memory. An execution satisfying some consistency model is defined as the existence of a visibility relation and a commit order obeying certain axioms. In our case, the write-read relation derived from the execution plays the role of the visibility relation. This simplification allows us to state a series of axioms defining these consistency models, which have a common shape. Intuitively, they define lower bounds on the set of transactions t_1 that *must* precede in commit order a transaction t_2 that is read in the execution. Besides shedding a new light on the differences between these consistency models, these axioms are essential for the algorithmic issues we investigate afterwards.

We establish that checking whether an execution satisfies RC, RA, or CC is polynomial time, while the same problem is NP-complete for PC and SI. Moreover, in the case of the NP-complete consistency models (PC, SI, SER), we show that their verification problem becomes polynomial

¹This is also used in Jepsen, e.g., checking dirty reads in Galera [jepsen-galera].

time provided that, roughly speaking, the number of sessions in the input executions is considered to be fixed (i.e., not counted for in the input size). In more detail, we establish that checking SER reduces to a search problem in a space that has polynomial size when the number of sessions is fixed. (This algorithm applies to arbitrary executions, but its complexity would be exponential in the number of sessions in general.) Then, we show that checking PC or SI can be reduced in polynomial time to checking SER using a transformation of executions that, roughly speaking, splits each transaction in two parts: one part containing all the reads, and one part containing all the writes (SI further requires adding some additional variables in order to deal with transactions writing on a common variable). We extend these results even further by relying on an abstraction of executions called *communication graphs* [DBLP:journals/pacmpl/ChalupaCPSV18]. Roughly speaking, the vertices of a communication graph correspond to sessions, and the edges represent the fact that two sessions access (read or write) the same variable. We show that all these criteria are polynomial-time checkable provided that the *biconnected* components of the communication graph are of fixed size.

We provide an experimental evaluation of our algorithms on executions of CockroachDB [cockroach], which claims to implement serializability [cockroach-claim] acknowledging however the possibility of anomalies, Galera [galera], whose documentation contains contradicting claims about whether it implements snapshot isolation [galera-claim, galera-notclaim], and AntidoteDB [antidote], which claims to implement causal consistency [antidote-claim]. Our implementation reports violations of these criteria in all cases. The consistency violations we found for AntidoteDB are novel and have been confirmed by its developers. We show that our algorithms are efficient and scalable. In particular, we show that, although the asymptotic complexity of our algorithms is exponential in general (w.r.t. the number of sessions), the worst-case behavior is not exercised in practice.

To summarize, the contributions of this work are fourfold:

- We develop a new specification framework for describing common transactional-consistency criteria (§3.3);
- We show that checking RC, RA, and CC is polynomial time while checking PC and SI is NP-complete (§3.4);
- We show that PC, SI, and SER are polynomial-time checkable assuming that the communication graph of the input execution has fixed-size biconnected components (§3.5 and §3.6);
- We perform an empirical evaluation of our algorithms on executions generated by production databases (§3.7);

Combined, these contributions form an effective algorithmic framework for the verification of transactional-consistency models. To the best of our knowledge, we are the first to investigate the asymptotic complexity for most of these consistency models, despite their prevalence in practice.

Additional material can be found in [DBLP:journals/corr/abs-1908-04509].

3.2 OVERVIEW

When a client sends a transactional request to a database, typically she receives an update to the request, *e.g.* for a write request, she may receive `true` or `false` representing whether the database was

3 Transactional Consistency

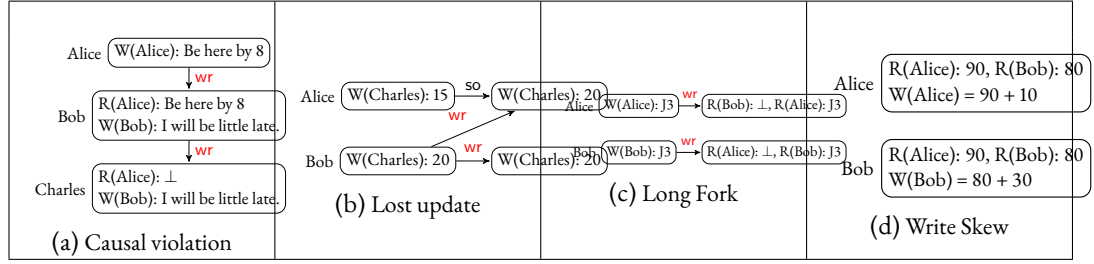


Table 3.1: Anomalies in database applications

able to make the changes. To process that request, the database may need to do very complicated work - there can be parallel requests or the data may need to be replicated in different locations. But to a client's perspective, the database should simply behave in a way, as if all requests are committed to the database in a *meaningful* sequence which is *consistent* to what she saw.

Whatever the clients saw, we call it a History. Given a history, we try to see if the database behaved in a way where the *commit sequence* is *consistent* to the History. We call a history and its possible commit sequence an Execution.

When we say, the execution is consistent, we mean, the commit sequence satisfies some *criteria* with the history. A programmer may want to have all the transactions to be committed in series - Serializability. But in practical, a programmer usually only need something weak, *e.g.* Read Committed, Snapshot Isolation - but it depends on what kind of behavior the programmer expects.

Let's look at some anomalies that can happen in database applications. Alice, Bob and Charles are in a Whatsapp group. Alice asks everyone to come at 8 o'clock. Bob sees it and replies, he will be late. But as illustrated in 3.1a, Charles cannot see Alice's message but sees Bob's message and thinks Bob is saying he will be late for nothing. This anomaly is called *Causality Violation*.

Even such anomalies can happen in bank operations. Say Alice and Bob are trying to send money to Charles. Alice sends \$15 and Bob sends \$20. But as illustrated in 3.1b, Alice's deposit is lost. This anomaly is called *Lost Update*.

Suppose Alice and Bob trying to book a cinema ticket. First Alice and Bob both publish they want to book a seat and check the response of other person so that they know they have a conflict. But as illustrated in 3.1c, they can't see each other's reponse and ends up booking the same seat. This anomaly is called *Long Fork*.

Consider Amazon is giving out cashback for purchases. But it has to also check sum of total cashback is not more than \$200 before giving away a cashback. Otherwise, it may end up giving away too much cashback to make profit. So if Alice and Bob both become eligible for a cashback, as illustrated in 3.1d, even if the total cashback value crosses the limit, the situation can let both of them have the cashback. This anomaly is called *Write Skew*.

Modern databases support different levels of consistency, allowing the programmer to modify the isolation level according to her concurrent and consistency need.

First we will formally define history. Then we will define some consistency axioms and criteria over history. Then we deduce the complexity classes to verify history against respective consistency criteria. We will see Serializability, Snapshot Isolation and Prefix consistency are in NP-Complete.

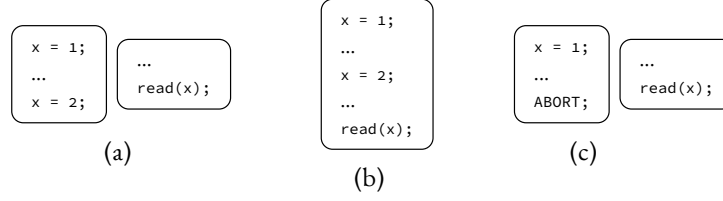


Figure 3.1: Examples of transactions used to justify our simplifying assumptions (each box represents a different transaction): (a) only the last written value is observable in other transactions, (b) reads following writes to the same variable return the last written value in the same transaction, and (c) values written in aborted transactions are not observable.

Then we will show for fixed number of sessions, the respective problems are in NL complexity class.

3.3 CONSISTENCY CRITERIA

3.3.1 HISTORIES

We consider a transactional database storing a set of variables $\text{Var} = \{x, y, \dots\}$. Clients interact with the database by issuing transactions formed of read and write operations. Assuming an unspecified set of values Val and a set of operation identifiers Old , we let

$$\text{Op} = \{\text{rd}[i]xv, \text{write}_i(x, v) : i \in \text{Old}, x \in \text{Var}, v \in \text{Val}\}$$

be the set of operations reading a value v or writing a value v to a variable x . We omit operation identifiers when they are not important.

Definition 3.3.1. A transaction $\langle O, \text{po} \rangle$ is a finite set of operations O along with a strict total order po on O , called program order.

We use t, t_1, t_2, \dots to range over transactions. The set of read, resp., write, operations in a transaction t is denoted by $\text{reads}(t)$, resp., $\text{writes}(t)$. The extension to sets of transactions is defined as usual. Also, we say that a transaction t writes a variable x , denoted by t writes x , when $\text{write}_i(x, v) \in \text{writes}(t)$ for some i and v . Similarly, a transaction t reads a variable x when $\text{rd}[i]xv \in \text{reads}(t)$ for some i and v .

To simplify the exposition, we assume that each transaction t contains at most one write operation to each variable², and that a read of a variable x cannot be preceded by a write to x in the same transaction³. If a transaction would contain multiple writes to the same variable, then only the last one should be visible to other transactions (w.r.t. any consistency criterion considered in practice). For instance, the $\text{read}(x)$ in Figure 3.1a should not return 1 because this is not the last value written to x by the other transaction. It can return the initial value or 2. Also, if a read would be preceded by a write to the same variable in the same transaction, then it should return a value written in the same transaction (i.e., the value written by the latest write to x in that transaction).

²That is, for every transaction t , and every $\text{write}(x, v), \text{write}(y, v') \in \text{writes}(t)$, we have that $x \neq y$.

³That is, for every transaction $t = \langle O, \text{po} \rangle$, if $\text{write}(x, v) \in \text{writes}(t)$ and there exists $\text{rd}[i]xv \in \text{reads}(t)$, then we have that $\langle \text{rd}[i]xv, \text{write}(x, v) \rangle \in \text{po}$

For instance, the $\text{read}(x)$ in Figure 3.1b can only return 2 (assuming that there are no other writes on x in the same transaction). These two properties can be verified easily (in a syntactic manner) on a given execution. Beyond these two properties, the various consistency criteria used in practice constrain only the last writes to each variable in each transaction and the reads that are not preceded by writes to the same variable in the same transaction.

Consistency criteria are formalized on an abstract view of an execution called *history*. A history includes only successful or committed transactions. In the context of databases, it is always assumed that the effect of aborted transactions should not be visible to other transactions, and therefore, they can be ignored. For instance, the $\text{read}(x)$ in Figure 3.1c should not return the value 1 written by the aborted transaction. The transactions are ordered according to a (partial) *session order* so which represents ordering constraints imposed by the applications using the database. Most often, so is a union of sequences, each sequence being called a *session*. We assume that the history includes a *write-read* relation that identifies the transaction writing the value returned by each read in the execution. As mentioned before, such a relation can be extracted easily from executions where each value is written at most once. Since in practice, databases are data-independent [DBLP:conf/popl/Wolper86], i.e., their behavior does not depend on the concrete values read or written in the transactions, any potential buggy behavior can be exposed in such executions.

Definition 3.3.2. A history $\langle T, so, wr \rangle$ is a set of transactions T along with a strict partial order so called session order, and a relation $wr \subseteq T \times \text{reads}(T)$ called write-read relation, s.t.

- the inverse of wr is a total function, and if $(t, {}^{\text{rd}}xv) \in wr$, then $\text{write}(x, v) \in t$, and
- $so \cup wr$ is acyclic.

To simplify the technical exposition, we assume that every history includes a distinguished transaction writing the initial values of all variables. This transaction precedes all the other transactions in so . We use h, h_1, h_2, \dots to range over histories.

We say that the read operation ${}^{\text{rd}}xv$ reads value v from variable x written by t when $(t, {}^{\text{rd}}xv) \in wr$. For a given variable x , wr_x denotes the restriction of wr to reads of variable x , i.e., $wr_x = wr \cap (T \times \{{}^{\text{rd}}xv \mid v \in \text{Val}\})$. Moreover, we extend the relations wr and wr_x to pairs of transactions as follows: $\langle t_1, t_2 \rangle \in wr$, resp., $\langle t_1, t_2 \rangle \in wr_x$, iff there exists a read operation ${}^{\text{rd}}xv \in \text{reads}(t_2)$ such that $\langle t_1, {}^{\text{rd}}xv \rangle \in wr$, resp., $\langle t_1, {}^{\text{rd}}xv \rangle \in wr_x$. We say that the transaction t_1 is *read* by the transaction t_2 when $\langle t_1, t_2 \rangle \in wr$, and that it is *read* when it is read by some transaction t_2 .

3.3.2 AXIOMATIC FRAMEWORK

We describe an axiomatic framework to characterize the set of histories satisfying a certain consistency criterion. The overarching principle is to say that a history satisfies a certain criterion if there exists a strict total order on its transactions, called *commit order* and denoted by co , which extends the write-read relation and the session order, and which satisfies certain properties. These properties are expressed by a set of axioms that relate the commit order with the session-order and the write-read relation in the history.

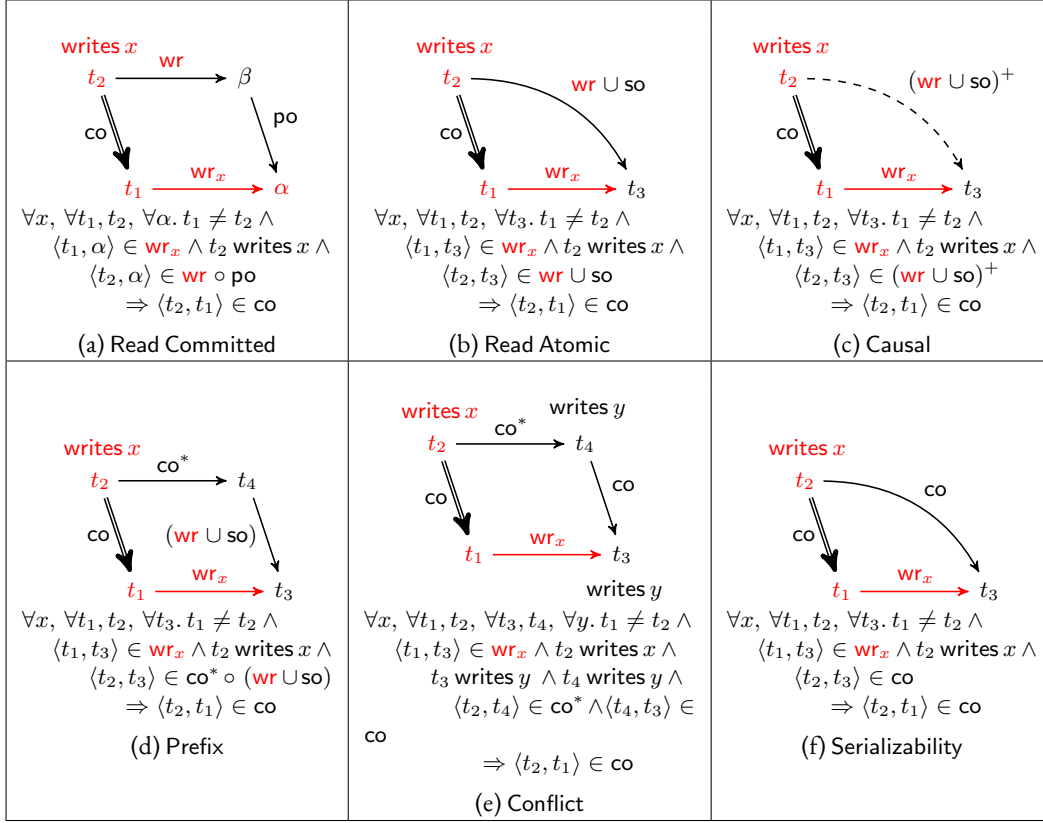


Figure 3.2: Definitions of consistency axioms. The reflexive and transitive, resp., transitive, closure of a relation rel is denoted by rel^* , resp., rel^+ . Also, \circ denotes the composition of two relations, i.e., $rel_1 \circ rel_2 = \{\langle a, b \rangle \mid \exists c. \langle a, c \rangle \in rel_1 \wedge \langle c, b \rangle \in rel_2\}$.

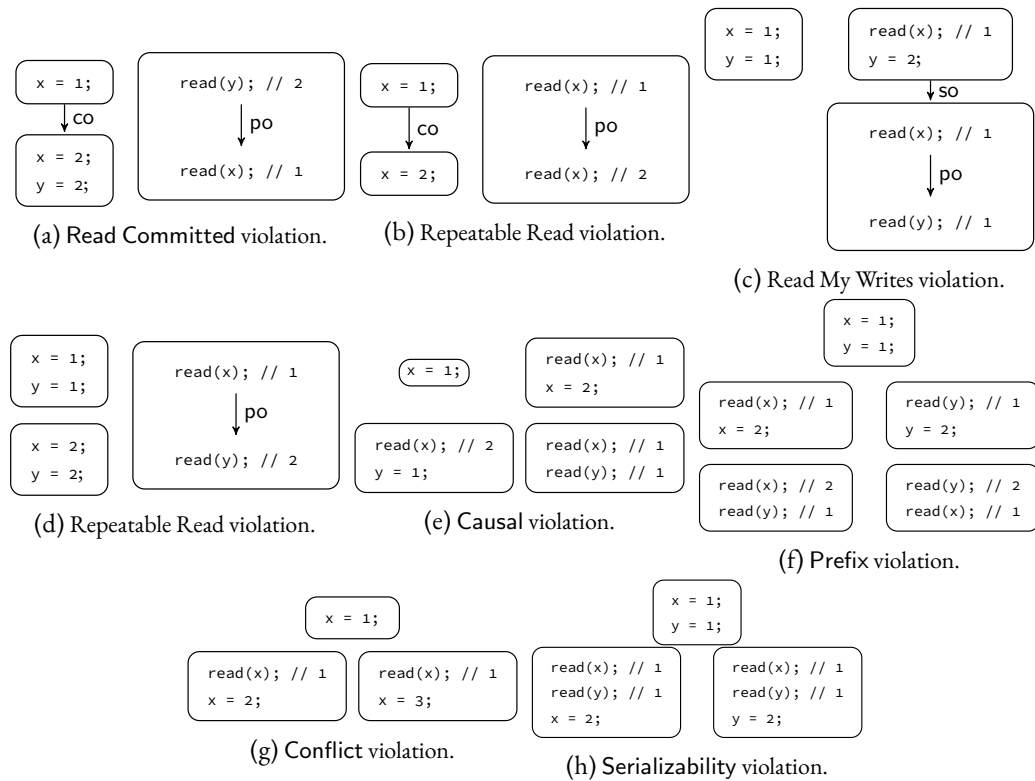


Figure 3.3: Examples of histories used to explain the axioms in Figure 3.2. For readability, the **wr** relation is defined by the values written in comments with each read.

The axioms we use have a uniform shape: they define mandatory co predecessors t_2 of a transaction t_1 that is read in the history. For instance, the criterion called **READ COMMITTED (RC)** [DBLP:conf/sigmod/BerensonBGMOO95] requires that every value read in the history was written by a committed transaction, and also, that the reads in the same transaction are “monotonic” in the sense that they do not return values that are older, w.r.t. the commit order, than other values read in the past⁴. While the first condition holds for any history (because of the surjectivity of wr), the second condition is expressed by the axiom **Read Committed** in Figure 3.2a. This axiom states that for any transaction t_1 writing a variable x that is read in a transaction t , the set of transactions t_2 writing x and read previously in the same transaction must precede t_1 in commit order. For instance, Figure 3.3a shows a history and a (partial) commit order that does not satisfy this axiom because $\text{read}(x)$ returns the value written in a transaction “older” than the transaction read in the previous $\text{read}(y)$. An example of a history and commit order satisfying this axiom is given in Figure 3.3b.

More precisely, the axioms are first-order formulas⁵ of the following form:

$$\forall x, \forall t_1, t_2, \forall \alpha. t_1 \neq t_2 \wedge \langle t_1, \alpha \rangle \in \text{wr}_x \wedge t_2 \text{ writes } x \wedge \phi(t_2, \alpha) \Rightarrow \langle t_2, t_1 \rangle \in \text{co}$$

where ϕ is a property relating t_2 and α (i.e., the read or the transaction reading from t_1) that varies from one axiom to another. Intuitively, this axiom schema states the following: in order for α to read specifically t_1 ’s write on x , it must be the case that every t_2 that also writes x and satisfies $\phi(t_2, \alpha)$ was committed before t_1 . Note that in all cases we consider, $\phi(t_2, \alpha)$ already ensures that t_2 is committed before the read α , so this axiom schema ensures that t_2 is furthermore committed before t_1 ’s write.

The axioms used throughout the paper are given in Figure 3.2. The property ϕ relates t_2 and α using the write-read relation and the session order in the history, and the commit order.

In the following, we explain the rest of the consistency criteria we consider and the axioms defining them. **READ ATOMIC (RA)** [DBLP:conf/concur/Cerone06] is a strengthening of **READ COMMITTED** defined by the axiom **Read Atomic**, which states that for any transaction t_1 writing a variable x that is read in a transaction t_3 , the set of wr or so predecessors of t_3 writing x must precede t_1 in commit order. The case of wr predecessors corresponds to the **Repeatable Read** criterion in [DBLP:conf/sigmod/BerensonBGMOO95], which requires that successive reads of the same variable in the same transaction return the same value, Figure 3.3b showing a violation, and also that every read of a variable x in a transaction t returns the value written by the maximal transaction t' (w.r.t. the commit order) that is read by t , Figure 3.3d showing a violation (for any commit order between the transactions on the left, either $\text{read}(x)$ or $\text{read}(y)$ will return a value not written by the maximal transaction). The case of so predecessors corresponds to the “read-my-writes” guarantee [DBLP:conf/pdis/TerryDPSTW94] concerning sessions, which states that a transaction t must observe previous writes in the same session. For instance, $\text{read}(y)$ returning 1 in Figure 3.3c shows that the last transaction on the right does not satisfy this guarantee: the transaction writing 1 to y was already visible to that session before it wrote 2 to y , and therefore

⁴This monotonicity property corresponds to the fact that in the original formulation of **READ COMMITTED** [DBLP:conf/sigmod/BerensonBGMOO95], every write is guarded by the acquisition of a lock on the written variable, that is held until the end of the transaction.

⁵These formulas are interpreted on tuples $\langle h, \text{co} \rangle$ of a history h and a commit order co on the transactions in h as usual.

Table 3.2: Consistency model definitions

Consistency model	Axioms
READ COMMITTED (RC)	Read Committed
READ ATOMIC (RA)	Read Atomic
CAUSAL CONSISTENCY (CC)	Causal
PREFIX CONSISTENCY (PC)	Prefix
SNAPSHOT ISOLATION (SI)	Prefix \wedge Conflict
SERIALIZABILITY (SER)	Serializability

the value 2 should have been read. Read Atomic requires that the so predecessor of the transaction reading y be ordered in co before the transaction writing 1 to y , which makes the union $\text{co} \cup \text{wr}$ cyclic.

The following lemma shows that for histories satisfying Read Atomic, the inverse of wr_x extended to transactions is a total function.

Lemma 3.3.1. *Let $h = \langle T, \text{so}, \text{wr} \rangle$ be a history. If $\langle h, \text{co} \rangle$ satisfies Read Atomic, then for every transaction t and two reads ${}^{\text{rd}}[i_1]xv_1, {}^{\text{rd}}[i_2]xv_2 \in \text{reads}(t)$, $\text{wr}^{-1}({}^{\text{rd}}[i_1]xv_1) = \text{wr}^{-1}({}^{\text{rd}}[i_2]xv_2)$ and $v_1 = v_2$.*

Proof. Let $\langle t_1, {}^{\text{rd}}[i_1]xv_1 \rangle, \langle t_2, {}^{\text{rd}}[i_2]xv_2 \rangle \in \text{wr}_x$. Then t_1, t_2 write to x . Let us assume by contradiction, that $t_1 \neq t_2$. By Read Atomic, $\langle t_2, t_1 \rangle \in \text{co}$ because $\langle t_1, {}^{\text{rd}}[i_1]xv_1 \rangle \in \text{wr}_x$ and t_2 writes to x . Similarly, we can also show that $\langle t_1, t_2 \rangle \in \text{co}$. This contradicts the fact that co is a strict total order. Therefore, $t_1 = t_2$. We also have that $v_1 = v_2$ because each transaction contains a single write to x . \square

CAUSAL CONSISTENCY (CC) [DBLP:journals/cacm/Lamport78] is defined by the axiom Causal, which states that for any transaction t_1 writing a variable x that is read in a transaction t_3 , the set of $(\text{wr} \cup \text{so})^+$ predecessors of t_3 writing x must precede t_1 in commit order ($(\text{wr} \cup \text{so})^+$ is usually called the *causal* order). A violation of this axiom can be found in Figure 3.3e: the transaction t_2 writing 2 to x is a $(\text{wr} \cup \text{so})^+$ predecessor of the transaction t_3 reading 1 from x because the transaction t_4 , writing 1 to y , reads x from t_2 and t_3 reads y from t_4 . This implies that t_2 should precede in commit order the transaction t_1 writing 1 to x , which again, is inconsistent with the write-read relation (t_2 reads from t_1).

PREFIX CONSISTENCY (PC) [DBLP:conf/ecoop/BurckhardtLPF15] is a strengthening of CC, which requires that every transaction observes a prefix of a commit order between all the transactions. With the intuition that the observed transactions are $\text{wr} \cup \text{so}$ predecessors, the axiom Prefix defining PC, states that for any transaction t_1 writing a variable x that is read in a transaction t_3 , the set of co^* predecessors of transactions observed by t_3 writing x must precede t_1 in commit order (we use co^* to say that even the transactions observed by t_3 must precede t_1). This ensures the prefix property stated above. An example of a PC violation can be found in Figure 3.3f: the two transactions on the bottom read from the three transactions on the top, but any serialization of those three transactions will imply that one of the combinations $x=1, y=2$ or $x=2, y=1$ cannot be produced at the end of a prefix in this serialization.

SNAPSHOT ISOLATION (SI) [DBLP:conf/sigmod/BerensonBGM0095] is a strengthening of PC that disallows two transactions to observe the same prefix of a commit order if they *conflict*, i.e., write to a common variable. It is defined by the conjunction of Prefix and another axiom called Conflict, which requires that for any transaction t_1 writing a variable x that is read in a transaction t_3 , the set of co^* predecessors writing x of transactions conflicting with t_3 and before t_3 in commit order, must precede t_1 in commit order. Figure 3.3g shows a Conflict violation.

Finally, SERIALIZABILITY (SER) [DBLP:journals/jacm/Papadimitriou79b] is defined by the axiom with the same name, which requires that for any transaction t_1 writing to a variable x that is read in a transaction t_3 , the set of co predecessors of t_3 writing x must precede t_1 in commit order. This ensures that each transaction observes the effects of all the co predecessors. Figure 3.3h shows a Serializability violation.

The next lemma states the relationship between these axioms.

Lemma 3.3.2. *The following entailments hold:*

$$\text{Causal} \Rightarrow \text{Read Atomic} \Rightarrow \text{Read Committed}$$

$$\text{Prefix} \Rightarrow \text{Causal}$$

$$\text{Serializability} \Rightarrow \text{Prefix} \wedge \text{Conflict}$$

Proof. We will show the contrapositive of each implication:

- If $\langle h, \text{co} \rangle$ does not satisfy Read Committed, then

$$\exists x, \exists t_1, t_2, \exists \alpha, \beta. \langle t_1, \alpha \rangle \in \text{wr}_x \wedge t_2 \text{ writes } x \wedge \langle t_2, \beta \rangle \in \text{wr} \wedge \langle \beta, \alpha \rangle \in \text{po} \wedge \langle t_1, t_2 \rangle \in \text{co}.$$

Let t_3 the transaction containing α and β . We have that $\langle t_2, t_3 \rangle \in \text{wr}$. But then we have t_1, t_2, t_3 such that $\langle t_1, t_3 \rangle \in \text{wr}_x$ and $\langle t_2, t_3 \rangle \in \text{wr}$ and t_2 writes x . So by Read Atomic, $\langle t_2, t_1 \rangle \in \text{co}$. This contradicts the fact that co is a strict total order. Therefore, $\langle h, \text{co} \rangle$ does not satisfy Read Atomic.

- If $\langle h, \text{co} \rangle$ does not satisfy Read Atomic, then

$$\exists x, \exists t_1, t_2, t_3. \langle t_1, t_3 \rangle \in \text{wr}_x \wedge t_2 \text{ writes } x \wedge \langle t_2, t_3 \rangle \in \text{wr} \cup \text{so} \wedge \langle t_1, t_2 \rangle \in \text{co}.$$

Then $\langle t_2, t_3 \rangle \in (\text{wr} \cup \text{so})^+$. Then, by Causal, we have $\langle t_2, t_1 \rangle \in \text{co}$, which contradicts the fact that co is a strict total order. Therefore, $\langle h, \text{co} \rangle$ does not satisfy Causal.

- If $\langle h, \text{co} \rangle$ does not satisfy Causal, then

$$\exists x, \exists t_1, t_2, t_3. \langle t_1, t_3 \rangle \in \text{wr}_x \wedge t_2 \text{ writes } x \wedge \langle t_2, t_3 \rangle \in (\text{wr} \cup \text{so})^+ \wedge \langle t_1, t_2 \rangle \in \text{co}.$$

But, $(\text{wr} \cup \text{so})^+ = (\text{wr} \cup \text{so})^* \circ (\text{wr} \cup \text{so}) \subseteq \text{co}^* \circ (\text{wr} \cup \text{so})$. Therefore, $\langle t_2, t_3 \rangle \in \text{co}^* \circ (\text{wr} \cup \text{so})$. Then, by Prefix, we have $\langle t_2, t_1 \rangle \in \text{co}$, which contradicts the fact that co is a strict total order. Therefore, $\langle h, \text{co} \rangle$ does not satisfy Prefix.

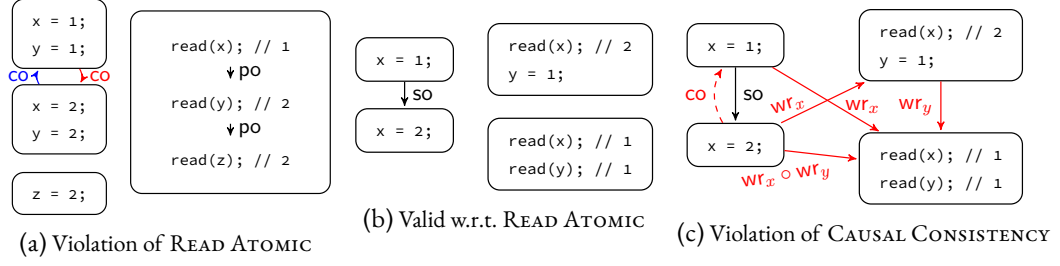


Figure 3.4: Applying the RA and CC checking algorithms.

- If $\langle h, co \rangle$ does not satisfy Prefix or Conflict, then

$$\exists x, \exists t_1, t_2, t_3, t_4. \langle t_1, t_3 \rangle \in \text{wr}_x \wedge t_2 \text{ writes } x \wedge \langle t_2, t_4 \rangle \in co^* \wedge \langle t_1, t_2 \rangle \in co$$

and

- $\langle t_4, t_3 \rangle \in co \wedge t_3 \text{ writes } y \wedge t_3 \text{ writes } y$ if it violates Conflict.
- $\langle t_4, t_3 \rangle \in (\text{wr} \cup so)$ if it violates Prefix.

In both cases, we have that $\langle t_4, t_3 \rangle \in co$. Because co is transitive, $\langle t_2, t_4 \rangle \in co^*$ and $\langle t_4, t_3 \rangle \in co$ imply that $\langle t_2, t_3 \rangle \in co$. Then by Serializability, we have $\langle t_2, t_1 \rangle \in co$, which contradicts the fact that co is a strict total order. Therefore, $\langle h, co \rangle$ does not satisfy Serializability.

□

Definition 3.3.3. Given a set of axioms X defining a criterion C like in Table 3.2, a history $h = \langle T, so, wr \rangle$ satisfies C iff there exists a strict total order co such that $wr \cup so \subseteq co$ and $\langle h, co \rangle$ satisfies X .

Definition 3.3.3 and Lemma 3.3.2 imply that each consistency criterion in Table 3.2 is stronger than its predecessors (reading them from top to bottom), e.g., CC is stronger than RA and RC. This relation is known to be strict [DBLP:conf/concur/Cerone0G15], e.g., RA is not stronger than CC.

3.4 CHECKING CONSISTENCY CRITERIA

This section establishes the complexity of checking the different consistency criteria in Table 3.2 for a given history. More precisely, we show that READ COMMITTED, READ ATOMIC, and CAUSAL CONSISTENCY can be checked in polynomial time while the problem of checking the rest of the criteria is NP-complete.

Intuitively, the polynomial time results are based on the fact that the axioms defining those consistency criteria do not contain the commit order (co) on the left-hand side of the entailment. Therefore, proving the existence of a commit order satisfying those axioms can be done using a saturation procedure that builds a “partial” commit order based on instantiating the axioms on the write-read relation and the session order in the given history. Since the commit order must

Input: A history $h = \langle T, \text{so}, \text{wr} \rangle$
Output: *true* iff h satisfies CAUSAL CONSISTENCY

```

1 if  $\text{so} \cup \text{wr}$  is cyclic then
2   | return false;
3  $\text{co} \leftarrow \text{so} \cup \text{wr}$ ;
4 foreach  $x \in \text{vars}(h)$  do
5   | foreach  $t_1 \neq t_2 \in T$  s.t.  $t_1$  and  $t_2$  write  $x$  do
6     | | if  $\exists t_3. \langle t_1, t_3 \rangle \in \text{wr}_x \wedge \langle t_2, t_3 \rangle \in (\text{so} \cup \text{wr})^+$  then
7       | | |  $\text{co} \leftarrow \text{co} \cup \{\langle t_2, t_1 \rangle\}$ ;
8 if  $\text{co}$  is cyclic then
9   | return false;
10 else
11   | return true;

```

Algorithm 4: Checking CAUSAL CONSISTENCY.

be an extension of the write-read relation and the session order, it contains those two relations from the beginning. This saturation procedure stops when the order constraints derived this way become cyclic. For instance, let us consider applying such a procedure corresponding to RA on the histories in Figure 3.4a and Figure 3.4b. Applying the axiom in Figure 3.2b on the first history, since the transaction on the right reads 2 from y , we get that its wr_x predecessor (i.e., the first transaction on the left) must precede the transaction writing 2 to y in commit order (the red edge). This holds because the wr_x predecessor writes on y . Similarly, since the same transaction reads 1 from x , we get that its wr_y predecessor must precede the transaction writing 1 to x in commit order (the blue edge). This already implies a cyclic commit order, and therefore, this history does not satisfy RA. On the other hand, for the history in Figure 3.4b, all the axiom instantiations are vacuous, i.e., the left part of the entailment is false, and therefore, it satisfies RA. Checking CC on the history in Figure 3.4c requires a single saturation step: since the transaction on the bottom right reads 1 from x , its $\text{wr}_x \circ \text{wr}_y$ predecessor that writes on x (the transaction on the bottom left) must precede in commit order the transaction writing 1 to x . Since this is already inconsistent with the session order, we get that this history violates CC.

Algorithm 4 lists our procedure for checking CC. As explained above, co is initially set to $\text{so} \cup \text{wr}$, and then, it is saturated with other ordering constraints implied by non-vacuous instantiations of the axiom Causal (where the left-hand side of the implication evaluates to true). The algorithms concerning RC and RA are defined in a similar way by essentially changing the test at line 6 so that it corresponds to the left-hand side of the implication in the corresponding axiom. Algorithm 4 can be rewritten as a Datalog program containing straightforward Datalog rules for computing transitive closures and relation composition, and a rule of the form⁶

$$\langle t_2, t_1 \rangle \in \text{co} :- t_1 \neq t_2, \langle t_1, t_3 \rangle \in \text{wr}_x, \langle t_2, t_3 \rangle \in (\text{so} \cup \text{wr})^+$$

to represent the Causal axiom. The following is a consequence of the fact that these algorithms run in polynomial time (or equivalently, the corresponding Datalog programs can be evaluated in polynomial time over a database that contains the wr and so relations in a given history).

⁶We write Datalog rules using a standard notation $\text{head} :- \text{body}$ where head is a relational atom (written as $\langle a, b \rangle \in R$ where a, b are elements and R a binary relation) and body is a list of relational atoms.

Theorem 3.4.1. *For any criterion $C \in \{\text{READ COMMITTED}, \text{READ ATOMIC}, \text{CAUSAL CONSISTENCY}\}$, the problem of checking whether a given history satisfies C is polynomial time.*

On the other hand, checking PC, SI, and SER is NP-complete in general. We show this using a reduction from boolean satisfiability (SAT) that covers uniformly all the three cases. In the case of SER, it provides a new proof of the NP-completeness result by [DBLP:journals/jacm/Papadimitriou79b], which uses a reduction from the so-called *non-circular* SAT and which cannot be extended to PC and SI.

Theorem 3.4.2. *For any criterion $C \in \{\text{PREFIX CONSISTENCY}, \text{SNAPSHOT ISOLATION}, \text{SERIALIZABILITY}\}$ the problem of checking whether a given history satisfies C is NP-complete.*

Proof. Given a history, any of these three criteria can be checked by guessing a total commit order on its transactions and verifying whether it satisfies the corresponding axioms. This shows that the problem is in NP.

To show NP-hardness, we define a reduction from boolean satisfiability. Therefore, let $\varphi = D_1 \wedge \dots \wedge D_m$ be a CNF formula over the boolean variables x_1, \dots, x_n where each D_i is a disjunctive clause with m_i literals. Let λ_{ij} denote the j -th literal of D_i .

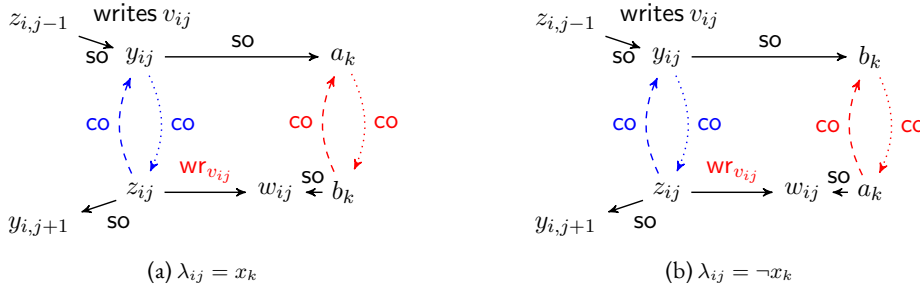
We construct a history h_φ such that φ is satisfiable if and only if h_φ satisfies PC, SI, or SER. Since $\text{SER} \Rightarrow \text{SI} \Rightarrow \text{PC}$, we show that (1) if h_φ satisfies PC, then φ is satisfiable, and (2) if φ is satisfiable, then h_φ satisfies SER.

CONSTRUCTION OF h_φ The main idea of the construction is to represent truth values of each of the variables and literals in φ with the polarity of the commit order between corresponding transaction pairs. For each variable x_k , h_φ contains a pair of transactions a_k and b_k , and for each literal λ_{ij} , h_φ contains a set of transactions w_{ij} , y_{ij} and z_{ij} ⁷. We want to have that x_k is false if and only if $\langle a_k, b_k \rangle \in \text{co}$, and λ_{ij} is false if and only if $\langle y_{ij}, z_{ij} \rangle \in \text{co}$ (the transaction w_{ij} is used to "synchronize" the truth value of the literals with that of the variables, which is explained later).

The history h_φ should ensure that the co ordering constraints corresponding to an assignment that falsifies the formula (*i.e.* one of its clauses) form a cycle. To achieve that, we add all pairs $\langle z_{ij}, y_{i,(j+1)\%m_i} \rangle$ in the session order so. An unsatisfied clause D_i , *i.e.* every λ_{ij} is false, leads to a cycle of the form $y_{i1} \xrightarrow{\text{co}} z_{i1} \xrightarrow{\text{so}} y_{i2} \xrightarrow{\text{co}} z_{i2} \dots z_{im_i} \xrightarrow{\text{so}} y_{i1}$.

The most complicated part of the construction is to ensure some consistency between the truth value of the literals and the truth value of the variables, e.g., $\lambda_{ij} = x_k$ is true iff x_k is true, for at least one literal λ_{ij} interpreted as true in every clause D_i (if such a literal exists). Figure 3.5a shows the sub-history associated to a positive literal $\lambda_{ij} = x_k$ while Figure 3.5b shows the case of a negative literal $\lambda_{ij} = \neg x_k$. For a positive literal $\lambda_{ij} = x_k$ (Figure 3.5a), (1) we enrich session order with the pairs $\langle y_{ij}, a_k \rangle$ and $\langle b_k, w_{ij} \rangle$, (2) we include writes to a variable v_{ij} in the transactions y_{ij} and z_{ij} , and (3) we make w_{ij} read from z_{ij} , *i.e.* $\langle z_{ij}, w_{ij} \rangle \in \text{wrr}_{v_{ij}}$. The case of a negative literal is similar, switching the roles of a_k and b_k .

⁷We assume that the transactions a_k and b_k associated to a variable x_k are distinct and different from the transactions associated to another variable $x_{k'} \neq x_k$ or to a literal λ_{ij} . Similarly, for the transactions w_{ij} , y_{ij} and z_{ij} associated to a literal λ_{ij} .


 Figure 3.5: Sub-histories included in h_φ for each literal λ_{ij} and variable x_k .

PC FOR h_φ IMPLIES SATISFIABILITY OF φ If h_φ satisfies PC, then there exists a total commit order co between the transactions described above, which together with h_φ satisfies Prefix. We show that the assignment of the variables x_k explained above (defined by the co order between a_k and b_k , for each k) satisfies the formula φ . For each clause D_i , the so constraints between the transactions y_{ij} , z_{ij} with $1 \leq j \leq m_i$ imply that there exist some z_{ij} that is committed before its corresponding y_{ij} . These two transactions are included in the sub-history corresponding to the literal λ_{ij} (Figure 3.5a or Figure 3.5b depending on the polarity of the literal).

The definition of this sub-history ensures that the interpretation to true of the literal λ_{ij} (given by the order in co between z_{ij} and y_{ij}) is consistent with the assignment of the variable it contains (defined by the co order between a_k , b_k). More precisely, it ensures that if the co goes upwards on the left-hand side ($\langle z_{ij}, y_{ij} \rangle \in \text{co}$) like in this case, then it must also go upwards on the right-hand side ($\langle b_k, a_k \rangle \in \text{co}$ in the case of a positive literal, and $\langle a_k, b_k \rangle \in \text{co}$ in the case of a negative literal) to satisfy Prefix. For instance, if $\lambda_{ij} = x_k$ is a positive literal and we assume by contradiction that $\langle a_k, b_k \rangle \in \text{co}$, then $\langle y_{ij}, w_{ij} \rangle \in \text{so} \circ \text{co} \circ \text{so}$. Therefore, for every commit order co such that $\langle h_\varphi, \text{co} \rangle$ satisfies Prefix, $\langle a_k, b_k \rangle \in \text{co}$ implies $\langle y_{ij}, z_{ij} \rangle \in \text{co}$, which contradicts the hypothesis. Indeed, if $\langle a_k, b_k \rangle \in \text{co}$, instantiating the Prefix axiom where y_{ij} plays the role of t_2 , z_{ij} plays the role of t_1 , and w_{ij} plays the role of t_3 , we obtain that $\langle y_{ij}, z_{ij} \rangle \in \text{co}$.

Therefore, the assignment of the variables x_k leads to at least one literal interpreted to true in each clause D_i , and the formula φ is satisfiable.

SATISFIABILITY OF φ IMPLIES SER FOR h_φ Let γ be a satisfying assignment for φ . Also, let co' be a binary relation that includes so and wr such that if $\gamma(x_k) = \text{false}$, then $\langle a_k, b_k \rangle \in \text{co}'$, $\langle y_{ij}, z_{ij} \rangle \in \text{co}'$ for each $\lambda_{ij} = x_k$, and $\langle z_{ij}, y_{ij} \rangle \in \text{co}'$ for each $\lambda_{ij} = \neg x_k$, and if $\gamma(x_k) = \text{true}$, then $\langle b_k, a_k \rangle \in \text{co}'$, $\langle z_{ij}, y_{ij} \rangle \in \text{co}'$ for each $\lambda_{ij} = x_k$, and $\langle y_{ij}, z_{ij} \rangle \in \text{co}'$ for each $\lambda_{ij} = \neg x_k$. Looking at the sub-histories corresponding to literals λ_{ij} (Figure 3.5a or Figure 3.5b), co' goes in the same direction (upwards or downwards) on both sides.

Note that co' is acyclic: no cycle can contain w_{ij} because w_{ij} has no “outgoing” dependency (i.e. co' contains no pair with w_{ij} as a first component), there is no cycle including some pair of transactions a_k , b_k and some pair y_{ij} , z_{ij} because there is no way to reach y_{ij} or z_{ij} from a_k or b_k , there is no cycle including only transactions a_k and b_k because a_{k_1} and b_{k_1} are not related to a_{k_2} and b_{k_2} , for $k_1 \neq k_2$, there is no cycle including transactions y_{i_1,j_1} , z_{i_1,j_1} and y_{i_2,j_2} , z_{i_2,j_2} for $i_1 \neq i_2$ since these are disconnected as well, and finally, there is no cycle including only transactions y_{ij} and z_{ij} , for a fixed i , because φ is satisfiable. It can be proved easily that the

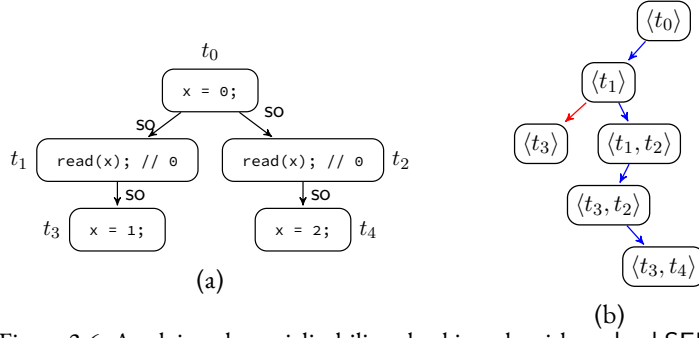


Figure 3.6: Applying the serializability checking algorithm `checkSER` (Algorithm 5) on the serializable history on the left. The right part pictures a search for valid extensions of serializable prefixes, represented by their boundaries. The red arrow means that the search is blocked (the prefix at the target is not a valid extension), while blue arrows mean that the search continues.

acyclic relation co' can be extended to a total commit order co which together with h_φ satisfies the Serializability axiom. Therefore, h_φ satisfies SER. \square

3.5 CHECKING CONSISTENCY OF BOUNDED-WIDTH HISTORIES

In this section, we show that checking prefix consistency, snapshot isolation, and serializability becomes polynomial time under the assumption that the *width* of the given history, i.e., the maximum number of mutually-unordered transactions w.r.t. the session order, is bounded by a fixed constant. If we consider the standard case where the session order is a union of transaction sequences (modulo the fictitious transaction writing the initial values), i.e., a set of sessions, then the width of the history is the number of sessions. We start by presenting an algorithm for checking serializability that is polynomial time when the width is bounded by a fixed constant. In general, the asymptotic complexity of this algorithm is exponential in the width of the history, but this worst-case behavior is not exercised in practice as shown in Section 3.7. Then, we prove that checking prefix consistency and snapshot isolation can be reduced in polynomial time to the problem of checking serializability.

3.5.1 CHECKING SERIALIZABILITY

We present an algorithm for checking serializability of a given history which constructs a valid commit order (satisfying Serialization), if any, by “linearizing” transactions one by one in an order consistent with the session order. At any time, the set of already linearized transactions is uniquely determined by an antichain of the session order (i.e., a set of mutually-unordered transactions w.r.t. so), and the next transaction to linearize is chosen among the immediate so successors of the transactions in this antichain. The crux of the algorithm is that the next transaction to linearize can be chosen such that it does not produce violations of Serialization in a way that does not depend on the order between the already linearized transactions. Therefore, the algorithm can be seen as a search in the space of so antichains. If the width of the history is bounded (by a fixed constant), then the number of possible so antichains is polynomial in the size of the history, which implies that the search can be done in polynomial time.

A *prefix* of a history $h = \langle T, \text{so}, \text{wr} \rangle$ is a set of transactions $T' \subseteq T$ such that all the so predecessors of transactions in T' are also in T' , i.e., $\forall t \in T. \text{so}^{-1}(t) \in T'$. A prefix T' is uniquely determined by the set of transactions in T' that are maximal w.r.t. so. This set of transactions forms an *antichain* of so, i.e., any two elements in this set are incomparable w.r.t. so. Given an antichain $\{t_1, \dots, t_n\}$ of so, we say that $\{t_1, \dots, t_n\}$ is the *boundary* of the prefix $T' = \{t : \exists i. \langle t, t_i \rangle \in \text{so} \vee t = t_i\}$. For instance, given the history in Figure 3.6a, the set of transactions $\{t_0, t_1, t_2\}$ is a prefix with boundary $\{t_1, t_2\}$ (the latter is an antichain of the session order).

A prefix T' of a history h is called *serializable* iff there exists a *partial* commit order co on the transactions in h such that the following hold:

- co does not contradict the session order and the write-read relation in h , i.e., $\text{wr} \cup \text{so} \cup \text{co}$ is acyclic,
- co is a total order on transactions in T' ,
- co orders transactions in T' before transactions in $T \setminus T'$, i.e., $\langle t_1, t_2 \rangle \in \text{co}$ for every $t_1 \in T'$ and $t_2 \in T \setminus T'$,
- co does not order any two transactions $t_1, t_2 \notin T'$
- the history h along with the commit order co satisfies the axiom defining serializability, i.e., $\langle h, \text{co} \rangle \models \text{Serialization}$.

For the history in Figure 3.6a, the prefix $\{t_0, t_1, t_2\}$ is serializable since there exists a partial commit order co that orders t_0, t_1, t_2 in this order, and both t_1 and t_2 before t_3 and t_4 . The axiom *Serialization* is satisfied trivially, since the prefix contains a single transaction writing x and all the transactions outside of the prefix do not read x .

A prefix $T' \uplus \{t\}$ of h is called a *valid extension*⁸ of a serializable prefix T' of h , denoted by $T' \triangleright T' \uplus \{t\}$ if:

- t does not read from a transaction outside of T' , i.e., for every $t' \in T \setminus T'$, $\langle t', t \rangle \notin \text{wr}$, and
- for every variable x written by t , there exists no transaction $t_2 \neq t$ outside of T' that reads a value of x written by a transaction t_1 in T' , i.e., for every x written by t and every $t_1 \in T'$ and $t_2 \in T \setminus (T' \uplus \{t\})$, $\langle t_1, t_2 \rangle \notin \text{wr}$.

For the history in Figure 3.6a, we have $\{t_0, t_1\} \triangleright \{t_0, t_1\} \uplus \{t_2\}$ because t_2 reads from t_0 and it does not write any variable. On the other hand $\{t_0, t_1\} \not\triangleright \{t_0, t_1\} \uplus \{t_3\}$ because t_3 writes x and the transaction t_2 , outside of this prefix, reads from the transaction t_0 included in the prefix.

Let \triangleright^* denote the reflexive and transitive closure of \triangleright .

The following lemma is essential in proving that iterative valid extensions of the initial empty prefix can be used to show that a given history is serializable.

Lemma 3.5.1. *For a serializable prefix T' of a history h , a prefix $T' \uplus \{t\}$ is serializable if it is a valid extension of T' .*

⁸We assume that $t \notin T'$ which is implied by the use of the disjoint union \uplus .

Input: A history $h = (T, \text{so}, \text{wr})$, a serializable prefix T' of h
Output: *true* iff $T' \triangleright^* h$

```

1 if  $T' = T$  then
2   return true;
3 foreach  $t \notin T'$  s.t.  $\forall t' \notin T'. \langle t', t \rangle \notin \text{wr} \cup \text{so}$  do
4   if  $T' \not\triangleright T' \uplus \{t\}$  then
5     continue;
6   if  $T' \uplus \{t\} \notin \text{seen} \wedge \text{checkSER}(h, T' \uplus \{t\})$  then
7     return true;
8    $\text{seen} \leftarrow \text{seen} \cup \{T' \uplus \{t\}\}$ ;
9 return false;

```

Algorithm 5: The algorithm `checkSER` for checking serializability. *seen* is a global variable storing a set of prefixes of h (which are not serializable). It is initialized as the empty set.

Proof. Let co' be the partial commit order for T' which satisfies the serializable prefix conditions. We extend co' to a partial order $\text{co} = \text{co}' \cup \{\langle t, t' \rangle \mid t' \notin T' \uplus \{t\}\}$. We show that $\langle h, \text{co} \rangle \models \text{Serialization}$. The other conditions for $T' \uplus \{t\}$ being a serializable prefix are satisfied trivially by co .

Assume by contradiction that $\langle h, \text{co} \rangle$ does not satisfy the axiom `Serialization`. Then, there exists $t_1, t_2, t_3, x \in \text{vars}(h)$ s.t. $\langle t_1, t_3 \rangle \in \text{wr}_x$ and t_2 writes on x and $\langle t_1, t_2 \rangle, \langle t_2, t_3 \rangle \in \text{co}$. Since $\langle h, \text{co}' \rangle$ satisfies this axiom, at least one of these two co ordering constraints are of the form $\langle t, t' \rangle$ where $t' \notin T' \uplus \{t\}$:

- the case $t_1 = t$ and $t_2 \notin T' \uplus \{t\}$ is not possible because co' contains no pair of the form $\langle t', _ \rangle \in \text{co}'$ with $t' \notin T'$ (recall that $\langle t_2, t_3 \rangle$ should be also included in co).
- If $t_2 = t$ then, $\langle t_1, t_2 \rangle \in \text{co}'$ and $\langle t_2, t_3 \rangle$ for some $t_3 \notin T' \uplus \{t\}$. But, by the definition of valid extension, for all variables x written by t , there exists no transaction $t_3 \notin T' \uplus \{t\}$ such that it reads x from $t_1 \in T'$. Therefore, this is also a contradiction. \square

Algorithm 5 lists our algorithm for checking serializability. It is defined as a recursive procedure that searches for a sequence of valid extensions of a given prefix (initially, this prefix is empty) until covering the whole history. Figure 3.6b pictures this search on the history in Figure 3.6a. The right branch (containing blue edges) contains only valid extensions and it reaches a prefix that includes all the transactions in the history.

Theorem 3.5.1. *A history h is serializable iff $\text{checkSER}(h, \emptyset)$ returns true.*

Proof. The “if” direction is a direct consequence of Lemma 3.5.1. For the reverse, assume that $h = \langle T, \text{so}, \text{wr} \rangle$ is serializable with a (total) commit order co . Let co_i be the set of transactions in the prefix of co of length i . Since co is consistent with so , we have that co_i is a prefix of h , for any i . We show by induction that co_{i+1} is a valid extension of co_i . The base case is trivial. For the induction step, let t be the last transaction in the prefix of co of length $i + 1$. Then,

- t cannot read from a transaction outside of co_i because co is consistent with the write-read relation wr ,

- also, for every variable x written by t , there exists no transaction $t_2 \neq t$ outside of co_i that reads a value of x written by a transaction $t_1 \in \text{co}_i$. Otherwise, $\langle t_1, t_2 \rangle \in \text{wr}_x$, $\langle t, t_2 \rangle \in \text{co}$, and $\langle t_1, t \rangle \in \text{co}$ which implies that $\langle h, \text{co} \rangle$ does not satisfy Serializability.

This implies that $\text{checkSER}(h, \emptyset)$ returns true. \square

Algorithm 5 enumerates prefixes of the given history h , each prefix being uniquely determined by an antichain of h containing the so-maximal transactions in that prefix. By definition, the size of each antichain of a history h is smaller than the width of h . Therefore, the number of possible antichains (prefixes) of a history h is $O(\text{size}(h)^{\text{width}(h)})$ where $\text{size}(h)$, resp., $\text{width}(h)$, is the number of transactions, resp., the width, of h . Since the valid extension property can be checked in quadratic time, the asymptotic time complexity of the algorithm defined by checkSER is upper bounded by $O(\text{size}(h)^{\text{width}(h)} \cdot \text{size}(h)^3)$. The following corollary is a direct consequence of these observations.

Corollary 3.5.1. *For an arbitrary but fixed constant $k \in \mathbb{N}$, the problem of checking serializability for histories of width at most k is polynomial time.*

3.5.2 REDUCING PREFIX CONSISTENCY TO SERIALIZABILITY

We describe a polynomial time reduction of checking prefix consistency of bounded-width histories to the analogous problem for serializability. Intuitively, as opposed to serializability, prefix consistency allows that two transactions read the same snapshot of the database and commit together even if they write on the same variable. Based on this observation, given a history h for which we want to check prefix consistency, we define a new history $h_{R|W}$ where each transaction t is split into a transaction performing all the reads in t and another transaction performing all the writes in t (the history $h_{R|W}$ retains all the session order and write-read dependencies of h). We show that if the set of read and write transactions obtained this way can be shown to be serializable, then the original history satisfies prefix consistency, and vice-versa. For instance, Figure 3.7 shows this transformation on the two histories in Figure 3.7a and Figure 3.7c, which represent typical anomalies known as “long fork” and “lost update”, respectively. The former is not admitted by PC while the latter is admitted. It can be easily seen that the transformed history corresponding to the “long fork” anomaly is not serializable while the one corresponding to “lost update” is serializable. We show that this transformation leads to a history of the same width, which by Corollary 3.5.1, implies that checking prefix consistency of bounded-width histories is polynomial time.

Thus, given a history $h = \langle T, \text{wr}, \text{so} \rangle$, we define the history $h_{R|W} = \langle T', \text{wr}', \text{so}' \rangle$ as follows:

- T' contains a transaction R_t , called a *read* transaction, and a transaction W_t , called a *write* transaction, for each transaction t in the original history, i.e., $T' = \{R_t | t \in T\} \cup \{W_t | t \in T\}$
- the write transaction W_t writes exactly the same set of variables as t , i.e., for each variable x , W_t writes to x iff t writes to x .
- the read transaction R_t reads exactly the same values and the same variables as t , i.e., for each variable x , $\text{wr}'_x = \{\langle W_{t_1}, R_{t_2} \rangle | \langle t_1, t_2 \rangle \in \text{wr}_x\}$

3 Transactional Consistency

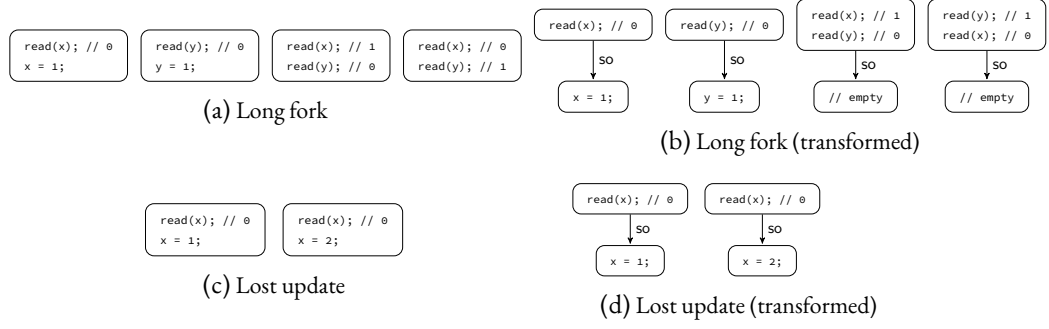


Figure 3.7: Reducing PC to SER. Initially, the value of every variable is 0.

- the session order between the read and the write transactions corresponds to that of the original transactions and read transactions precede their write counterparts, i.e.,

$$so' = \{\langle R_t, W_t \rangle | t \in T\} \cup \{\langle R_{t_1}, R_{t_2} \rangle, \langle R_{t_1}, W_{t_2} \rangle, \langle W_{t_1}, R_{t_2} \rangle, \langle W_{t_1}, W_{t_2} \rangle | \langle t_1, t_2 \rangle \in so\}$$

The following lemma is a straightforward consequence of the definitions.

Lemma 3.5.2. *The histories h and $h_{R|W}$ have the same width.*

Next, we show that $h_{R|W}$ is serializable if h is prefix consistent. Formally, we show that

$$\forall co. \exists co'. \langle h, co \rangle \models \text{Prefix} \Rightarrow \langle h_{R|W}, co' \rangle \models \text{Serializability}$$

Thus, let co be a commit (total) order on transactions of h which together with h satisfies the prefix consistency axiom. We define two *partial* commit orders co'_1 and co'_2 , co'_2 a strengthening of co'_1 , which we prove that they are acyclic and that any linearization co' of co'_2 is a valid witness for $h_{R|W}$ satisfying serializability.

Thus, let co'_1 be a *partial* commit order on transactions of $h_{R|W}$ defined as follows:

$$co'_1 = \{\langle R_t, W_t \rangle | t \in T\} \cup \{\langle W_{t_1}, W_{t_2} \rangle | \langle t_1, t_2 \rangle \in co\} \cup \{\langle W_{t_1}, R_{t_2} \rangle | \langle t_1, t_2 \rangle \in \text{wr} \cup so\}$$

We show that if co'_1 were to be cyclic, then it contains a minimal cycle with one read transaction, and at least one but at most two write transactions. Then, we show that such cycles cannot exist.

Lemma 3.5.3. *The relation co'_1 is acyclic.*

PROOF. We first show that if co'_1 were to be cyclic, then it contains a minimal cycle with one read transaction, and at least one but at most two write transactions. Then, we show that such cycles cannot exist. Therefore, let us assume that co'_1 is cyclic. Then,

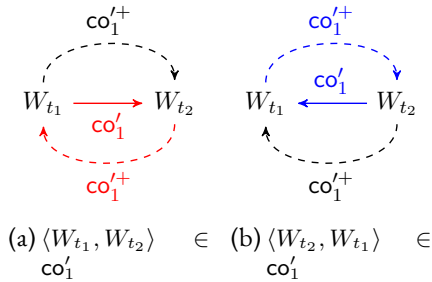


Figure 3.8: Cycles with non-consecutive write transactions.

- Since $\langle W_{t_1}, W_{t_2} \rangle \in co'_1$ implies $\langle t_1, t_2 \rangle \in co$, for every t_1 and t_2 , a cycle in co'_1 cannot contain only write transactions. Otherwise, it will imply a cycle in the original commit order co . Therefore, a cycle in co'_1 must contain at least one read transaction.

- Assume that a cycle in co'_1 contains two write transactions W_{t_1} and W_{t_2} which are not consecutive, like in Figure 3.8. Since either $\langle W_{t_1}, W_{t_2} \rangle \in co'_1$ or $\langle W_{t_2}, W_{t_1} \rangle \in co'_1$, there exists a smaller cycle in co'_1 where these two write transactions are consecutive. If $\langle W_{t_1}, W_{t_2} \rangle \in co'_1$, then co'_1 contains the smaller cycle on the lower part of the original cycle (Figure 3.8a), and if $\langle W_{t_2}, W_{t_1} \rangle \in co'_1$, then co'_1 contains the cycle on the upper part of the original cycle (Figure 3.8b). Thus, all the write transactions in a minimal cycle of co'_1 must be consecutive.
- If a minimal cycle were to contain three write transactions, then all of them cannot be consecutive unless they all three form a cycle, which is not possible. So a minimal cycle contains at most two write transactions.
- Since co'_1 contains no direct relation between read transactions, it cannot contain a cycle with two consecutive read transactions, or only read transactions.

This shows that a minimal cycle of co'_1 would include a read transaction and a write transaction, and at most one more write transaction. We prove that such cycles are however impossible:

- if the cycle is of size 2, then it contains two transactions W_{t_1} and R_{t_2} such that $\langle W_{t_1}, R_{t_2} \rangle \in co'_1$ and $\langle R_{t_2}, W_{t_1} \rangle \in co'_1$. Since all the $\langle R_-, W_- \rangle$ dependencies in co'_1 are of the form $\langle R_t, W_t \rangle$, it follows that $t_1 = t_2$. Then, we have $\langle W_{t_1}, R_{t_1} \rangle \in co'_1$ which implies $\langle t_1, t_1 \rangle \in \text{wr} \cup \text{so}$, a contradiction.
- if the cycle is of size 3, then it contains three transactions W_{t_1} , W_{t_2} , and R_{t_3} such that $\langle W_{t_1}, W_{t_2} \rangle \in co'_1$, $\langle W_{t_2}, R_{t_3} \rangle \in co'_1$, and $\langle R_{t_3}, W_{t_1} \rangle \in co'_1$. Using a similar argument as in the previous case, $\langle R_{t_3}, W_{t_1} \rangle \in co'_1$ implies $t_3 = t_1$. Therefore, $\langle t_1, t_2 \rangle \in \text{co}$ and $\langle t_2, t_1 \rangle \in \text{wr} \cup \text{so}$, which contradicts the fact that $\text{wr} \cup \text{so} \subseteq \text{co}$. \square

We define a strengthening of co'_1 where intuitively, we add all the dependencies from read transactions t_3 to write transactions t_2 that “overwrite” values read by t_3 . Formally, $co'_2 = co'_1 \cup \text{RW}(co'_1)$ where

$$\text{RW}(co'_1) = \{ \langle t_3, t_2 \rangle \mid \exists x \in \text{vars}(h). \exists t_1 \in T'. \langle t_1, t_3 \rangle \in \text{wr}_x', \langle t_1, t_2 \rangle \in co'_1, t_2 \text{ writes } x \}$$

It can be shown that any cycle in co'_2 would correspond to a Prefix violation in the original history. Therefore,

Lemma 3.5.4. *The relation co'_2 is acyclic.*

Proof. Assume that co'_2 is cyclic. Any minimal cycle in co'_2 still satisfies the properties of minimal cycles of co'_1 proved in Lemma 3.5.3 (because all write transactions are still totally ordered and co'_2 doesn't relate directly read transactions). So, a minimal cycle in co'_2 contains a read transaction and a write transaction, and at most one more write transaction.

Since co'_1 is acyclic, a cycle in co'_2 , and in particular a minimal one, must necessarily contain a dependency from $\text{RW}(co'_1)$. Note that a minimal cycle cannot contain two such dependencies since this would imply that it contains two non-consecutive write transactions. The red edges in

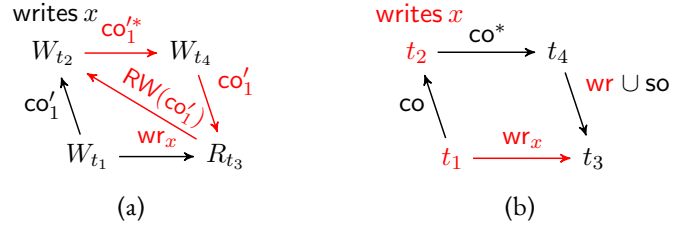


Figure 3.9: Cycles in co'_2 correspond to Prefix violations: (a) Minimal cycle in co'_2 , (b) Prefix violation in $\langle h, co \rangle$.

Figure 3.9a show a minimal cycle of co'_2 satisfying all the properties mentioned above. This cycle contains a dependency $\langle R_{t3}, W_{t2} \rangle \in RW(co'_1)$ which implies the existence of a write transaction W_{t1} in $h_{R|W}$ s.t. $\langle W_{t1}, R_{t3} \rangle \in wr_x'$ and $\langle W_{t1}, W_{t2} \rangle \in co'_1$ and W_{t1}, W_{t2} write on x (these dependencies are represented by the black edges in Figure 3.9a). The relations between these transactions of $h_{R|W}$ imply that the corresponding transactions of h are related as shown in Figure 3.9b: $\langle W_{t1}, W_{t2} \rangle \in co'_1$ and $\langle W_{t2}, W_{t4} \rangle \in co'_1$ imply $\langle t1, t2 \rangle \in co$ and $\langle t2, t4 \rangle \in co^*$, respectively, $\langle W_{t1}, W_{t3} \rangle \in wr_x'$ implies $\langle t1, t3 \rangle \in wr_x$, and $\langle W_{t4}, R_{t3} \rangle \in co'_1$ implies $\langle t4, t3 \rangle \in wr \cup so$. This implies that $\langle h, co \rangle$ doesn't satisfy the Prefix axiom, a contradiction. \square

Lemma 3.5.5. *If a history h satisfies prefix consistency, then $h_{R|W}$ is serializable.*

Proof. Let co' be any total order consistent with co'_2 . Assume by contradiction that $\langle h_{R|W}, co' \rangle$ doesn't satisfy Serializability. Then, there exist $t'_1, t'_2, t'_3 \in T'$ such that $\langle t'_1, t'_2 \rangle, \langle t'_2, t'_3 \rangle \in co'$ and t'_1, t'_2 write on some variable x and $\langle t'_1, t'_3 \rangle \in wr_x'$. But then t'_1, t'_2 are write transactions and co'_1 must contain $\langle t'_1, t'_2 \rangle$. Therefore, $RW(co'_1)$ and co'_2 should contain $\langle t'_3, t'_2 \rangle$, a contradiction with co' being consistent with co'_2 . \square

Finally, it can be proved that any linearization co' of co'_2 satisfies Serializability (together with $h_{R|W}$). Moreover, it can also be shown that the serializability of $h_{R|W}$ implies that h satisfies PC. Therefore,

Theorem 3.5.2. *A history h satisfies prefix consistency iff $h_{R|W}$ is serializable.*

PROOF. The “only-if” direction is proven by Lemma 3.5.5. For the reverse, we show that

$$\forall co'. \exists co. \langle h_{R|W}, co' \rangle \models \text{Serializability} \Rightarrow \langle h, co \rangle \models \text{Prefix}$$

Thus, let co' be a commit (total) order on transactions of $h_{R|W}$ which together with $h_{R|W}$ satisfies the serializability axiom. Let co be a commit order on transactions of h defined by $co = \{ \langle t1, t2 \rangle \mid \langle W_{t1}, W_{t2} \rangle \in co' \}$ (co is clearly a total order). If co were not to be consistent with $wr \cup so$, then there would exist transactions $t1$ and $t2$ such that $\langle t1, t2 \rangle \in wr \cup so$ and $\langle t2, t1 \rangle \in co$, which would imply that $\langle W_{t1}, R_{t2} \rangle, \langle R_{t2}, W_{t2} \rangle \in wr \cup so$ and $\langle W_{t2}, W_{t1} \rangle \in co'$, which violates the acyclicity of co' . We show that $\langle h, co \rangle$ satisfies Prefix. Assume by contradiction that there exists a Prefix violation between $t1, t2, t3, t4$ (shown in Figure 3.10a), i.e., for some $x \in \text{vars}(h)$, $\langle t1, t3 \rangle \in wr_x$ and $t2$ writes x , $\langle t1, t2 \rangle \in co$, $\langle t2, t4 \rangle \in co^*$ and $\langle t4, t3 \rangle \in wr \cup so$. Then, the corresponding transactions $W_{t1}, W_{t2}, W_{t4}, R_{t3}$ in $h_{R|W}$ would be related as follows:

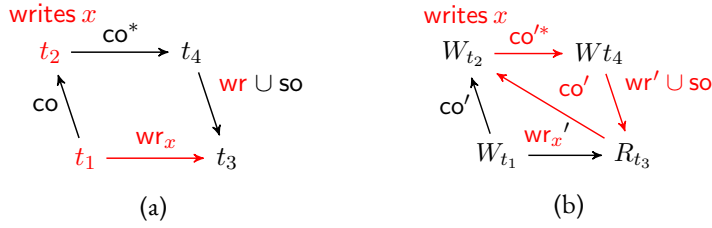


Figure 3.10: Prefix violations correspond to cycles in co' : (a) Prefix violation in (h, co) , (b) Cycle in co' .

$\langle W_{t_1}, W_{t_2} \rangle \in co'$ and $\langle W_{t_1}, R_{t_3} \rangle \in wr_x'$ because $\langle t_1, t_3 \rangle \in wr_x$ and $\langle t_1, t_2 \rangle \in co$. Since co' satisfies Serializability, then $\langle R_{t_3}, W_{t_2} \rangle \in co'$. But $\langle t_2, t_4 \rangle \in co^*$ and $\langle t_4, t_3 \rangle \in wr \cup so$ imply that $\langle W_{t_2}, W_{t_4} \rangle \in co'^*$ and $\langle W_{t_4}, R_{t_3} \rangle \in wr' \cup so'$, which show that co' is cyclic (the red cycle in Figure 3.10b), a contradiction. \square

Since the history $h_{R|W}$ can be constructed in linear time, Lemma 3.5.2, Theorem 3.5.2, and Corollary 3.5.1 imply the following result.

Corollary 3.5.2. *For an arbitrary but fixed constant $k \in \mathbb{N}$, the problem of checking prefix consistency for histories of width at most k is polynomial time.*

3.5.3 REDUCING SNAPSHOT ISOLATION TO SERIALIZABILITY

We extend the reduction of prefix consistency to serializability to the case of snapshot isolation. Compared to prefix consistency, snapshot isolation disallows transactions that read the same snapshot of the database to commit together if they write on a common variable (stated by the Conflict axiom). More precisely, for any pair of transactions t_1 and t_2 writing to a common variable, t_1 must observe the effects of t_2 or vice-versa. We refine the definition of $h_{R|W}$ such that any “serialization” (i.e., commit order satisfying Serializability) disallows that the read transactions corresponding to two such transactions are ordered both before their write counterparts. We do this by introducing auxiliary variables that are read or written by these transactions. For instance, Figure 3.11 shows this transformation on the two histories in Figure 3.11a and Figure 3.11c, which represent the anomalies known as “lost update” and “write skew”, respectively. The former is not admitted by SI while the latter is admitted. Concerning “lost update”, the read counterpart of the transaction on the left writes to a variable x_{12} that is read by its write counterpart, but also written by the write counterpart of the other transaction. This forbids that the latter is serialized in between the read and write counterparts of the transaction on the left. A similar scenario is imposed on the transaction on the right, which makes that the transformed history is not serializable. Concerning the “write skew” anomaly, the transformed history is exactly as for the PC reduction since the two transactions don’t write on a common variable. It is clearly serializable.

For a history $h = \langle T, wr, so \rangle$, the history $h_{R|W}^c = \langle T', wr', so' \rangle$ is defined as $h_{R|W}$ with the following additional construction: for every two transactions t_1 and $t_2 \in T$ that write on a common variable,

- R_{t_1} and W_{t_2} (resp., R_{t_2} and W_{t_1}) write on a variable $x_{1,2}$ (resp., $x_{2,1}$),

3 Transactional Consistency

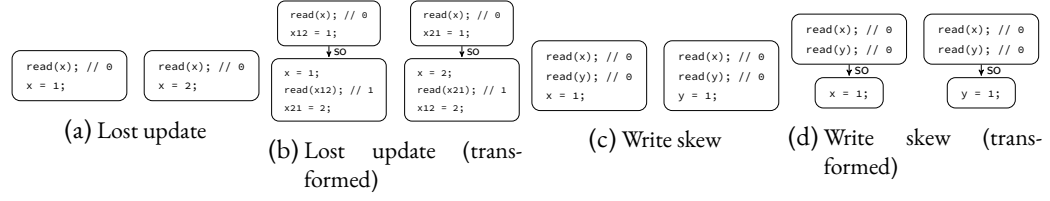


Figure 3.11: Reducing SI to SER.

- the write transaction of t_i reads $x_{i,j}$ from the read transaction of t_i , for all $i \neq j \in \{1, 2\}$, i.e., $\text{wr}_{x_{1,2}} = \{\langle R_{t_1}, W_{t_1} \rangle\}$ and $\text{wr}_{x_{2,1}} = \{\langle R_{t_2}, W_{t_2} \rangle\}$.

Note that $h_{R|W}$ and $h_{R|W}^c$ have the same width (the session order is defined exactly in the same way), which implies, by Lemma 3.5.2, that h and $h_{R|W}^c$ have the same width.

The following result can be proved using similar reasoning as in the case of prefix consistency.

Theorem 3.5.3. *A history h satisfies snapshot isolation iff $h_{R|W}^c$ is serializable.*

Note that $h_{R|W}^c$ and h have the same width, and that $h_{R|W}^c$ can be constructed in linear time. Therefore, Theorem 3.5.3, and Corollary 3.5.1 imply the following result.

Corollary 3.5.3. *For an arbitrary but fixed constant $k \in \mathbb{N}$, the problem of checking snapshot isolation for histories of width at most k is polynomial time.*

3.6 COMMUNICATION GRAPHS

In this section, we present an extension of the polynomial time results for PC, SI, and SER, which allows to handle histories where the sharing of variables between different sessions is *sparse*. For the results in this section, we take the simplifying assumption that the session order is a union of transaction sequences (modulo the fictitious transaction writing the initial values), i.e., each transaction sequence corresponding to the standard notion of *session*⁹. We represent the sharing of variables between different sessions using an undirected graph called a *communication graph*. For instance, the communication graph of the history in Figure 3.12a is given in Figure 3.12b. For readability, the edges are marked with the variables accessed by the two sessions.

We show that the problem of checking PC, SI, or SER is polynomial time when the size of every *biconnected* component of the communication graph is bounded by a fixed constant. This is stronger than the results in Section 3.5 because the number of biconnected components can be arbitrarily large which means that the total number of sessions is unbounded. In general, we prove that the time complexity of these consistency criteria is exponential only in the maximum size of such a biconnected component, and not the whole number of sessions.

An undirected graph is biconnected if it is connected and if any one vertex were to be removed, the graph will remain connected, and a biconnected component of a graph G is a maximal biconnected subgraph of G . Figure 3.12b shows the decomposition in biconnected compo-

⁹The results can be extended to arbitrary session orders by considering maximal transaction sequences in session order instead of sessions.

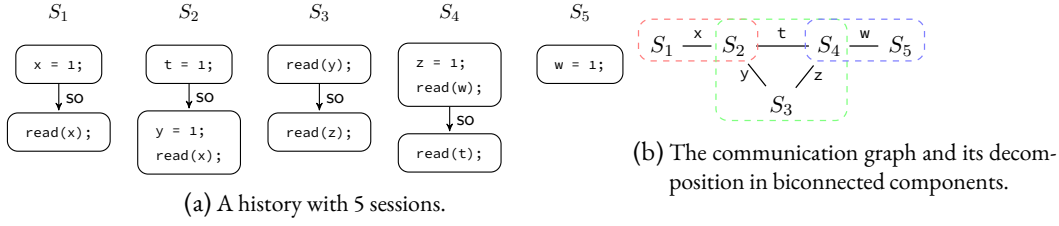


Figure 3.12: A history and its communication graph.

nents of a communication graph. This graph contains 5 sessions while every biconnected component is of size at most 3. Intuitively, if a history h is a violation to some consistency criterion $C \in \{\text{PC}, \text{SI}, \text{SER}\}$, then there exists a projection of h on sessions from the *same* biconnected component which is also a violation to C (the reverse is trivially true). Therefore, checking any of these criteria can be done in isolation for each biconnected component (more precisely, on sub-histories that contain only sessions in the same biconnected component). Actually, this decomposition argument works even for RC, RA, and CC. For instance, in the case of the history in Figure 3.12a, any consistency criterion can be checked looking in isolation at three sub-histories: a sub-history with S_1 and S_2 , a sub-history with S_2 , S_3 , and S_4 , and a sub-history with S_4 and S_5 .

Formally, a *communication graph* of a history h is an undirected graph $\text{Comm}(h) = (V, E)$ where the set of vertices V is the set of sessions¹⁰ in h , and $(v, v') \in E$ iff the sessions v and v' contain two transactions t_1 and t_2 , respectively, such that t_1 and t_2 read or write a common variable x .

We begin with a technical lemma showing that *minimal* paths of certain form in the graph representing a history h and a relation co (on the transactions of h) lie within a single biconnected component of the underlying communication graph. This is used to show that any consistency violation can be exposed by looking at a single biconnected component at a time. The graph representing a history h and a relation co on the transactions of h is denoted by $G(h, \text{co})$ ¹¹.

Given a graph $G(h, \text{co})$ and a relation r on its vertices, a term over the relations so , wr , and co , e.g., $(\text{wr} \cup \text{so})^+$, a path of the form r (or an r -path) is a sequence of edges representing so , wr , or co dependencies as specified by the term r , e.g., a sequence of wr or so dependencies.

Lemma 3.6.1. *Let B_1, \dots, B_n be the biconnected components of $\text{Comm}(h)$ for a history $h = \langle T, \text{wr}, \text{so} \rangle$. For each B_i , let co_i be a total order on the transactions of B_i ¹² extending the session order so on the transactions of B_i . Also, let $\text{co} = \bigcup_i \text{co}_i$. Then, for every term $r \in \{\text{co}^+, (\text{wr} \cup \text{so})^+\}$, any minimal r -path in the graph $G(h, \text{co})$ between two transactions from the same biconnected component includes only transactions of that biconnected component.*

PROOF. We consider the case $r = \text{co}^+$. Consider a *minimal* co^+ -path $\pi = t_0, \dots, t_n$ between two transactions t_0 and t_n from the same biconnected component B of $\text{Comm}(h)$ (i.e., from sessions in B). Assume by contradiction, that π traverses multiple biconnected components.

¹⁰The transaction writing the initial values is considered as a distinguished session.

¹¹The nodes of $G(h, \text{co})$ correspond to transactions in h and the edges connect pairs of transactions in so , wr , or co .

¹²That is, transactions that are included in the sessions in B_i .

We define a path $\pi_s = v_0, \dots, v_m$ between sessions, i.e., vertices of $\text{Comm}(h)$, which contains an edge (v_j, v_{j+1}) iff π contains an edge (t_i, t_{i+1}) with t_i a transaction of session v_j and t_{i+1} a transaction of session $v_{j+1} \neq v_j$. Since any graph decomposes to a forest of biconnected components, this path must necessarily leave and enter some biconnected component B_1 to and from the same biconnected component B_2 , i.e., π_s must contain two vertices v_{j_1} and v_{j_2} in B_1 such that the successor v_{j_1+1} of v_{j_1} and the predecessor v_{j_2-1} of v_{j_2} are from B_2 . Let t_1, t_2, t_3, t_4 be the transactions in the path π corresponding to $v_{j_1}, v_{j_2}, v_{j_1+1}$, and v_{j_2-1} , respectively. Now, since any two biconnected components share at most one vertex, it follows that t_3 and t_4 are from the same session and

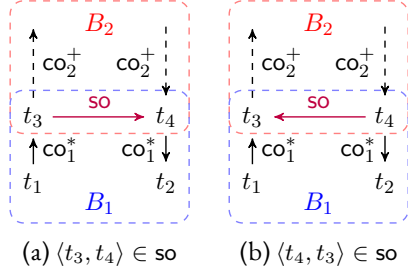


Figure 3.13: Minimal paths between transactions in the same biconnected component.

- if $\langle t_3, t_4 \rangle \in \text{so}$, then there exists a shorter path between t_0 and t_1 that uses the so relation between $\langle t_3, t_4 \rangle$ (we recall that $\text{so} \subseteq \bigcup_i \text{co}_i$) instead of the transactions in B_2 , pictured in Figure 3.13a, which is a contradiction to the minimality of π ,
- if $\langle t_4, t_3 \rangle \in \text{so}$, then, we have a cycle in $\bigcup_i \text{co}_i \cup \text{so}$, pictured in Figure 3.13b, which is also a contradiction.

The case $r = (\text{wr} \cup \text{so})^+$ can be proved in a similar manner since the reasoning outlined in Figure 3.13 reduces to short-circuiting a path using a single so edge

(and so is included in $(\text{wr} \cup \text{so})^+$). \square

Now we prove our final claim. For a history $h = (T, \text{so}, \text{wr})$ and biconnected component B of $\text{Comm}(h)$, the projection of h over transactions in sessions of B is denoted by $h \downarrow B$, i.e., $h \downarrow B = (T', \text{so}', \text{wr}')$ where T' is the set of transactions in sessions of B , so' and wr' are the projections of so and wr , respectively, on T' .

Theorem 3.6.1. *For any criterion $C \in \{\text{RA}, \text{RC}, \text{CC}, \text{PC}, \text{SI}, \text{SER}\}$, a history h satisfies C iff for every biconnected component B of $\text{Comm}(h)$, $h \downarrow B$ satisfies C .*

Proof. The “only-if” direction is obvious. For the “if” direction, we first consider the cases $C \in \{\text{RA}, \text{RC}, \text{CC}, \text{SER}\}$. The proof concerning PC and SI is based on the reduction to SER outlined in Section 3.5.2 and Section 3.5.3, respectively, and it is given afterwards. Let B_1, \dots, B_n be the biconnected components of $\text{Comm}(h)$.

Let $C \in \{\text{RA}, \text{RC}, \text{CC}, \text{SER}\}$, and let co_i be the commit order that witnesses that $h \downarrow B_i$ satisfies C , for each $1 \leq i \leq n$. The union $\bigcup_i \text{co}_i$ is acyclic since otherwise, any minimal cycle would be a minimal path between transactions of the same biconnected component B_j , and, by Lemma 3.6.1, it will include only transactions of B_j which is a contradiction to co_j being a total order. We show that any linearization co of $\bigcup_i \text{co}_i$ along with h satisfies the axioms of C . The axioms defining RA, RC, CC, and SER involve transactions that write or read a common variable, which implies that they belong to the same biconnected component (we refer to the transactions

t_1 , t_2 , and t_3 in Figure 3.2). Furthermore, by Lemma 3.6.1, minimal paths witnessing the dependencies in those axioms, e.g., $(\text{wr} \cup \text{so})^+$ for CC, are also formed of transactions included in the same biconnected component. Therefore, co satisfies any of those axioms provided that each co_i does.

We now consider the case where $C = \text{PC}$. Assume that each B_i satisfies PC. Based on the reduction in Section 3.5.2, h satisfies PC iff $h_{R|W}$ satisfies SER. Moreover, since $h_{R|W}$ is obtained from h by splitting each transaction t into a read transaction R_t and a write transaction W_t while keeping all session order dependencies, each session in h corresponds to a session in $h_{R|W}$ that reads or writes exactly the same set of variables. Therefore, $\text{Comm}(h)$ is isomorphic to $\text{Comm}(h_{R|W})$. Since B_i satisfies PC, we get that the corresponding biconnected component B'_i of $\text{Comm}(h_{R|W})$ satisfies SER, for every i . Therefore, $h_{R|W}$ satisfies SER, which implies that h satisfies PC. The case of SI is proved in a similar way using the reduction to the serializability of $h_{R|W}^c$ presented in Section 3.5.3 (note that two transactions of $h_{R|W}^c$ may read or write an additional common variable only if they were writing a common variable in the original history and therefore, $\text{Comm}(h)$ is still isomorphic to $\text{Comm}(h_{R|W}^c)$). \square

Since the decomposition of a graph into biconnected components can be done in linear time, Theorem 3.6.1 implies that any of the criteria PC, SI, or SER can be checked in time $O(\text{size}(h)^{\text{bi-size}(h)} \cdot \text{bi-nb}(h))$ where $\text{bi-size}(h)$ and $\text{bi-nb}(h)$ are the maximum size of a biconnected component in $\text{Comm}(h)$ and the number of biconnected components of $\text{Comm}(h)$, respectively. The following corollary is a direct consequence of this observation.

Corollary 3.6.1. *For an arbitrary but fixed constant $k \in \mathbb{N}$ and any criterion $C \in \{\text{PC}, \text{SI}, \text{SER}\}$, the problem of checking if a history h satisfies C is polynomial time, provided that the size of every biconnected component of $\text{Comm}(h)$ is bounded by k .*

3.7 EXPERIMENTAL EVALUATION

To demonstrate the practical value of the theory developed in the previous sections, we argue that our algorithms:

- are efficient and scalable,
- enable an effective testing framework allowing to expose consistency violations in production databases.

We focus on three of the criteria introduced in Section 3.3: *serializability* which is NP-complete in general and polynomial time when the number of sessions is considered to be a constant, *snapshot isolation* which can be reduced in linear time to serializability, and *causal consistency* which is polynomial time in general. As benchmark, we consider histories extracted from three distributed databases: CockroachDB [cockroach], Galera [galera], and AntidoteDB [antidote]. Following the approach in Jepsen [jepson], histories are generated with random clients. For the experiments described hereafter, the randomization process is parametrized by: (1) the number of sessions (**#sess**), (2) the number of transactions per session (**#trs**), (3) the number of operations per transaction

3 Transactional Consistency

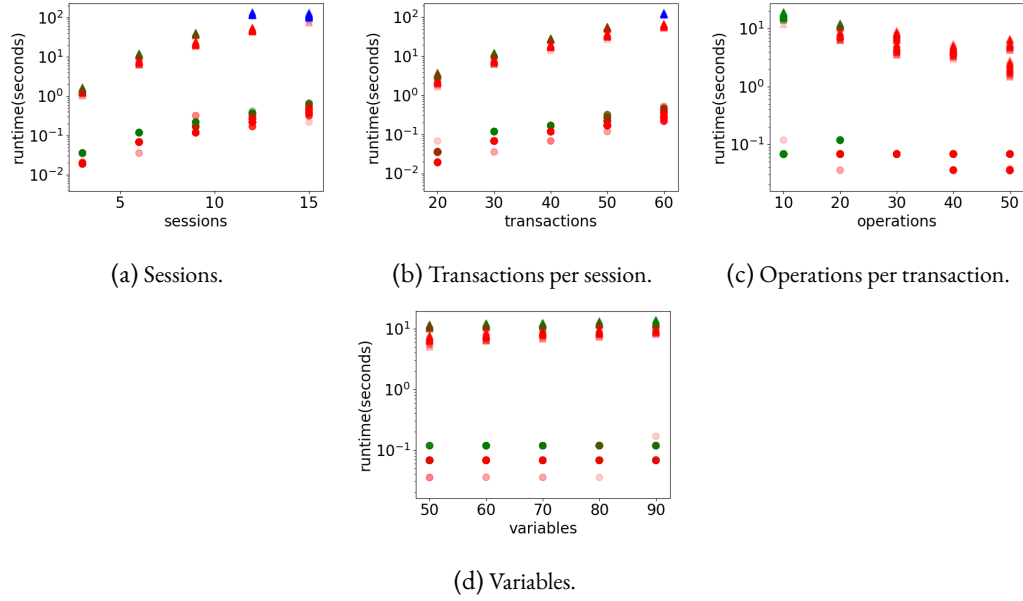


Figure 3.14: Scalability of our algorithm for checking SERIALIZABILITY (Algorithm 5) with comparison to a SAT encoding. The x-axis represents the varying parameter while the y-axis represents the wall-clock time in logarithmic scale. The circular, resp., triangular, dots represent wall-clock times of our algorithm, resp., the SAT encoding. The red, green, and blue dots represent invalid, valid and resource-exhausted instances, respectively.

(#ops), and (4) an upper bound on the number of used variables (#vars)¹³. For any valuation of these parameters, half of the histories generated with CockroachDB and Galera are restricted such that the sets of variables written by any two sessions are disjoint (the sets of read variables are not constrained). This restriction is used to increase the frequency of valid histories.

In a first experiment, we investigated the efficiency of our serializability-checking algorithm (Algorithm 5) and we compared its performance with a direct SAT encoding¹⁴ of the serializability definition in Section 3.3 (we used MiniSAT [DBLP:conf/sat/EenS03] to solve the SAT queries). We used histories extracted from CockroachDB which claims to implement serializability, acknowledging however the possibility of anomalies [cockroach-claim]. The sessions of a history are uniformly distributed among 3 nodes of a single cluster. To evaluate scalability, we fix a reference set of parameter values: #sess=6, #trs=30, #ops=20, and #vars = 60 × #sess, and vary only one parameter at a time. For instance, the number of sessions varies from 3 to 15 in increments of 3. We consider 100 histories for each combination of parameter values. The experimental data is reported in Figure 3.14. Our algorithm scales well even when increasing the number of sessions, which is not guaranteed by its worst-case complexity (in general, this is exponential in the number

¹³We ensure that every value is written at most once.

¹⁴For each ordered pair of transactions t_1, t_2 we add two propositional variables representing $\langle t_1, t_2 \rangle \in (\text{wr} \cup \text{so})^+$ and $\langle t_1, t_2 \rangle \in \text{co}$, respectively. Then we generate clauses corresponding to: (1) singleton clauses defining the relation $\text{wr} \cup \text{so}$ (extracted from the input history), (2) $\langle t_1, t_2 \rangle \in \text{wr} \cup \text{so}$ implies $\langle t_1, t_2 \rangle \in \text{co}$, (3) co being a total order, and (4) the axioms corresponding to the considered consistency model. This is an optimization that does not encode wr and so separately, which is sound because of the shape of our axioms (and because these relations are fixed apriori).

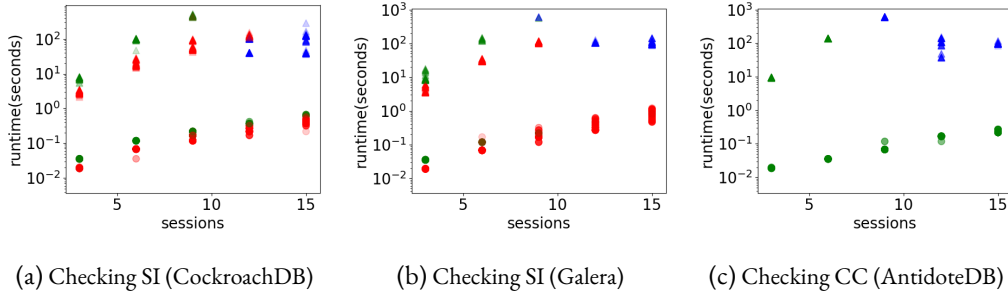


Figure 3.15: Scalability of our algorithms for checking SNAPSHOT ISOLATION (Section 3.5.3) and CAUSAL CONSISTENCY (Algorithm 4) with comparison to a SAT encoding. The x-axis represents the varying parameter while the y-axis represents the wall-clock time in logarithmic scale. The circular, resp., triangular, dots represent wall-clock times of our algorithm, resp., the SAT encoding. The red, green, and blue dots represent invalid, valid and resource-exhausted instances, respectively.

of sessions). Also, our algorithm is at least two orders of magnitude more efficient than the SAT encoding. While the performance of SAT solvers is known to be heavily affected by the specific encoding of the problem, we strove to make the SAT formula as succinct as possible and optimize its construction. We have fixed a 10 minutes timeout, a limit of 10GB of memory, and a limit of 10GB on the files containing the formulas to be passed to the SAT solver. The blue dots represent resource-exhausted instances. The SAT encoding reaches the file limit for 148 out of 200 histories with at least 12 sessions (Figure 3.14a) and for 50 out of 100 histories with 60 transactions per session (Figure 3.14b), the other parameters being fixed as explained above.

We have found a large number of violations, whose frequency increases with the number of sessions, transactions per session, or operations per transaction, and decreases when allowing more variables. This is expected since increasing any of the former parameters increases the chance of interference between different transactions while increasing the latter has the opposite effect. The second and third column of Table 3.3 give a more precise account of the kind of violations we found by identifying for each criterion X , the number of histories that violate X but no other criterion weaker than X , e.g., there is only one violation to SI that satisfies PC.

The second experiment measures the scalability of the SI checking algorithm obtained by applying the reduction to SER described in Section 3.5.3 followed by the SER checking algorithm in Algorithm 5, and its performance compared to a SAT encoding of SI. Actually, the reduction to SER is performed on-the-fly, while traversing the history and checking for serializability (of the transformed history). The SAT encoding follows the same principles as in the case of serializability. We focus on its behavior when increasing the number of sessions (varying the other parameters leads to similar results). As benchmark, we used the same CockroachDB histories as in Figure 3.14a and a number of histories extracted from Galera¹⁵ whose documentation contains contradicting claims about whether it implements snapshot isolation [**galera-claim**, **galera-notclaim**]. We use 100 histories per combination of parameter values as in the previous experiment. The results are reported in Figure 3.15a and Figure 3.15b. We observe the same behavior as in the case of SER.

¹⁵In order to increase the frequency of valid histories, all sessions are executed on a single node.

Table 3.3: Violation statistics. The “disjoint writes” columns refer to histories where the set of variables written by any two sessions are disjoint.

Weakest criterion violated	Serializability checking		Snapshot Isolation checking	
	CockroachDB (disjoint writes)	CockroachDB (no constraints)	Galera (disjoint writes)	Galera (no constraints)
Read Committed			19	50
Read Atomic	180	547	91	139
Causal Consistency	339	382	88	43
Prefix Consistency	2	7		
Snapshot Isolation		1		1
Serializability	25			
Total number of violations	546/1000	937/1000	198/250	233/250

In particular, the SAT encoding reaches the file limit for 150 out of 200 histories with at least 12 sessions in the case of the CockroachDB histories, and for 162 out of 300 histories with at least 9 sessions in the case of the Galera histories. The last two columns in Table 3.3 classify the set of violations depending on the weakest criterion that they violate.

We also evaluated the performance of the CC checking algorithm in Section 3.4 when increasing the number of sessions, on histories extracted from AntidoteDB, which claims to implement causal consistency [antidote-claim]. The results are reported in Figure 3.15c. In this case, the SAT encoding reaches the file limit for 150 out of 300 histories with at least 9 sessions. All the histories considered in this experiment are valid. However, when experimenting with other parameter values, we have found several violations. The smallest parameter values for which we found violations were 3 sessions, 14 transactions per session, 14 operations per transaction, and 5 variables. The violations we found are also violations of Read Atomic. For instance, one of the violations contains two transactions t_1 and t_2 , each of them writing to two variables x_1 and x_2 , and another transaction t_3 that reads x_1 from t_1 and x_2 from t_2 (t_1 and t_2 are from different sessions while t_3 is an so successor of t_1 in the same session). These violations are novel and they were confirmed by the developers of AntidoteDB.

The refinement of the algorithms above based on communication graphs, described in Section 3.6, did not have a significant impact on their performance. The histories we generated contained few biconnected components (many histories contained just a single biconnected component) which we believe is due to our proof of concept deployment of these databases on a single machine that did not allow to experiment with very large number of sessions and variables.

3.8 RELATED WORK

[DBLP:conf/concur/Cerone0G15] give the first formalization of the criteria we consider in this paper, using the specification methodology of [DBLP:conf/popl/BurckhardtGYZ14]. This formalization uses two auxiliary relations, a *visibility* relation which represents the fact that a transaction “observes” the effects of another transaction and a *commit order*, also called arbitration order, like in our case. Executions are abstracted using a notion of history that includes only a session order and the adherence to some consistency criterion is defined as the existence of a vis-

ibility relation and a commit order satisfying certain axioms. Motivated by practical goals, our histories include a write-read relation, which enables more uniform and in our opinion, more intuitive, axioms to characterize consistency criteria. Our formalizations are however equivalent with those of [DBLP:conf/concur/Cerone0G15] (a formal proof of this equivalence is presented in the extended version of this paper [DBLP:journals/corr/abs-1908-04509]). Moreover, [DBLP:conf/concur/Cerone0G15] do not investigate algorithmic issues as in our paper.

[DBLP:journals/jacm/Papadimitriou79b] showed that checking serializability of an execution is NP-complete. Moreover, it identifies a stronger criterion called *conflict serializability* which is polynomial-time checkable. Conflict serializability assumes that histories are given as sequences of operations and requires that the commit order be consistent with a *conflict-order* between transactions defined based on this sequence (roughly, a transaction t_1 is before a transaction t_2 in the conflict order if it accesses some variable x before t_2 does). This result is not applicable to distributed databases where deriving such a sequence between operations submitted to different nodes in a network is impossible.

[DBLP:conf/popl/BouajjaniEGH17] showed that checking several variations of causal consistency on executions of a *non-transactional* distributed database is polynomial time (they also assume that every value is written at most once). Assuming singleton transactions, our notion of CC corresponds to the causal convergence criterion in [DBLP:conf/popl/BouajjaniEGH17]. Therefore, our result concerning CC can be seen as an extension of this result concerning causal convergence to transactions.

There are some works that investigated the problem of checking consistency criteria like sequential consistency and linearizability in the case of shared-memory systems. [DBLP:journals/siamcomp/GibbonsK97] showed that checking linearizability of the single-value register type is NP-complete in general, but polynomial time for executions where every value is written at most once. Using a reduction from serializability, they showed that checking sequential consistency is NP-complete even when every value is written at most once. [DBLP:journals/pacmpl/EmmiE18] extended the result concerning linearizability to a series of abstract data types called collections, that includes stacks, queues, key-value maps, etc. Sequential consistency reduces to serializability for histories with singleton transactions (i.e., formed of a single read or write operation). Therefore, our polynomial-time result for checking serializability of bounded-width histories (Corollary 3.5.1) implies that checking sequential consistency of histories with a bounded number of threads is polynomial time. The latter result has been established independently by [DBLP:journals/pacmpl/AbdullaAJLNS19].

The notion of *communication graph* is inspired by the work of [DBLP:journals/pacmpl/ChalupaCPSV18], which investigates partial-order reduction (POR) techniques for multi-threaded programs. In general, the goal of partial-order reduction [DBLP:conf/popl/FlanaganG05] is to avoid exploring executions which are equivalent w.r.t. some suitable notion of equivalence, e.g., Mazurkiewicz trace equivalence [DBLP:conf/ac/Mazurkiewicz86]. They use the acyclicity of communication graphs to define a class of programs for which their POR technique is optimal. The algorithmic issues they explore are different than ours and they don't investigate biconnected components of this graph as in our results.

4 APPLICATIONS IMPLEMENTED ON DISTRIBUTED DATASTORES

4.1 INTRODUCTION

Data storage is no longer about writing data to a single disk with a single point of access. Modern applications require not just data reliability, but also high-throughput concurrent accesses. Applications concerning supply chains, banking, etc. use traditional relational databases for storing and processing data, whereas applications such as social networking software and e-commerce platforms use cloud-based storage systems (such as Azure CosmosDb [cosmosdb], Amazon DynamoDb [DBLP:conf/sosp/DeCandiaHJKLPSVV07], Facebook TAO [DBLP:conf/usenix/BronsonACCDDFGKLMPPSV1], etc.). We use the term *storage system* in this paper to refer to any such database system/service.

Providing high-throughput processing, unfortunately, comes at an unavoidable cost of weakening the guarantees offered to users. Concurrently-connected clients may end up observing different views of the same data. These “anomalies” can be prevented by using a strong *isolation level* such as *serializability*, which essentially offers a single view of the data. However, serializability requires expensive synchronization and incurs a high performance cost. As a consequence, most storage systems use weaker isolation levels, such as *Causal Consistency* [DBLP:journals/cacm/Lamport78, DBLP:conf/sosp/LloydFKA11, antidote-white-paper], *Snapshot Isolation* [DBLP:conf/sigmod/BerensonBGMOO95], *Read Committed* [DBLP:conf/sigmod/BerensonBGMOO95], etc. for better performance. In a recent survey of database administrators [DBLP:conf/sigmod/Pavlo17], 86% of the participants responded that most or all of the transactions in their databases execute at read committed isolation level.

A weaker isolation level allows for more possible behaviors than stronger isolation levels. It is up to the developers then to ensure that their application can tolerate this larger set of behaviors. Unfortunately, weak isolation levels are hard to understand or reason about [DBLP:conf/popl/BrutschyDOV17, adya-thesis] and resulting application bugs can cause loss of business [DBLP:conf/sigmod/WarszawskiB17]. Consider a simple shopping cart application that stores a per-client shopping cart in a key-value store (*key* is the client ID and *value* is a multi-set of items). Figure 4.1 shows procedures for adding an item to the cart (`AddItem`) and deleting *all* instances of an item from the cart (`DeleteItem`). Each procedure executes in a transaction, represented by the calls to `Begin` and `Commit`. Suppose that initially, a user u has a single instance of item I in their cart. Then the user connects to the application via two different sessions (for instance, via two browser windows), adds I in one session (`AddItem(I , u)`) and deletes I in the other session (`DeleteItem(I , u)`). With serializability, the cart can either be left in the state $\{I\}$ (delete happened first, followed by the add) or \emptyset (delete happened second). However, with causal consistency (or read committed), it is possible that with two sequential reads of the shopping cart, the cart is empty in the first read (signaling that the delete has

4 Applications implemented on Distributed Datastores

```
// Append item to cart          // Fetch cart and delete item
AddItem(item i, userId) {      DeleteItem(item i, userId) {
  Begin()                      Begin()
  key = "cart:" + userId       key = "cart:" + userId
  cart = read(key)             cart = read(key)
  cart.append(i)               cart.remove(i)
  write(key, cart)             write(key, cart)
  Commit()                     Commit()
}                               }
```

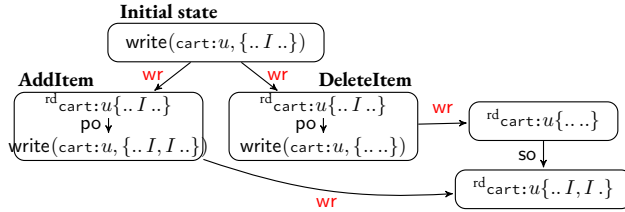


Figure 4.1: A simple shopping cart service.

succeeded), but there are *two* instances of I in the second read! Such anomalies, of deleted items reappearing, have been noted in previous work [DBLP:conf/sosp/DeCandiaHJKLPSVV07].

TESTING STORAGE-BASED APPLICATIONS This paper addresses the problem of *testing* code for correctness against weak behaviors: a developer should be able to write a test that runs their application and then asserts for correct behavior. The main difficulty today is getting coverage of weak behaviors during the test. If one runs the test against the actual production storage system, it is very likely to only result in serializable behaviors because of their optimized implementation. For instance, only 0.0004% of all reads performed on Facebook’s TAO storage system were not serializable [DBLP:conf/sosp/LuVAHSTKL15]. Emulators, offered by cloud providers for local development, on the other hand, do not support weaker isolation levels at all [cosmosdb-local]. Another option, possible when the storage system is available open-source, is to set it up with a tool like Jepsen [jepsen] to inject noise (bring down replicas or delay packets on the network). This approach is unable to provide good coverage at the level of client operations [DBLP:journals/pacmpl/RahmaniNDJ19] (§4.7). Another line of work has focussed on finding anomalies by identifying non-serializable behavior (§4.8). Anomalies, however, do not always correspond to bugs [DBLP:conf/pldi/BrutschyD0V18, DBLP:journals/pvldb/GanRRB020]; they may either not be important (e.g., gather statistics) or may already be handled in the application (e.g., checking and deleting duplicate items).

We present MonkeyDB, a mock in-memory storage system meant for testing correctness of storage-backed applications. MonkeyDB supports common APIs for accessing data (key-value updates, as well as SQL queries), making it an easy substitute for an actual storage system. MonkeyDB can be configured with one of several isolation levels. On a read operation, MonkeyDB computes the set of all possible return values allowed under the chosen isolation level, and randomly returns one of them. The developer can then simply execute their test multiple times to get coverage of possible weak behaviors. For the program in Figure 4.1, if we write a test asserting that two sequential reads cannot return empty-cart followed by $\{I, I\}$, then it takes only 20 runs of the test (on average) to fail the assert. In contrast, the test does not fail when using MySQL with read committed, even after 100k runs.

DESIGN OF MONKEYDB MonkeyDB does not rely on stress generation, fault injection, or data replication. Rather, it works directly with a formalization of the given isolation level in order to compute allowed return values.

The theory behind MonkeyDB builds on the axiomatic definitions of isolation levels introduced by Biswas et al. [DBLP:journals/pacmpl/BiswasE19]. These definitions use logical constraints (called *axioms*) to characterize the set of executions of a key-value store that conform to a particular isolation level (we discuss SQL queries later). These constraints refer to a specific set of relations between events/transactions in an execution that describe control-flow or data-flow dependencies: a program order *po* between events in the same transaction, a session order *so* between transactions in the same session¹, and a write-read *wr* (read-from) relation that associates each read event with a transaction that writes the value returned by the read. These relations along with the events (also called, operations) in an execution are called a *history*. The history corresponding to the shopping cart anomaly explained above is given on the bottom of Figure 4.1. Read operations include the read value, and boxes group events from the same transaction. A history describes only the interaction with the key-value store, omitting application side events (e.g., computing the value to be written to a key).

MonkeyDB implements a *centralized* operational semantics for key-value stores, which is based on these axiomatic definitions. Transactions are executed *serially*, one after another, the concurrency being simulated during the handling of read events. This semantics maintains a history that contains all the past events (from all transactions/sessions), and write events are simply added to the history. The value returned by a read event is established based on a non-deterministic choice of a write-read dependency (concerning this read event) that satisfies the axioms of the considered isolation level. Depending on the weakness of the isolation level, this makes it possible to return values written in arbitrarily “old” transactions, and simulate any concurrent behavior. For instance, the history in Figure 4.1 can be obtained by executing `AddItem`, `DeleteItem`, and then the two reads (serially). The read in `DeleteItem` can take its value from the initial state and “ignore” the previously executed `AddItem`, because the obtained history validates the axioms of causal consistency (or read committed). The same happens for the two later reads in the same session, the first one being able to read from `DeleteItem` and the second one from `AddItem`.

We formally prove that this semantics does indeed simulate any concurrent behavior, by showing that it is equivalent to a semantics where transactions are allowed to interleave. In comparison with concrete implementations, this semantics makes it possible to handle a wide range of isolation levels in a uniform way. It only has two sources of non-determinism: the order in which entire transactions are submitted, and the choice of write-read dependencies in read events. This enable better coverage of possible behaviors, the penalty in performance not being an issue in safety testing workloads which are usually small (see our evaluation).

We also extend our semantics to cover SQL queries as well, by compiling SQL queries down to transactions with multiple key-value reads/writes. A table in a relational database is represented using a set of primary key values (identifying uniquely the set of rows) and a set of keys, one for each cell in the table. The set of primary key values is represented using a set of Boolean key-value pairs that simulate its characteristic function (adding or removing an element corresponds to updating one of these keys to true or false). Then, SQL queries are compiled to read or write

¹A session is a sequential interface to the storage system. It corresponds to what is also called a connection.

accesses to the keys representing a table. For instance, a `SELECT` query that retrieves the set of rows in a table that satisfy a `WHERE` condition is compiled to (1) reading Boolean keys to identify the primary key values of the rows contained in the table, (2) reading keys that represent columns used in the `WHERE` condition, and (3) reading all the keys that represent cells in a row satisfying the `WHERE` condition. This rewriting contains the minimal set of accesses to the cells of a table that are needed to ensure the conventional specification of SQL. It makes it possible to “export” formalizations of key-value store isolation levels to SQL transactions.

CONTRIBUTIONS This paper makes the following contributions:

- We define an operational semantics for key-value stores under various isolation levels, which simulates all concurrent behaviors with executions where transactions execute serially (§4.3) and which is based on the axiomatic definitions in [DBLP:journals/pacmpl/BiswasE19] (and outlined in §3.3),
- We broaden the scope of the key-value store semantics to SQL transactions using a compiler that rewrites SQL queries to key-value accesses (§4.4),
- The operational semantics and the SQL compiler are implemented in a tool called MonkeyDB (§4.5). It randomly resolves possible choices to provide coverage of weak behaviors. It supports both a key-value interface as well as SQL, making it readily compatible with any storage-backed application.
- We present an evaluation of MonkeyDB on several applications, showcasing its superior coverage of weak behaviors as well as bug-finding abilities (§4.6, §4.7).²

4.2 PROGRAMMING LANGUAGE

Figure 4.2 lists the definition of two simple programming languages that we use to represent applications running on top of Key-Value or SQL stores, respectively. A program is a set of *sessions* running in parallel, each session being composed of a sequence of *transactions*. Each transaction is delimited by `begin` and `commit` instructions³, and its body contains instructions that access the store, and manipulate a set of local variables `Vars` ranged over using x, y, \dots

In case of a program running on top of a Key-Value store, the instructions can be reading the value of a key and storing it to a local variable x ($x := \text{read}(k)$), writing the value of a local variable x to a key (`write(k, x)`), or an assignment to a local variable x . The set of values of keys or local variables is denoted by `Vals`. Assignments to local variables use expressions interpreted as values whose syntax is left unspecified. Each of these instructions can be guarded by a Boolean condition $\phi(\vec{x})$ over a set of local variables \vec{x} (their syntax is not important). Other constructs like `while` loops can be defined in a similar way. Let \mathcal{P}_{KV} denote the set of programs where a transaction body can contain only such instructions.

²Source code of our benchmarks is available as supplementary material.

³For simplicity, we assume that all the transactions in the program commit. Aborted transactions can be ignored when reasoning about safety because their effects should be invisible to other transactions.

$$\begin{aligned}
k &\in \text{Keys} & x &\in \text{Vars} & \text{tab} &\in \mathbb{T} & \vec{c}, \vec{c}_1, \vec{c}_2 &\in \mathbb{C}^* \\
\\
\text{Prog} &\stackrel{\text{def}}{=} \text{Sess} \mid \text{Sess} \parallel \text{Prog} \\
\text{Sess} &\stackrel{\text{def}}{=} \text{Trans} \mid \text{Trans}; \text{Sess} \\
\text{Trans} &\stackrel{\text{def}}{=} \text{begin}; \text{Body}; \text{commit} \\
\text{Body} &\stackrel{\text{def}}{=} \text{Instr} \mid \text{Instr}; \text{Body} \\
\text{Instr} &\stackrel{\text{def}}{=} \text{InstrKV} \mid \text{InstrSQL} \mid x := e \mid \text{if}(\phi(\vec{x}))\{\text{Instr}\} \\
\text{InstrKV} &\stackrel{\text{def}}{=} x := \text{read}(k) \mid \text{write}(k, x) \\
\text{InstrSQL} &\stackrel{\text{def}}{=} \text{SELECT } \vec{c}_1 \text{ AS } x \text{ FROM } \text{tab} \text{ WHERE } \phi(\vec{c}_2) \mid \\
&\quad \text{INSERT INTO } \text{tab} \text{ VALUES } \vec{x} \mid \\
&\quad \text{DELETE FROM } \text{tab} \text{ WHERE } \phi(\vec{c}) \mid \\
&\quad \text{UPDATE } \text{tab} \text{ SET } \vec{c}_1 = \vec{x} \text{ WHERE } \phi(\vec{c}_2)
\end{aligned}$$

Figure 4.2: Program syntax. The set of all keys is denoted by Keys , Vars denotes the set of local variables, \mathbb{T} the set of table names, and \mathbb{C} the set of column names. We use ϕ to denote Boolean expressions, and e to denote expressions interpreted as values. We use $\vec{\cdot}$ to denote vectors of elements.

For programs running on top of SQL stores, the instructions include simplified versions of standard SQL instructions and assignments to local variables. These programs run in the context of a *database schema* which is a (partial) function $\mathcal{S} : \mathbb{T} \rightarrow 2^{\mathbb{C}}$ mapping table names in \mathbb{T} to sets of column names in \mathbb{C} . The SQL store is an *instance* of a database schema \mathcal{S} , i.e., a function $\mathcal{D} : \text{dom}(\mathcal{S}) \rightarrow 2^{\mathbb{R}}$ mapping each table tab in the domain of \mathcal{S} to a set of *rows* of tab , i.e., functions $r : \mathcal{S}(\text{tab}) \rightarrow \text{Vals}$. We use \mathbb{R} to denote the set of all rows. The **SELECT** instruction retrieves the columns \vec{c}_1 from the set of rows of tab that satisfy $\phi(\vec{c}_2)$ (\vec{c}_2 denotes the set of columns used in this Boolean expression), and stores them into a variable x . **INSERT** adds a new row to tab with values \vec{x} , and **DELETE** deletes all rows from tab that satisfy a condition $\phi(\vec{c})$. The **UPDATE** instruction assigns the columns \vec{c}_1 of all rows of tab that satisfy $\phi(\vec{c}_2)$ with values in \vec{x} . Let \mathcal{P}_{SQL} denote the set of programs where a transaction body can contain only such instructions.

4.3 OPERATIONAL SEMANTICS FOR \mathcal{P}_{KV}

We define a small-step operational semantics for Key-Value store programs, which is parametrized by an isolation level I . Transactions are executed *serially* one after another, and the values returned by read operations are decided using the axiomatic definition of I . The semantics maintains a history of previously executed operations, and the value returned by a read is chosen non-deterministically as long as extending the current history with the corresponding write-read dependency satisfies the axioms of I . We show that this semantics is sound and complete for any natural isolation level I , i.e., it generates precisely the same set of histories as a *baseline* semantics where transactions can interleave arbitrarily and the read operations can return arbitrary values as long as they can be proved to be correct at the end of the execution.

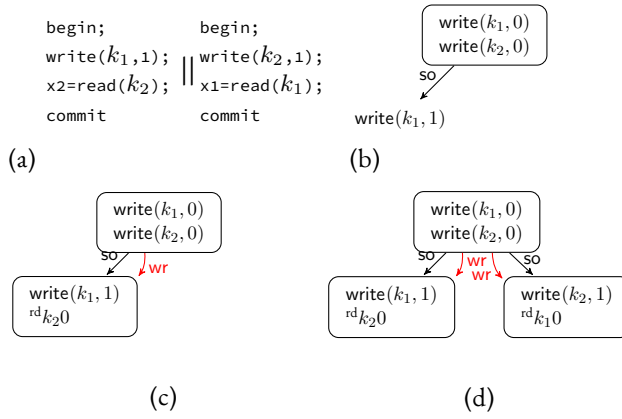


Figure 4.3: The Causal semantics on the program in (a), assuming that the transaction on the left is scheduled first.

4.3.1 DEFINITION OF THE OPERATIONAL SEMANTICS

We use the program in Figure 4.3a to give an overview of our semantics, assuming Causal Consistency. This program has two concurrent transactions whose reads can both return the initial value 0, which is not possible under Serializability.

Our semantics executes transactions in their entirety one after another (without interleaving them), maintaining a history that contains all the executed operations. We assume that the transaction on the left executes first. Initially, the history contains a fictitious transaction log that writes the initial value 0 to all keys, and that will precede all the transaction logs created during the execution in session order.

Executing a write instruction consists in simply appending the corresponding write operation to the log of the current transaction. For instance, executing the first write (and `begin`) in our example results in adding a transaction log that contains a write operation (see Figure 4.3b). The execution continues with the read instruction from the same transaction, and it cannot switch to the other transaction.

The execution of a read instruction consists in choosing non-deterministically a write-read dependency that validates Causal when added to the current history. In our example, executing `read(k2)` results in adding a write-read dependency from the transaction log writing initial values, which determines the return value of the read (see Figure 4.3c). This choice makes the obtained history satisfy Causal.

The second transaction executes in a similar manner. When executing its read instruction, the chosen write-read dependency is again related to the transaction log writing initial values (see Figure 4.3d). This choice is valid under Causal. Since a read must not read from the preceding transaction, this semantics is able to simulate all the “anomalies” of a weak isolation level (this execution being an example).

Formally, the operational semantics is defined as a transition relation \Rightarrow_I between *configurations*, which are defined as tuples containing the following:

- history h storing the operations executed in the past,
- identifier j of the current session,

$$\begin{array}{c}
\text{SPAWN} \\
\frac{t \text{ fresh} \quad P(j) = \text{begin}; \text{Body}; \text{commit}; S}{h, _, _, \epsilon, P \Rightarrow_I h \oplus_j \langle t, \emptyset, \emptyset \rangle, j, \emptyset, \text{Body}, P[j \mapsto S]} \\
\\
\text{IF-TRUE} \\
\frac{\varphi(\vec{x})[x \mapsto \gamma(x) : x \in \vec{x}] \text{ true}}{h, j, \gamma, \text{if}(\phi(\vec{x}))\{\text{Instr}\}; B, P \Rightarrow_I h, j, \gamma, \text{Instr}; B, P} \\
\\
\text{IF-FALSE} \\
\frac{\varphi(\vec{x})[x \mapsto \gamma(x) : x \in \vec{x}] \text{ false}}{h, j, \gamma, \text{if}(\phi(\vec{x}))\{\text{Instr}\}; B, P \Rightarrow_I h, j, \gamma, B, P} \\
\\
\text{WRITE} \\
\frac{v = \gamma(x) \quad i \text{ fresh}}{h, j, \gamma, \text{write}(k, x); B, P \Rightarrow_I h \oplus_j \text{write}_i(k, v), j, \gamma, B, P} \\
\\
\text{READ-LOCAL} \\
\frac{\text{write}(k, v) \text{ is the last write on } k \text{ in } t \text{ w.r.t. po} \quad i \text{ fresh}}{h, j, \gamma, x := \text{read}(k); B, P \Rightarrow_I h \oplus_j \text{rd}[i]kv, j, \gamma[x \mapsto v], B, P} \\
\\
\text{READ-EXTERN} \\
\frac{\begin{array}{c} h = (T, \text{so}, \text{wr}) \\ t \text{ is the id of the last transaction log in so}(j) \quad \text{write}(k, v) \in \text{writes}(t') \text{ with } t' \in T \text{ and } t' \neq t \\ i \text{ fresh} \quad h' = (h \oplus_j \text{rd}[i]kv) \oplus \text{wr}(t', \text{rd}[i]kv) \quad h' \text{ satisfies } I \end{array}}{h, j, \gamma, x := \text{read}(k); B, P \Rightarrow_I h', j, \gamma[x \mapsto v], B, P}
\end{array}$$

Figure 4.4: Operational semantics for \mathcal{P}_{KV} programs under isolation level I . For a function $f : A \rightarrow B$, $f[a \mapsto b]$ denotes the function $f' : A \rightarrow B$ defined by $f'(c) = f(c)$, for every $c \neq a$ in the domain of f , and $f'(a) = b$.

- local variable valuation γ for the current transaction,
- code B that remains to be executed from the current transaction, and
- sessions/transactions P that remain to be executed from the original program.

For readability, we define a program as a partial function $P : \text{SessId} \rightarrow \text{Sess}$ that associates session identifiers in SessId with concrete code as defined in Figure 4.2 (i.e., sequences of transactions). Similarly, the session order so in a history is defined as a partial function $\text{so} : \text{SessId} \rightarrow \text{Tlogs}^*$ that associates session identifiers with sequences of transaction logs. Two transaction logs are ordered by so if one occurs before the other in some sequence $\text{so}(j)$ with $j \in \text{SessId}$.

Before presenting the definition of \Rightarrow_I , we introduce some notation. Let h be a history that contains a representation of so as above. We use $h \oplus_j \langle t, O, \text{po} \rangle$ to denote a history where $\langle t, O, \text{po} \rangle$ is appended to $\text{so}(j)$. Also, for an operation o , $h \oplus_j o$ is the history obtained from h by adding o to the last transaction log in $\text{so}(j)$ and as a last operation in the program order of this log (i.e., if $\text{so}(j) = \sigma; \langle t, O, \text{po} \rangle$, then the session order so' of $h \oplus_j o$ is defined by $\text{so}'(k) = \text{so}(k)$ for all $k \neq j$ and $\text{so}(j) = \sigma; \langle t, O \cup o, \text{po} \cup \{(o', o) : o' \in O\} \rangle$). Finally, for a history $h = \langle T, \text{so}, \text{wr} \rangle$, $h \oplus \text{wr}(t, o)$ is the history obtained from h by adding (t, o) to the write-read relation.

$$\begin{array}{c}
 \text{SPAWN*} \\
 \frac{t \text{ fresh} \quad P(j) = \text{begin}; \text{Body}; \text{commit}; S \quad \vec{B}(j) = \epsilon}{h, \vec{\gamma}, \vec{B}, P \Rightarrow h \oplus_j \langle t, \emptyset, \emptyset \rangle, \vec{\gamma}[j \mapsto \emptyset], \vec{B}[j \mapsto \text{Body}], P[j \mapsto S]} \\
 \\
 \text{READ-EXTERN*} \\
 \frac{\vec{B}(j) = x := \text{read}(k); B \quad h = (T, \text{so}, \text{wr}) \quad t \text{ is the id of the last transaction log in } \text{so}(j) \quad \text{write}(k, v) \in \text{writes}(t') \text{ with } t' \in \text{compTrans}(h, \vec{B}) \text{ and } t \neq t' \quad i \text{ fresh} \quad h' = (h \oplus_j^{\text{rd}}[i]kv) \oplus \text{wr}(t', \text{rd}[i]kv)}{h, \vec{\gamma}, \vec{B}, P \Rightarrow h', \vec{\gamma}[(j, x) \mapsto v], \vec{B}[j \mapsto B], P}
 \end{array}$$

Figure 4.5: A baseline operational semantics for \mathcal{P}_{KV} programs. Above, $\text{compTrans}(h, \vec{B})$ denotes the set of transaction logs in h that excludes those corresponding to live transactions, i.e., transaction logs $t'' \in T$ such that t'' is the last transaction log in some $\text{so}(j')$ and $\vec{B}(j') \neq \epsilon$.

Figure 4.4 lists the rules defining \Rightarrow_I . The SPAWN rule starts a new transaction, provided that there is no other live transaction ($B = \epsilon$). It adds an empty transaction log to the history and schedules the body of the transaction. IF-TRUE and IF-FALSE check the truth value of a Boolean condition of an if conditional. WRITE corresponds to a write instruction and consists in simply adding a write operation to the current history. READ-LOCAL and READ-EXTERN concern read instructions. READ-LOCAL handles the case where the read follows a write on the same key k in the same transaction: the read returns the value written by the last write on k in the current transaction. Otherwise, READ-EXTERN corresponds to reading a value written in another transaction t' (t is the id of the log of the current transaction). The transaction t' is chosen non-deterministically as long as extending the current history with the write-read dependency associated to this choice leads to a history that still satisfies I .

An *initial* configuration for program P contains the program P along with a history $h = \langle \{t_0\}, \emptyset, \emptyset \rangle$, where t_0 is a transaction log containing only writes that write the initial values of all keys, and empty current transaction code ($B = \epsilon$). An execution of a program P under an isolation level I is a sequence of configurations $c_0 c_1 \dots c_n$ where c_0 is an initial configuration for P , and $c_m \Rightarrow_I c_{m+1}$, for every $0 \leq m < n$. We say that c_n is *I-reachable* from c_0 . The history of such an execution is the history h in the last configuration c_n . A configuration is called *final* if it contains the empty program ($P = \emptyset$). Let $\text{hist}_I(P)$ denote the set of all histories of an execution of P under I that ends in a final configuration.

4.3.2 CORRECTNESS OF THE OPERATIONAL SEMANTICS

We define the correctness of \Rightarrow_I in relation to a *baseline* semantics where transactions can interleave arbitrarily, and the values returned by read operations are only constrained to come from committed transactions. This semantics is represented by a transition relation \Rightarrow , which is defined by a set of rules that are analogous to \Rightarrow_I . Since it allows transactions to interleave, a configuration contains a history h , the sessions/transactions P that remain to be executed, and:

- a valuation map $\vec{\gamma}$ that records local variable values in the current transaction of each session ($\vec{\gamma}$ associates identifiers of sessions that have live transactions with valuations of local variables),

- a map \vec{B} that stores the code of each live transaction (associating session identifiers with code).

Figure 4.5 lists some rules defining \Rightarrow (the others can be defined in a similar manner). SPAWN^* starts a new transaction in a session j provided that this session has no live transaction ($\vec{B}(j) = \epsilon$). Compared to SPAWN in Figure 4.4, this rule allows unfinished transactions in other sessions. READ-EXTERN^* does not check conformance to I , but it allows a read to only return a value written in a completed (committed) transaction. In this work, we consider only isolation levels satisfying this constraint. Executions, initial and final configurations are defined as in the case of \Rightarrow_I . The history of an execution is still defined as the history in the last configuration. Let $\text{hist}_*(P)$ denote the set of all histories of an execution of P w.r.t. \Rightarrow that ends in a final configuration.

Practical isolation levels satisfy a “prefix-closure” property saying that if the axioms of I are satisfied by a pair $\langle h_2, \text{co}_2 \rangle$, then they are also satisfied by every *prefix* of $\langle h_2, \text{co}_2 \rangle$. A prefix of $\langle h_2, \text{co}_2 \rangle$ contains a prefix of the sequence of transactions in h_2 when ordered according to co_2 , and the last transaction log in this prefix is possibly incomplete. In general, this prefix-closure property holds for isolation levels I that are defined by axioms as in (??), provided that the property $\phi(t_2, \alpha)$ is *monotonic*, i.e., the set of models in the context of a pair $\langle h_2, \text{co}_2 \rangle$ is a *superset* of the set of models in the context of a prefix $\langle h_1, \text{co}_1 \rangle$ of $\langle h_2, \text{co}_2 \rangle$. For instance, the property ϕ in the axiom defining Causal is $(t_2, \alpha) \in (\text{wr} \cup \text{so})^+$, which is clearly monotonic. In general, standard isolation levels are defined using a property α of the form $(t_2, \alpha) \in R$ where R is an expression built from the relations po , so , wr , and co using (reflexive and) transitive closure and composition of relations [DBLP:journals/pacmpl/BiswasE19]. Such properties are monotonic in general (they would not be if those expressions would use the negation/complement of a relation). An axiom as in (??) is called *monotonic* when the property ϕ is monotonic.

The following theorem shows that $\text{hist}_I(P)$ is precisely the set of histories under the baseline semantics, which satisfy I (the validity of the reads is checked at the end of an execution), provided that the axioms of I are monotonic.

Theorem 4.3.1. *For any isolation level I defined by a set of monotonic axioms, $\text{hist}_I(P) = \{h \in \text{hist}_*(P) : h \text{ satisfies } I\}$.*

The \subseteq direction follows mostly from the fact that \Rightarrow_I is more constrained than \Rightarrow . For the opposite direction, given a history h that satisfies I , i.e., there exists a commit order co such that $\langle h, \text{co} \rangle$ satisfies the axioms of I , we can show that there exists an execution under \Rightarrow_I with history h , where transactions execute serially in the order defined by co . The prefix closure property is used to prove that READ-EXTERN transitions are enabled (these transitions get executed with a prefix of h). See the supplementary material for more details.

It can also be shown that \Rightarrow_I is *deadlock-free* for every natural isolation level (e.g., Read Committed, Causal Consistency, Snapshot Isolation, and Serializability), i.e., every read can return some value satisfying the axioms of I at the time when it is executed (independently of previous choices).

Table:			Intermediate representation:	
A			A = { 1, 2, 3 }	
Id	Name	City		
1	Alice	Paris	A.1.Id: 1, A.1.Name: Alice,	A.1.City: Paris
2	Bob	Bangalore	A.2.Id: 2, A.2.Name: Bob,	A.2.City: Bangalore
3	Charles	Bucharest	A.3.Id: 3, A.3.Name: Charles, A.3.City: Bucharest	

Figure 4.6: Representing tables with set variables and key-value pairs. We write a key-value pair as key:value.

4.4 COMPILING SQL TO KEY-VALUE API

We define an operational semantics for SQL programs (in \mathcal{P}_{SQL}) based on a compiler that rewrites SQL queries to Key-Value `read` and `write` instructions. For presentation reasons, we use an intermediate representation where each table of a database instance is represented using a *set* variable that stores values of the primary key⁴ (identifying uniquely the rows in the table) and a set of key-value pairs, one for each cell in the table. In a second step, we define a rewriting of the API used to manipulate set variables into Key-Value `read` and `write` instructions.

INTERMEDIATE REPRESENTATION Let $\mathcal{S} : \mathbb{T} \rightarrow 2^{\mathbb{C}}$ be a database schema (recall that \mathbb{T} and \mathbb{C} are the set of table names and column names, resp.). For each table *tab*, let *tab.pkey* be the name of the primary key column. We represent an instance $\mathcal{D} : \text{dom}(\mathcal{S}) \rightarrow 2^{\mathbb{R}}$ using:

- for each table *tab*, a set variable *tab* (with the same name) that contains the primary key value $r(\text{tab.pkey})$ of every row $r \in \mathcal{D}(\text{tab})$,
- for each row $r \in \mathcal{D}(\text{tab})$ with primary key value $\text{pkeyVal} = r(\text{tab.pkey})$, and each column $c \in \mathcal{S}(\text{tab})$, a key *tab.pkeyVal.c* associated with the value $r(c)$.

Example 4.4.1. The table *A* on the left of Figure 4.6, where the primary key is defined by the *Id* column, is represented using a set variable *A* storing the set of values in the column *Id*, and one key-value pair for each cell in the table.

Figure 4.7 lists our rewriting of SQL queries over a database instance \mathcal{D} to programs that manipulate the set variables and key-value pairs described above. This rewriting contains the minimal set of accesses to the cells of a table that are needed to implement an SQL query according to its conventional specification. To manipulate set variables, we use `add` and `remove` for adding and removing elements, respectively (returning true or false when the element is already present or deleted from the set, respectively), and `elements` that returns all of the elements in the input set⁵.

`SELECT`, `DELETE`, and `UPDATE` start by reading the contents of the set variable storing primary key values and then, for every row, the columns in \vec{c}_2 needed to check the Boolean condition ϕ (the keys corresponding to these columns). For every row satisfying this Boolean condition, `SELECT` continues by reading the keys associated to the columns that need to be returned,

⁴For simplicity, we assume that primary keys correspond to a single column in the table.

⁵`add(s, e)` and `remove(s, e)` add and remove the element *e* from *s*, respectively. `elements(s)` returns the content of *s*.

SELECT/DELETE/UPDATE

```

rows := elements(tab)
for ( let pkeyVal of rows ) {
  for ( let c of  $\vec{c}_2$  ) {
    val[c] := read(tab.pkeyVal.c)
    if (  $\phi[c \mapsto \text{val}[c] : c \in \vec{c}_2]$  true )
      // SELECT  $\vec{c}_1$  AS  $x$  FROM  $tab$  WHERE  $\phi(\vec{c}_2)$ 
      for ( let c of  $\vec{c}_1$  )
        out[c] := read(tab.pkeyVal.c)
      x := x  $\cup$  out
      // DELETE FROM  $tab$  WHERE  $\phi(\vec{c}_2)$ 
      remove(tab, pkeyVal);
      // UPDATE  $tab$  SET  $\vec{c}_1 = \vec{x}$  WHERE  $\phi(\vec{c}_2)$ 
      for ( let c of  $\vec{c}_1$  )
        write( tab.pkeyVal.c,  $\gamma(\vec{x}[c])$  )
  }
}

```

INSERT INTO tab VALUES \vec{x}

```

pkeyVal :=  $\gamma(\vec{x}[0])$ 
if ( add(tab, pkeyVal) ) {
  for ( let c of  $\mathcal{S}(tab)$  ) {
    write( tab.pkeyVal.c,  $\gamma(\vec{x}[c])$  )
  }
}

```

Figure 4.7: Compiling SQL queries to the intermediate representation. Above, γ is a valuation of local variables. Also, in the case of INSERT, we assume that the first element of \vec{x} represents the value of the primary key.

```

add(tab, pkeyVal):      elements(tab):

if (read(tab.has.pkeyVal))  ret :=  $\emptyset$ 
  return false;           for ( let pkeyVal of Vals )
write(tab.has.pkeyVal, true)  if (read(tab.has.pkeyVal))
return true;                ret := ret  $\cup$  {pkeyVal}
                           return ret;

```

Figure 4.8: Manipulating set variables using key-value pairs.

DELETE removes the primary key value associated to this row from the set tab , and UPDATE writes to the keys corresponding to the columns that need to be updated. In the case of UPDATE, we assume that the values of the variables in \vec{x} are obtained from a valuation γ (this valuation would be maintained by the operational semantics of the underlying Key-Value store). INSERT adds a new primary key value to the set variable tab (the call to `add` checks whether this value is unique) and then writes to the keys representing columns of this new row.

MANIPULATING SET VARIABLES Based on the standard representation of a set using its characteristic function, we implement each set variable tab using a set of keys $tab.has.pkeyVal$, one for each value $pkeyVal \in \text{Vals}$. These keys are associated with Boolean values, indicating whether $pkeyVal$ is contained in tab . In a concrete implementation, this set of keys need not be fixed a-priori, but can grow during the execution with every new instance of an INSERT. Figure 4.8 lists the implementations of `add/elements`, which are self-explanatory (`remove` is analogous).

4.5 IMPLEMENTATION

We implemented MonkeyDB⁶ to support an interface common to most storage systems. Operations can be either key-value (KV) updates (to access data as a KV map) or SQL queries (to access data as a relational database). MonkeyDB supports transactions as well; a transaction can include multiple operations. Figure 4.9 shows the architecture of MonkeyDB. A client can connect to MonkeyDB over a TCP connection, as is standard for SQL databases⁷. This offers a plug-and-play experience when using standard frameworks such as JDBC [**jdbc**]. Client applications can also use MonkeyDB as a library in order to directly invoke the storage APIs, or interact with it via HTTP requests, with JSON payloads.

MonkeyDB contains a SQL-To-KV compiler that parses an input query⁸, builds its Abstract Syntax Tree (AST) and then applies the rewriting steps described in Section 4.4 to produce an equivalent sequence of KV API calls (`read()` and `write()`). It uses a hashing routine (`hash`) to generate unique keys corresponding to each cell in a table. For instance, in order to insert a value v for a column c in a particular row with primary key value $pkeyVal$, of a table tab , we invoke `write(hash(tab, pkeyVal, c), v)`. We currently support only a subset of the standard SQL operators. For instance, nested queries or join operators are unsupported; these can be added in the future with more engineering effort.

MonkeyDB schedules transactions from different sessions one after the other using a single global lock. Internally, it maintains execution state as a history consisting of a set of transaction logs, write-read relations and a partial session order (as discussed in §??). On a `read()`, MonkeyDB first collects a set of possible writes present in transaction log that can potentially form write-read (read-from) relationships, and then invokes the consistency checker (Figure 4.9) to confirm validity under the chosen isolation level. Finally, it randomly returns one of the values associated with valid writes. A user can optionally instruct MonkeyDB to only select from the set of *latest* valid write per session. This option helps limit weak behaviors for certain reads.

The implementation of our consistency checker is based on prior work [**DBLP:journals/pacmpl/BiswasE19**]. It maintains the write-read relation as a graph, and detects cycles (isolation-level violations) using DFS traversals on the graph. The consistency checker is an independent and pluggable module: we have one for Read Committed and one for Causal Consistency, and more can be added in the future.

4.6 EVALUATION: MICROBENCHMARKS

We consider a set of micro-benchmarks inspired from real-world applications (§4.6.1) and evaluate the number of test iterations required to fail an invalid assertion (§4.6.2). We also measure the *coverage* of weak behaviors provided by MonkeyDB (§4.6.3). Each of these applications were implemented based on their specifications described in prior work; they all use MonkeyDB as a library, via its KV interface.

⁶We plan to make MonkeyDB available open-source soon.

⁷We support the MySQL client-server protocol using <https://github.com/jonhoo/msql-srv>.

⁸We use <https://github.com/ballista-compute/sqlparser-rs>

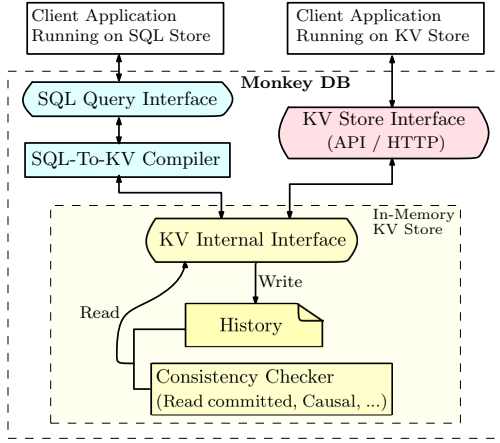


Figure 4.9: Architecture of MonkeyDB

```

// Get user's tweets
Timeline(user u) {
  Begin()
  key = "tweets:" + u.id
  T = read(key)
  Commit()
  return sortByTime(T)
}

// Get following users' tweets
NewsFeed(user u) {
  Begin()
  FW = read("following:" + u.id)
  NF = {}
  foreach v ∈ FW:
    T = read("tweets:" + v.id)
    NF = NF ∪ T
  Commit()
  return sortByTime(NF)
}

```

Figure 4.10: Example operations of the Twitter app

4.6.1 APPLICATIONS

TWITTER [twissandra] This is based on a social-networking application that allows users to create a new account, follow, unfollow, tweet, browse the newsfeed (tweets from users you follow) and the timeline of any particular user. Figure 4.10 shows the pseudo code for two operations.

A user can access twitter from multiple clients (sessions), which could lead to unexpected behavior under weak isolation levels. Consider the following scenario with two users, A and B where user A is accessing twitter from two different sessions, S_1 and S_2 . User A views the timeline of user B from one session ($S_1: \text{Timeline}(B)$) and decides to follow B through another session ($S_2: \text{Follow}(A, B)$). Now when user A visits their timeline or newsfeed ($S_2: \text{NewsFeed}(A)$), they expect to see all the tweets of B that were visible via `Timeline` in session S_1 . But under weak isolation levels, this does not always hold true and there could be missing tweets.

SHOPPING CART [DBLP:conf/pldi/Sivaramakrishnan15] This application allows a user to add, remove and change quantity of items from different sessions. It also allows the user to view all items present in the shopping cart. The pseudo code and an unexpected behavior under weak isolation levels were discussed in §4.1, Figure 4.1.

Application	Assertion	Avg. time to fail	
		(Iters)	(sec)
Stack	Element popped more than once	3.7	0.02
Courseware	Course registration overflow	10.6	0.09
Courseware	Removed course registration	57.5	0.52
Shopping	Item reappears after deletion	20.2	0.14
Twitter	Missing tweets in feed	6.3	0.03

Table 4.1: Assertions checking results in microbenchmarks

COURSEWARE [DBLP:conf/esop/NairP020] This is an application for managing students and courses, allowing students to register, de-register and enroll for courses. Courses can also be created or deleted. Courseware maintains the current status of students (registered, de-registered), courses (active, deleted) as well as enrollments. Enrollment can contain only registered students and active courses, subject to the capacity of the course.

Under weak isolation, it is possible that two different students, when trying to enroll concurrently, will both succeed even though only one spot was left in the course. Another example that breaks the application is when a student is trying to register for a course that is being concurrently removed: once the course is removed, no student should be seen as enrolled in that course.

TREIBER STACK [DBLP:conf/cav/NagarMJ20] Treiber stack is a concurrent stack data structure that uses compare-and-swap (CAS) instructions instead of locks for synchronization. This algorithm was ported to operate on a kv-store in prior work [DBLP:conf/cav/NagarMJ20] and we use that implementation. Essentially, the stack contents are placed in a kv-store, instead of using an in-memory linked data structure. Each row in the store contains a pair consisting of the stack element and the key of the next row down in the stack. A designated key “head” stores the key of the top of the stack. CAS is implemented as a transaction, but the pop and push operations do not use transactions, i.e., each read/write/CAS is its own transaction.

When two different clients try to pop from the stack concurrently, under serializability, each pop would return a unique value, assuming that each pushed value is unique. However, under causal consistency, concurrent pops can return the same value.

4.6.2 ASSERTION CHECKING

We ran the above applications with MonkeyDB to find out if assertions, capturing unexpected behavior, were violated under causal consistency. Table 4.1 summarizes the results. For each application, we used 3 client threads and 3 operations per thread. We ran each test with MonkeyDB for a total of 10,000 times; we refer to a run as an iteration. We report the average number of iterations (Iters) before an assertion failed, and the corresponding time taken (sec). All the assertions were violated within 58 iterations, in half a second or less. In contrast, running with an actual database almost never produces an assertion violation.

4.6.3 COVERAGE

The previous section only checked for a particular set of assertions. As an additional measure of test robustness, we count the number of distinct *client-observable states* generated by a test. A

client-observable state, for an execution, is the vector of values returned by read operations. For instance, a stack’s state is defined by return values of `pop` operations; a shopping cart’s state is defined by the return value of `GetCart` and so on.

For this experiment, we randomly generated test harnesses; each harness spawns multiple threads that each execute a sequence of operations. In order to compute the absolute maximum of possible states, we had to limit the size of the tests: either 2 or 3 threads, each choosing between 2 to 4 operations.

Note that any program that concurrently executes operations against a store has two main sources of non-determinism: the first is the interleaving of operations (i.e., the order in which operations are submitted to the store) and second is the choice of read-from (i.e., the value returned by the store under its configured isolation level). MonkeyDB only controls the latter; it is up to the application to control the former. There are many tools that systematically enumerate interleavings (such as `CHESS` [DBLP:conf/pldi/MusuvathiQ08], `COYOTE` [coyote-web]), but we use a simple trick instead to avoid imposing any burden on the application: we included an option in MonkeyDB to deliberately add a small random delay (sleep between 0 to 4 ms) before each transaction begins. This option was sufficient in our experiments, as we show next.

We also implemented a special setup using the `COYOTE` tool [coyote-web] to enumerate all sources of non-determinism, interleavings as well as read-from, in order to explore the entire state space of a test. We use this to compute the total number of states. Figure 4.11 shows the number of distinct states observed under different isolation levels, averaged across multiple (50) test harnesses. For each of serializability and causal consistency, we show the max (as computed by `COYOTE`) and versions with and without the delay option in MonkeyDB.

Each of these graphs show similar trends: the number of states with causal consistency are much higher than with serializability. Thus, testing with a store that is unable to generate weak behaviors will likely be ineffective. Furthermore, the “delay” versions of MonkeyDB are able to approach the maximum within a few thousand attempts, implying that MonkeyDB’s strategy of per-read randomness is effective for providing coverage to the application.

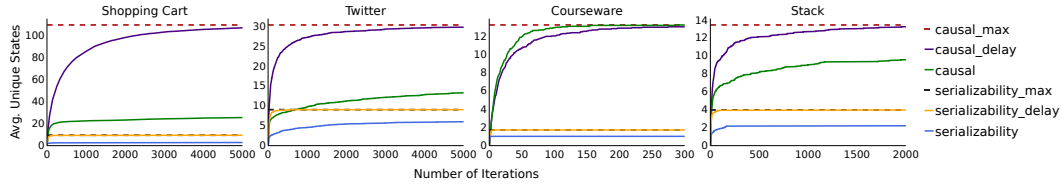


Figure 4.11: State coverage obtained with MonkeyDB for various microbenchmarks

4.7 EVALUATION: OLTP WORKLOADS

OLTPBench [DBLP:journals/pvldb/DifallahPCC13] is a benchmark suite of representative OLTP workloads for relational databases. We picked a subset of OLTPBench for which we had reasonable assertions. Table 4.2 lists basic information such as the number of database tables, the number of static transactions, how many of them are read-only, and the number of different assertions corresponding to system invariants for testing the benchmark. We modified OLTPBench by rewrit-

Benchmark	#Tables	#Txns	#Read-only	#Assertions
TPC-C	9	5	2	12
SmallBank	3	6	1	1
Voter	3	1	0	1
Wikipedia	12	5	2	3

Table 4.2: OLTP benchmarks tested with MonkeyDB

ing SQL join and aggregation operators into equivalent application-level loops, following a similar strategy as prior work [DBLP:journals/pacmpl/RahmaniNDJ19]. Except for this change, we ran OLTPBench unmodified.

For TPC-C, we obtained a set of 12 invariants from its specification document [tpcc-spec]. For all other benchmarks, we manually identified invariants that the application should satisfy. We asserted these invariants by issuing a read-only transaction to MonkeyDB at the end of the execution of the benchmark. None of the assertions fail under serializability; they are indeed invariants under serializability.⁹ When using weaker isolation, we configured MonkeyDB to use latest reads only (§4.5) for the assertion-checking transactions in order to isolate the weak behavior to only the application.

We ran each benchmark 100 times and report, for each assertion, the number of runs in which it was violated. Note that OLTPBench runs in two phases. The first is a loading phase that consists of a big initial transaction to populate tables with data, and then the execution phase issues multiple concurrent transactions. With the goal of testing correctness, we *turn down* the scale factor to generate a small load and limit the execution phase time to ten seconds with just two or three sessions. A smaller test setup has the advantage of making debugging easier. With MonkeyDB, there is no need to generate large workloads.

TPC-C TPC-C emulates a wholesale supplier transactional system that delivers orders for a warehouse company. This benchmark deals with customers, payments, orders, warehouses, deliveries, etc. We configured OLTPBench to issue a higher proportion ($> 85\%$) of update transactions, compared to read-only ones. Further, we considered a small input workload constituting of one warehouse, two districts per warehouse and three customers per district.

TPC-C has twelve assertions (A1 to A12) that check for consistency between the database tables. For example, A12 checks: for any customer, the sum of delivered order-line amounts must be equal to the sum of balance amount and YTD (Year-To-Date) payment amount of that customer.

Figure 4.12 shows the percentage of test runs in which an assertion failed. It shows that all the twelve assertions are violated under Read Committed isolation level. In fact, 9 out of the 12 assertions are violated in more than 60% of the test runs. In case of causal, all assertions are violated with three sessions, except for A4 and A11. We manually inspected TPC-C and we believe that both these assertions are valid under causal consistency. For instance, A4 checks for consistency between two tables, both of which are only updated within the same transaction, thus causal consistency is enough to preserve consistency between them.

⁹We initially observed two assertions failing under serializability. Upon analyzing the code, we identified that the behavior is due to a bug in OLTPBench that we have reported to the authors (link omitted).

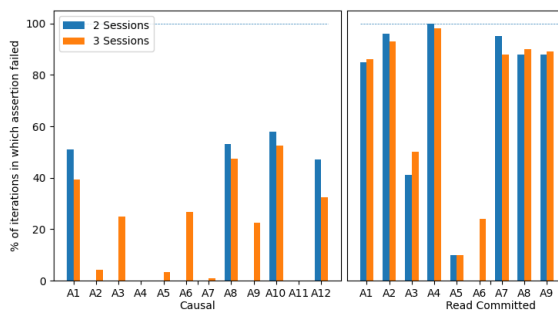


Figure 4.12: Assertion checking: TPC-C

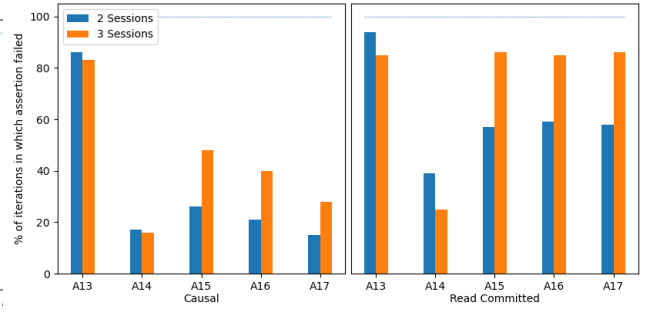


Figure 4.13: Assertion checking: SmallBank, Voter, and Wikipedia

These results demonstrate the effectiveness of MonkeyDB in breaking (invalid) assertions. Running with MySQL, under read committed, was unable to violate any assertion except for two (A10 and A12), even when increasing the number of sessions to 10. We used the same time limit of 10 seconds for the execution phase. We note that MySQL is much faster than MonkeyDB and ends up processing up to $50\times$ more transactions in the same time limit, yet is unable to violate most assertions. Prior work [DBLP:journals/pacmpl/RahmaniNDJ19] attempted a more sophisticated test setup where TPC-C was executed on a Cassandra cluster, while running Jepsen [jepsen] for fault injection. This setup also was unable to violate all assertions, even when running without transactions, and on a weaker isolation level than read committed. Only six assertions were violated with 10 sessions, eight assertions with 50 sessions, and ten assertions with 100 sessions. With MonkeyDB, there is no need to set up a cluster, use fault injection or generate large workloads that can make debugging very difficult.

SMALLBANK, VOTER, AND WIKIPEDIA SmallBank is a standard financial banking system, dealing with customers, saving and checking accounts, money transfers, etc. Voter emulates the voting system of a television show and allows users to vote for their favorite contestants. Wikipedia is based on the popular online encyclopedia. It deals with a complex database schema involving page revisions, page views, user accounts, logging, etc. It allows users to edit its pages and maintains a history of page edits and user actions.

We identified a set of five assertions, A13 to A17, that should be satisfied by these systems. For SmallBank, we check if the total money in the bank remains the same while it is transferred from one account to another (A13). Voter requires that the number of votes by a user is limited to a fixed threshold (A14). For Wikipedia, we check if for a given user and for a given page, the number of edits recorded in the user information, history, and logging tables are consistent (A15-A17). As before, we consider small work loads: (1) five customers for SmallBank, (2) one user for Voter, and (3) two pages and two users for Wikipedia.

Figure 4.13 shows the results. MonkeyDB detected that all the assertions are invalid under the chosen isolation levels. Under causal, MonkeyDB could break an assertion in 26.7% (geo-mean) runs given 2 sessions and in 37.2% (geo-mean) runs given 3 sessions. Under read committed, the corresponding numbers are 56.1% and 65.4% for 2 and 3 sessions, respectively.

4.8 TESTING SQL DATABASE

There have been several directions of work addressing the correctness of database-backed applications. We directly build upon one line of work concerned with the logical formalization of isolation levels [ansi, DBLP:conf/icde/AdyaLO00, DBLP:conf/sigmod/BerensonBGM0095, DBLP:conf/concur/Cerone00, DBLP:journals/pacmpl/BiswasE19]. Our work relies on the axiomatic definitions of isolation levels, as given in [DBLP:journals/pacmpl/BiswasE19], which also investigated the problem of checking whether a given history satisfies a certain isolation level. Our kv-store implementation relies on these algorithms to check the validity of the values returned by read operations. Working with a logical formalization allowed us to avoid implementing an actual database with replication or sophisticated synchronization.

Another line of work concentrates on the problem of finding “anomalies”: behaviors that are not possible under serializability. This is typically done via a static analysis of the application code that builds a static dependency graph that over-approximates the data dependencies in all possible executions of the application [DBLP:journals/jacm/CeroneG18, DBLP:journals/jacm/CeroneG18, DBLP:conf/concur/0002G16, DBLP:journals/tods/FeketeLOOS05, DBLP:conf/vldb/JorwekarFRS07, DBLP:conf/sigmod/WarszawskiB17, DBLP:journals/pvldb/GanRRB020]. Anomalies with respect to a given isolation level then corresponds to a particular class of cycles in this graph. Static dependency graphs turn out to be highly imprecise in representing feasible executions, leading to false positives. Another source of false positives is that an anomaly might not be a bug because the application may already be designed to handle the non-serializable behavior [DBLP:conf/pldi/BrutschyD0V18, DBLP:journals/pvldb/GanRRB020]. Recent work has tried to address these issues by using more precise logical encodings of the application, e.g. [DBLP:conf/popl/BrutschyD0V17, DBLP:conf/pldi/BrutschyD0V18] or by using user-guided heuristics [DBLP:journals/pvldb/GanRRB020].

Another approach consists of modeling the application logic and the isolation level in first-order logic and relying on SMT solvers to search for anomalies [DBLP:journals/pacmpl/KakiESJ18, DBLP:conf/concur/NagarJ18, burcu-netys], or defining specialized reductions to assertion checking [DBLP:conf/concur/BeillahiBE19, DBLP:conf/cav/BeillahiBE19]. The CLOTHO tool [DBLP:journals/pacmpl/Rah00] for instance, uses a static analysis of the application to generate test cases with plausible anomalies, which are deployed in a concrete testing environment for generating actual executions.

Our approach, based on testing with MonkeyDB, has several practical advantages. There is no need for analyzing application code; we can work with any application. There are no false positives because we directly run the application and check for user-defined assertions, instead of looking for application-agnostic anomalies. The limitation, of course, is the inherent incompleteness of testing.

Several works have looked at the problem of reasoning about the correctness of applications executing under weak isolation and introducing additional synchronization when necessary [DBLP:conf/eurosys/BalegasD00, DBLP:conf/popl/GotsmanYFNS16, DBLP:conf/esop/NairP020, DBLP:conf/usenix/0001LCPRV14]. As in the previous case, our work based on testing has the advantage that it can scale to real sized applications (as opposed to these techniques which are based on static analysis or logical proof arguments), but it cannot prove that an application is correct. Moreover, the issue of repairing applications is orthogonal to our work.

From a technical perspective, our operational semantics based on recording past operations and certain data-flow and control-flow dependencies is similar to recent work on stateless model check-

ing in the context of weak memory models, e.g. [DBLP:journals/pacmpl/Kokologiannakis18, DBLP:conf/tacas/AbdullaAAJLS]. This work, however, does not consider transactions. Furthermore, their focus is on avoiding enumerating equivalent executions, which is beyond the scope of our work (but an interesting direction for future work).

4.9 CONCLUSION

Our goal is to enable developers to test the correctness of their storage-backed applications under weak isolation levels. Such bugs are hard to catch because weak behaviors are rarely generated by real storage systems, but failure to address them can lead to loss of business [DBLP:conf/sigmod/WarszawskiB17]. We present MonkeyDB, an easy-to-use mock storage system for weeding out such bugs. MonkeyDB uses a logical understanding of isolation levels to provide (randomized) coverage of all possible weak behaviors. Our evaluation reveals that using MonkeyDB is very effective at breaking assertions that would otherwise hold under a strong isolation level.

