

Formal Modeling and Analysis of RAMP Transaction Systems*

Si Liu
University of Illinois at
Urbana-Champaign

Jatin Ganhotra
University of Illinois at
Urbana-Champaign

Peter Csaba Ölveczky
University of Oslo

Indranil Gupta
University of Illinois at
Urbana-Champaign

Muntasir Raihan Rahman
University of Illinois at
Urbana-Champaign

José Meseguer
University of Illinois at
Urbana-Champaign

ABSTRACT

To cope with large data sets, distributed data stores partition their data across servers. However, real-world systems usually do not provide useful transactional semantics for operations accessing multiple partitions due to the delays involved in achieving multi-partition consistency. Read Atomic Multi-Partition (RAMP) transactions have recently been proposed as efficient light-weight multi-partition transactions that guarantee read atomicity: either all updates or no updates of a transaction are visible to other transactions. In this paper we formalize RAMP transactions in rewriting logic and perform model checking verification of key properties using the Maude tool. In particular, we develop detailed formal models—and formally analyze—a number of extensions and optimizations of RAMP that are only briefly mentioned by the RAMP developers.

Keywords

multi-partition transactions; formal analysis; rewriting logic

1. INTRODUCTION

The success of cloud computing relies on software systems that store large amounts of data correctly and efficiently. It is hard to satisfy both these requirements in a distributed storage system. While traditional relational databases (like MySQL) offer strong notions of correctness (such as ACID properties) when multiple clients access data (via transactions), these systems are considered too slow for today's

*This work was supported in part by NSF CNS 1409416, NSF CNS 1319527, NSF CCF 0964471, AFOSR/AFRL FA8750-11-2-0084, and gifts from Microsoft and Google.

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

SAC'16, April 4-8, 2016, Pisa, Italy

©2016 ACM. ISBN 978-1-4503-3739-7/16/04...\$15.00

DOI: <http://dx.doi.org/10.1145/2851613.2851838>

workloads and applications. As a result, a new generation of storage systems called NoSQL (“Not only SQL”) have emerged, and are already a multi-billion dollar industry [15]. The NoSQL era was kicked off by systems from Google [6] and Facebook [11], and includes open-source systems that are wildly popular in industry today [17, 5, 10]. These systems offer weak notions of correctness, such as “eventual consistency,” while enabling operations on stored data to execute orders of magnitude faster than relational databases.

At the heart of this dichotomy is the CAP theorem [4, 14], which says that it is impossible to have both efficiency (low latency) and strong consistency (correctness) in a distributed storage system.¹ Several efforts have recently emerged to bridge this gap. One of the most promising is the RAMP system of transactions proposed by Bailis *et al.* [1, 2]. RAMP allows clients to execute transactions (like in relational databases) in NoSQL-like storage systems. It offers a correctness property called “Read Atomicity” (RA) which ensures that a given transaction’s updates are either all visible or not visible at all, to other transactions.² RAMP also allows data to be partitioned across multiple partitions.

The RAMP paper [2] presented experimental results that demonstrated its efficiency properties of low latency and high throughput; it also gave “hand proofs” that RAMP provides read atomicity. The paper [2] also briefly mentions a number of extensions and optimizations, such as faster commit and one-phase writes; however, no details or correctness proofs for these extensions are given. There is therefore a clear need for formally specifying and analyzing RAMP and its proposed extensions. A main contribution of our work is to develop detailed models of, and to formally analyze, a number of those briefly-mentioned extensions and variants of RAMP. Providing correctness in the protocol and its extensions is a critical step towards making RAMP a production-capable system. Without such analysis, developers and deployers might enable or disable particular extensions at will, without knowing the correctness ramifications.

In this paper, we use the expressive Maude formal specification language [7] to provide the first executable formal model of RAMP and a number of its extensions. This approach al-

¹When there are partitions—which are endemic in distributed systems in today’s Internet—present.

²While RA is stronger than eventual consistency, it is still not equivalent to ACID as it is not serializable.

lows us to analyze the correctness properties of RAMP in a fully-automated manner, rather than relying only on hand proofs. However, explicit-state model checking can only analyze the system from *single* initial configurations, which provide limited coverage. One of the new contributions of this paper is therefore a general technique in Maude for model checking a protocol for *all possible* initial configurations up to certain bounds. We use this technique to analyze some of the (unproved) conjectures in the original RAMP paper—in particular, how some of the extensions and optimizations of RAMP affect its behavior. We model and analyze the effect of the following RAMP building blocks: fast commit, one-phase write, and two-phase commit.

We explore not only the “strong” notion of read atomicity correctness but also a weaker consistency model called “read-your-writes” (whereby a client is assured of being able to read at least its own latest writes). For instance, we find that two-phase commit is a necessary building block for ensuring read atomicity, and that for one-phase writes even the weak correctness notion of read-your-writes cannot be guaranteed!

The rest of this paper is structured as follows. Section 2 gives some background on RAMP and Maude. Section 3 presents our formal model of RAMP and its extensions. Section 4 explains how the different consistency levels and other key properties of RAMP can be formally specified as reachability properties, which can then be analyzed using Maude. Finally, Section 5 discusses related work and Section 6 gives some concluding remarks.

2. PRELIMINARIES

2.1 RAMP Transactions

To deal with large amounts of data, distributed databases *partition* their data across multiple servers. However, many systems do not provide useful transactional semantics for operations accessing multiple partitions, since the latency needed to ensure correct multi-partition transactional access is often high. Therefore, trade-offs that combine efficiency with weaker transactional guarantees are needed.

In [2], Bailis *et al.* propose a new isolation model, called *read atomic* (RA) isolation, and *Read Atomic Multi-Partition* (RAMP) transactions, that together combine efficient multi-partition operations and partial fault tolerance with some transactional guarantee: either all or none of a transaction’s updates are visible to other transactions. For example, if *A* and *B* become “friends” in a transaction, then other transactions should *not* see that *A* is a friend of *B* but that *B* is not a friend of *A*; either both or no relationships are visible.

RAMP writers attach metadata to each write and the reads use this metadata to get the correct version. There are three versions of RAMP: RAMP-Fast, RAMP-Small, and RAMP-Hybrid. The write protocols in these algorithms only differ in the amount of attached metadata. To guarantee that all partitions perform a transaction successfully or that none do, RAMP performs two-phase writes by using the two-phase commit protocol (2PC): In the *prepare* phase, each timestamped write is sent to its partition, which adds the write to its local database. In the *commit* phase, each partition updates an index which contains the highest-

timestamped committed version of each item. The RAMP algorithms described in [2] only deal with read-only and write-only transactions. Although we have formalized and analyzed both RAMP-Fast and RAMP-Small and their extensions and optimizations, we only deal with RAMP-Fast in this paper due to space limitations, and refer to our longer report [12] for a treatment of RAMP-Small.

In RAMP-Fast, read transactions first fetch the highest-timestamped committed *version* of each requested data item from the corresponding partition, and then decide if they have missed any version that has been prepared but not yet committed. The timestamp and metadata from each version read produces a mapping from items to timestamps that represent the highest-timestamped write for each transaction, appearing in the first-round read set. If the reader has a lower timestamp version than indicated in the mapping for that item, a second-round read will be issued to fetch the missing version. Once all the missing versions have been fetched, the client can return the resulting set of versions, which include both the first-round reads as well as any missing versions fetched in the second round of reads. The detailed specification of RAMP-Fast is given as follows in [2]:

Algorithm 1 RAMP-Fast

Server-side Data Structures

- 1: *versions*: set of versions (*item, value, timestamp ts_v , metadata md*)
- 2: *latestCommit*[*i*]: last committed timestamp for item *i*

Server-side Methods

- 3: **procedure** PREPARE(*v* : version)
- 4: *versions.add*(*v*)
- 5: **return**
- 6: **procedure** COMMIT(*ts_c* : timestamp)
- 7: $I_{ts} \leftarrow \{w.item \mid w \in versions \wedge w.ts_v = ts_c\}$
- 8: $\forall i \in I_{ts}, latestCommit[i] \leftarrow \max\{latestCommit[i], ts_c\}$
- 9: **procedure** GET(*i* : item, *ts_{req}* : timestamp)
- 10: **if** $ts_{req} = \emptyset$ **then**
- 11: **return** $v \in versions : v.item = i \wedge v.ts_v = latestCommit[item]$
- 12: **else**
- 13: **return** $v \in versions : v.item = i \wedge v.ts_v = ts_{req}$

Client-side Methods

- 14: **procedure** PUT_ALL(*W* : set of (*item, value*))
 - 15: $ts_{rx} \leftarrow$ generate new timestamp
 - 16: $I_{rx} \leftarrow$ set of items in *W*
 - 17: **parallel-for** $\langle i, v \rangle \in W$
 - 18: $v \leftarrow \langle item = i, value = v, ts_v = ts_{rx}, md = (I_{rx} - \{i\}) \rangle$
 - 19: invoke PREPARE(*v*) on respective server (i.e., partition)
 - 20: **parallel-for** server *s* : *s* contains an item in *W*
 - 21: invoke COMMIT(ts_{rx}) on *s*
 - 22: **procedure** GET_ALL(*I* : set of items)
 - 23: $ret \leftarrow \{\}$
 - 24: **parallel-for** *i* $\in I$
 - 25: $ret[i] \leftarrow$ GET(*i*, \emptyset)
 - 26: $v_{latest} \leftarrow \{\}$ (default value: -1)
 - 27: **for** response *r* $\in ret$ **do**
 - 28: **for** $i_{tx} \in r.md$ **do**
 - 29: $v_{latest}[i_{tx}] \leftarrow \max\{v_{latest}[i_{tx}], r.ts_v\}$
 - 30: **parallel-for** item *i* $\in I$
 - 31: **if** $v_{latest}[i] > ret[i].ts_v$ **then**
 - 32: $ret[i] \leftarrow$ GET(*i*, $v_{latest}[i]$)
 - 33: **return** *ret*
-

Extensions of RAMP. The paper [2] also briefly discusses a number of extensions and optimizations of the RAMP algorithms, but without giving details. These include:

- RAMP with one-phase writes. If a client does not wish to read her own writes, writes only require one *prepare* phase, since the client can execute the *commit* phase asynchronously.
- RAMP with faster commit. If a server returns a version with the timestamp fresher than the highest committed version of the item, then the server can mark the version as committed. This allows faster updates and thus fewer round trip time delays.
- RAMP without two-phase commit. We have also experimented with the possibility of decoupling two-phase commit. (This is *not* a variant proposed in [2].)

2.2 Rewriting Logic and Maude

Maude [7] is a rewriting-logic-based formal specification language and high-performance simulation and model checking tool for concurrent systems. Maude specifications are executable, and the tool provides a variety of formal analysis methods, including simulation, reachability analysis, and linear temporal logic (LTL) model checking.

Specification. A Maude module specifies a *rewrite theory* $(\Sigma, E \cup A, R)$, where:

- Σ is an algebraic *signature*; that is, a set of *sorts*, *subsorts*, and *function symbols*.
- $(\Sigma, E \cup A)$ is a *membership equational logic theory* [7], with E a set of possibly conditional equations and membership axioms, and A a set of equational axioms such as associativity, commutativity, and identity, so that equational deduction is performed *modulo* the axioms A . The theory $(\Sigma, E \cup A)$ specifies the system's state space as an algebraic data type.
- R is a collection of *labeled conditional rewrite rules* $[l] : t \rightarrow t' \text{ if } \textit{cond}$, specifying the system's local transitions.

We briefly summarize the syntax of Maude and refer to [7] for more details. Operators are introduced with the `op` keyword: `op f : s1 ... sn -> s`. They can have user-definable syntax, with underbars ‘`_`’ marking the argument positions. Equations and rewrite rules are introduced with, respectively, keywords `eq`, or `ceq` for conditional equations, and `rl` and `cr1`. The mathematical variables in such statements are declared with the keywords `var` and `vars`, or can be introduced on the fly, in which case they have the form `var : sort`. An equation $f(t_1, \dots, t_n) = t$ with the `owise` (“otherwise”) attribute can be applied to a subterm $f(\dots)$ only if no other equation with left-hand side $f(u_1, \dots, u_n)$ can be applied.

A *class* declaration `class C | att1 : s1, ..., attn : sn` declares a class C with attributes att_1 to att_n of sorts s_1 to s_n . An *object instance* of class C is represented as a term $\langle O : C \mid att_1 : val_1, \dots, att_n : val_n \rangle$, where O , of sort `Objid`, is the object's *identifier*, and where val_1 to val_n are the current

values of the attributes att_1 to att_n . A *message* is a term of sort `Msg`. The state is a term of the sort `Configuration`, and has the structure of a *multiset* made up of objects and messages. The dynamic behavior of a system is axiomatized by specifying each of its transition patterns by a rewrite rule. For example, the rule

```

rl [1] : m(0,w)
        < 0 : C | a1 : x, a2 : 0', a3 : z >
=>
        < 0 : C | a1 : x + w, a2 : 0', a3 : z >
        m'(0',x) .

```

defines a family of transitions in which a message m , with parameters 0 and w , is read and consumed by an object 0 of class C , the attribute $a1$ of the object 0 is changed to $x + w$, and a new message $m'(0',x)$ is generated. Attributes whose values do not change and do not affect the next state, such as $a3$, need not be mentioned in a rule.

Reachability Analysis. Maude's *search* command

```
(search [n] t0 =>* pattern such that condition .)
```

uses a breadth-first strategy to search for at most n states that are reachable from the initial state t_0 in zero or more steps, match the search *pattern*, and satisfy the search *condition*. To search for *final* states (i.e., states that cannot be further rewritten) matching *pattern* and satisfying *condition*, the arrow `=>!` is used instead of `=>*`.

3. MODELING RAMP IN MAUDE

This section presents formal models of RAMP and its variants. Since we focus on analyzing the correctness of RAMP, and not its performance, our model is an *untimed* one, so that all possible interleavings are analyzed by Maude model checking. The executable Maude specifications are available at <https://sites.google.com/site/siliunobi/ramp>.

We show a few rewrite rules of a few models for illustration purposes, and refer to our longer report [12] for more details.

3.1 Data Types, Objects, and Messages

We formalize RAMP in an object-oriented style, where the distributed state consists of a number of `Partition` objects modeling partitions of the database, a number of `Client` objects modeling clients that issue transactions, and a number of messages traveling between the objects.

Although a transaction (request) is a *list* of operations, read-only and write-only transactions can be seen as *sets* of operations under the assumption that transactions do not have conditional reads and that a transaction does not write the same item multiple times.

Basic Data Types. A *version* of a data item is modeled as a 4-tuple `version(item, value, timestamp, metadata)` consisting of the data item, its value, and the version's timestamp and metadata. A timestamp is modeled as a pair `timestamp(id, sqn)` consisting of a client's identifier id and a sequence number sqn that together identify a writing transaction. Metadata (which in RAMP-Fast are write sets) are

modeled as a set of data items, denoting, for each item, the other items that are written in the same transaction.

A transaction is modeled as a ‘,’-separated set of write operations `write(id, item, value)` and read operations `read(id, item)`, where `id` denotes the identity of the operation. The sort `TrList` denotes lists of transactions.

Objects and Messages. A *client* issues transactions and collects responses for analysis purposes. A transaction from a client is issued when its preceding transaction has been committed. Concurrent transactions can be modeled by multiple transactions issued by different clients.

A client is modeled as an object instance of the class `Client` with the following attributes: a list of transactions it wants to issue (`transac`); the sequence number that together with the client identifier determine timestamps (`sqn`); a flag indicating whether a transaction has been committed (`trflag`); several sets of operations buffering intermediate information (`pendingOps`, denoting pending reads/writes, `pendingPrep`, denoting pending writes in the *prepare* phase, `1stGets`, denoting pending first-round reads); a mapping from each item to its latest committed timestamp from a client’s perspective (`latest`); and the returned value and operations for each transaction (`result`):

```
class Client |
  transac : TrList, sqn : Int, trflag : Bool,
  pendingOps : OpsInfo, pendingPrep : Set{Nat},
  1stGets : Set{Nat}, latest : Latest, result : Results .
```

A *partition* stores a set of versions for its data items. We model a partition as an object of class `Partition` with attributes: `versions`, denoting a set of versions of the items in the partition, and `latestCommit`, denoting the timestamp of the last commit of each item:

```
class Partition |
  versions : Versions, latestCommit : Latest .
```

For example, the state fragment

```
< c1 : Client |
  transac : (read(3,x), read(4,y)), sqn : 3,
  trflag : true, pendingOps : {1,2}, pendingPrep : {2},
  1stGets : empty, latest : empty, result : nil >
< p1 : Partition |
  versions : version(x,v,timestamp(c2,3),y),
            version(y,w,timestamp(c1,1),x),
  latestCommit : (x |-> timestamp(c2,3) ,
                 y |-> timestamp(c1,1)) >
```

models a setting where partition `p1` holds a version for each of the items `x` and `y`, with respective values `v` and `w`, and associated timestamps `timestamp(c2,3)` and `timestamp(c1,1)`. Client `c1` intends to issue a transaction of two reads, reading items `x` and `y`; `trflag` is `true`, meaning that one of `c1`’s transactions is being executed, and `pendingOps` has two pending operations 1 and 2. We can tell from `pendingPrep` that the ongoing transaction is a write-only transaction with a write 1 that has already committed. The remaining three attributes are `empty` or `nil`, since they are only updated when a read-only transaction is issued.

Messages travel between clients and partitions, and have the form `msg msgContent from sender to receiver`, where `msgContent` is either `prepare(id, version)` (placing a write on its partition), `commit(id, timestamp)` (marking versions as committed using the *timestamp* generated by a client), `prepared(id)` (reply to `prepare` for write `id`), `committed(id)` (reply to `commit`), `get(id, item, timestamp)` (fetching the highest-timestamped committed version or any missing version for an *item* by *timestamp*), `1st-response(id, version)` (returned *version* for first-round `get` for read `id`), or `2nd-response(id, version)` (returned version for second-round `get` for read `id`), where `id` is the operation’s identifier.

3.2 Formalizing RAMP-Fast

This section formalizes the dynamic behaviors of RAMP-Fast using rewrite rules. We also refer to the corresponding lines of code in the description in Section 2.1. We show some of the rules for writes, and refer to [12] for more details.

Starting a New Transaction (Lines 14–19 for writes and lines 22–26 for reads). Whenever a client is not executing a transaction (the attribute `trflag` has the value `false` and the buffers storing intermediate information are empty) and there is a pending transaction in the client’s transaction list `transac` (the pattern³ `TR TRL` denotes a list with a single transaction `TR` followed by a (possibly empty) list of `TRL` of the other remaining transactions), the client starts issuing the transaction `TR` by sending a `prepare` message for each write in a write-only transaction, or by sending a `get` message for each read in a read-only transaction. The function `genOps` generates these messages.

The client will then wait for intermediate response messages from the partitions by setting `pendingOps` to the operations in `TR`, `pendingPrep` to the writes (`w`) in `TR` if `TR` is a write-only transaction, `1stGets` to the reads (`r`) in `TR` if `TR` is a read-only transaction, and `latest` to default timestamps for each item requested in `TR`’s reads, respectively. Furthermore, to record the final committed versions, the client also adds to `result` the default version `rs(TR)` (with `null` value and timestamp, and `empty` metadata) for each read in `TR`. Finally, the client sets `trflag` to `true` to indicate that it is currently executing a transaction:

```
r1 [init-transaction] :
  < 0 : Client | transac : TR TRL, sqn : SQN, trflag : false,
                pendingOps : empty, pendingPrep : empty,
                1stGets : empty, latest : VL, result : RS >
=>
  < 0 : Client | transac : TRL, trflag : true,
                pendingOps : rw(TR), pendingPrep : w(TR),
                1stGets : r(TR), latest : vl(TR, VL),
                result : RS rs(TR) >
  genOps(TR, SQN, 0) .
```

Receiving Prepare Messages (Lines 3–5). When a partition `0` receives a `prepare` message with identifier `ID` from the client `0’` for a version `v`, it adds `v` to its local storage `versions` and sends a `prepared` message back to the client:

```
r1 [receive-prepare] :
```

³We do not show variable declarations, but follow the Maude convention that variables are written in capital letters.

```

msg prepare(ID, version(X, V, timestamp(O', SQN), MD))
  from O' to O
< O : Partition | versions : VS >
=>
< O : Partition | versions :
  version(X, V, timestamp(O', SQN), MD), VS >
msg prepared(ID) from O to O' .

```

Receiving Prepared Messages (Lines 20–21). When a client receives a `prepared` message for the current write ID from a partition `O'`, it deletes ID from the set `pendingPrep`. If the resulting set is empty, meaning that all `prepared` messages have been received, the client starts committing the transaction using the function `startCommit`, which generates a `commit` message for each write.

```

cr1 [receive-prepared] :
  msg prepared(ID) from O' to O
  < O : Client | pendingPrep : NS, pendingOps : OI, sqn : SQN >
=>
  < O : Client | pendingPrep : NS' >
  (if NS' == empty then startCommit(OI, SQN, O) else none fi)
  if NS' := delete(ID, NS) .

```

Receiving Commit Messages (Lines 6–8). When a partition receives a `commit` message with timestamp `timestamp(O', SQN)`, it invokes the function `cmt` to update the latest commit timestamp in the set `latestCommit` with the fresher timestamp of the incoming one and the local one, provided those two can be matched; it then notifies the client to commit the write:

```

r1 [receive-commit] :
  msg commit(ID, timestamp(O', SQN)) from O' to O
  < O : Partition | versions : VS, latestCommit : LC >
=>
  < O : Partition | latestCommit :
    cmt(LC, VS, timestamp(O', SQN)) >
  msg committed(ID) from O to O' .

```

Receiving Committed Message. When a client receives a `committed` message, it removes the committed read/write ID from its set `pendingOps`. If all pending reads/writes have been committed (`pendingOps` is empty), the client gets ready to issue its next transaction by setting `trflag` to `false` and increasing the sequence number:

```

r1 [receive-committed] :
  msg committed(ID) from O' to O
  < O : Client | pendingOps : OI >
=>
  < O : Client | pendingOps : remove(ID, OI) > .

r1 [commit-transaction] :
  < O : Client | trflag : true, pendingOps : empty,
    sqn : SQN >
=>
  < O : Client | trflag : false, sqn : SQN + 1 > .

```

3.3 Formalizing Variants of RAMP

We have formalized the following optimizations of RAMP: RAMP without two-phase commit, RAMP with faster commit, and RAMP with one-phase writes. We explain the

modifications for one-phase writes, and refer to [12] for details on the other modifications.

A client that does not wish to read her writes can “return after issuing its `prepare` round. The client can subsequently execute the `commit` phase asynchronously.” Specifically, after collecting all `prepared` messages (no more `pendingPreps`), a client starts to execute the commit phase by invoking `startCommit` to generate `commit` messages, and kicks off its next transaction, unless it has no more transactions to issue (the rule for this final case not shown):

```

cr1 [receive-prepared-1pw-1] :
  msg prepared(ID) from O' to O
  < O : Client | pendingPrep : (ID , NS) >
=>
  < O : Client | pendingPrep : NS >   if NS != empty .

r1 [receive-prepared-1pw-next-trans] :
  msg prepared(ID) from O' to O
  < O : Client | pendingPrep : ID, pendingOps : OI,
    latest : VL, sqn : SQN,
    transac : TR TRS, result : RS >
=>
  < O : Client | pendingOps : rw(TR), pendingPrep : w(TR),
    1stGets : r(TR), latest : vl(TR, VL),
    result : RS rs(TR), sqn : SQN + 1 >
  startCommit(OI, SQN, O)
  genOps(TR, SQN + 1, O) .

```

4. FORMAL ANALYSIS

We use reachability analysis to analyze whether RAMP and its variants satisfy the properties in [2]. Explicit-state model checkers like Maude are quite expressive but only analyze the system from a *single* initial configuration. To increase coverage, we want to model check RAMP for *all possible* configurations up to certain bounds, such as k operations and j clients. Despite the wealth of Maude applications, we are not aware of any such comprehensive model checking in Maude. Section 4.1 therefore presents a general technique in Maude for model checking a system from a range of different initial configurations. Section 4.2 formalizes the properties that RAMP transactions should satisfy and analyzes them for *all* configurations with four operations and two clients, as well as for a number of configurations with six operations.

4.1 Model Checking Many Initial States

The idea behind model checking many initial configurations is to introduce a new operator `init`, have a *one-step* rewrite `init(parameters) → t0` for *any* possible initial configuration t_0 , and start the analysis from `init(parameters)`. However, there are two things to take into account:

1. The analysis must take into account the additional rewrite step before the actual analysis of all behaviors from a concrete initial state begins. A search command (`search t0 =>* pattern .`), which searches for states reachable in zero or more steps from t_0 , must be replaced by (`init(parameters) =>+ pattern .`), which searches for states reachable in *one* or more steps. Likewise, in temporal logic model checking, an LTL formula φ should be replaced by the formula $\bigcirc\varphi$ (which means that φ holds in the next state).

- The property to analyze may depend on the particular initial state. When generating multiple initial states and model checking them in one command, it might be necessary to record (parts of) the initial state, and carry this record in the state throughout the analysis.

To generate all possible initial states, we declare a new sort denoting *sets* of configurations:

```
sort ConfigSet . subsort Configuration < ConfigSet .
op empty : -> ConfigSet .
op _;- : ConfigSet ConfigSet -> ConfigSet
      [assoc comm id: empty] .
```

We assume that there is a function

```
op initAux : s1 ... sn -> ConfigSet .
```

such that `initAux(params, params')` generates all possible initial states, and add the following rewrite rule to our model:

```
var C : Configuration . var CS : ConfigSet .
crl [init] : init(params) => C
           if C ; CS := initAux(params, params') .
```

If the state needs to carry a record $f(t_0)$ of the initial state t_0 , we use the following rule instead:

```
crl [init] :
  init(params) => C < r : Record | record : f(C) >
  if C ; CS := initAux(params, params') .
```

where `Record` is a new class which stores selected data from the chosen initial state throughout the computations.

For RAMP, `init` (see [12] for its definition) has the parameters: `#ops`, the total number of (read or write) operations to be performed; `#clients`, the number of clients; and `dataItems`, the set of data items in the system, with each partition storing one item. We also store in the `record` object the list of transactions to be issued by each client.

One of 2764 initial states defined by `init(4,2,(x,y))` is

```
< 1 : Client | transac : (read(1,x) , read(2,x)) , sqn : 1 ,
  1stGets : empty , latest : empty , pendingOps : empty ,
  pendingPrep : empty , result : nil , trflag : false >
< 2 : Client | transac : (write(4,y,4) write(3,x,3)) , ... >
< x : Partition | latestCommit : x |-> timestamp(0,0) ,
  versions : version(x,null,timestamp(0,0) , empty) >
< y : Partition | latestCommit : y |-> timestamp(0,0) ,
  versions : version(y,null,timestamp(0,0) , empty) >
< 100 : Record | record : 1 |->(read(1,x) , read(2,x)) ,
  2 |-> write(4,y,4) write(3,x,3) >
```

where client 1 has one transaction with two reads, and client 2 has two transactions, each having one write operation.

4.2 Analyzing the Correctness Properties

We formalize the correctness requirements of RAMP as reachability properties and analyze them using Maude.

During the execution of a transaction with multiple reads, one read will be committed before the other, leading to intermediate states where key properties do not hold. Since

RAMP is terminating if each client issues a finite number of transactions, we therefore analyze most properties on final states, when all transactions are committed.

We have performed our analysis from `init(4,2,(x,y))`, which means that we consider all possible initial configurations with a total of 4 operations, 2 clients, and 2 data items. There are 2764 such initial configurations. Each analysis command took about 20 seconds to execute on a 2.9 GHz Intel 4-Core i7-3520M CPU with 3.7 GB memory.

For 6 total operations, there are too many initial states for model checking them all in one shot. However, if we ignore: transactions in which the same item is read twice or written twice, single-operation transactions (neither a single read transaction nor a single write transaction can lead to fractured reads), and symmetric scenarios, we should be left with 6 different scenarios different to analyze (see [12]). Although we believe that those 6 initial states “cover” all possible scenarios with 6 operations, this remains to be verified. The analysis of each scenario took about a second.

Read Atomic Isolation. The main correctness property, *read atomic isolation*, that RAMP should satisfy is defined as follows in [2]:

“A system provides *Read Atomic* isolation if it prevents fractured reads anomalies and also prevents transactions from reading uncommitted, aborted, or intermediate data.”

where *fractured reads* are described as follows:

“A transaction T_j exhibits *fractured reads* if transaction T_i writes version x_m and y_n (in any order, with x possibly but not necessarily equal to y), T_j reads version x_m and version y_k , and $k < n$.”

We analyze this property by searching for a reachable *final* state where the property does *not* hold; i.e., a state where one of the clients’ committed results, RES1 or RES2, has a fractured read:

```
(search [1] init(4,2,(x,y)) =>!
  C:Configuration
  < r:Oid : Record | record : (01:Oid |-> TL1:TrList)
  (02:Oid |-> TL2:TrList) >
  < 01:Oid : Client | result : RES1:Result >
  < 02:Oid : Client | result : RES2:Result >
  such that fracRead(RES1:Result, TL1:TrList, TL2:TrList)
  or fracRead(RES2:Result, TL2:TrList, TL1:TrList) .)
```

where the function `fracRead` checks whether there are fractured reads on each client’s committed result, based on the initial transactions of all clients (see [12] for details).

Our analysis results are consistent with the analytic and predicted results in [2]: all versions of RAMP, except the one *without* two-phase commit, satisfy read atomicity.

Companions Present. Another property that is verified in [2] is the following invariant:

“If a version x_i is referenced by *lastCommit* (that is, $lastCommit[x] = i$), then each of x_i ’s sibling

versions⁴ are present in *versions* on their respective partitions.”

This invariant can be analyzed by searching for a state that does not satisfy the property:

```
(search [1] init(4,2,(x,y)) =>+ C:Configuration
  such that not comp-present(C:Configuration) .)
```

where **comp-present** holds if and only if for each committed item, the sibling version of its associated version is present in the other partition (see [12] for details).

Our analysis shows that all versions, except the ones without 2PC, satisfy this property.

Synchronization Independence. This property is described as follows in [2]:

“Synchronization Independence ensures that one client’s transactions cannot cause another client’s to block and that, if a client can contact the partition responsible for each item in its transaction, the transaction will eventually commit (or abort of its own volition).”

We check the stronger property that all transactions are eventually always committed on all clients’ sides: since we assume that there are no uncommitted or aborted transactions, a client-side method can return only if the corresponding partition-side methods have returned. The following function **syn-indep** holds in a state if and only if there are no pending transactions (**transac** is **nil** and all buffers are **empty**) and no ongoing transactions (**trflag** is **false**):

```
op syn-indep : Configuration -> Bool .
eq syn-indep(
  < O1 : Client | transac : nil, trflag : false,
    pendingOps : empty, pendingPrep : empty,
    1stGets : empty >
  < O2 : Client | transac : nil, trflag : false,
    pendingOps : empty, pendingPrep : empty,
    1stGets : empty > REST)
  = true .
eq syn-indep(SYSTEM) = false [owise] .
```

Again, we analyze the property by searching for a *final* state which does *not* satisfy **syn-indep**. The property is satisfied by all versions of RAMP in our analysis.

Read Your Writes. This property says that a client’s writes are visible to her subsequent reads. It can be analyzed by searching for a final state in which a client’s recorded reads contain a read that read the client’s write which was *not* the write immediately preceding the read:

```
(search [1] init(4,2,(x,y)) =>+
  C:Configuration
  < r:0id : Record | record : (C1:0id |-> TL1:TrList)
    (C2:0id |-> TL2:TrList) >
  < C1:0id : Client | result : RES1:Result >
```

⁴We call the set of versions produced by a (write-only) transaction *sibling versions*.

```
< C2:0id : Client | result : RES2:Result >
  such that tooOldRead(RES1:Result, TL1:TrList)
    or tooOldRead(RES2:Result, TL2:TrList) .)
```

A client has read a too old value if either it reads a *null* value, or a value it has written earlier, *and* it has a (later) write transaction writing the item.

The function **tooOldRead** returns **true** if a read ID has fetched either: (i) a value **V1** of an earlier write **ID1** although the client had a later write **ID2** preceding the read ID (first equation); or (ii) the initial value **null** even though the client had an earlier write **ID1** that has not written **null** (second equation). In the first equation below, the **result** value contains a read ID **|-> version(X,V1,TS,MD)**, and the transactions that the client should execute match the pattern

```
TRL (OPS1 , write(ID1,X,V1) , OPS1')
TRL' (OPS2 , write(ID2,X,V2) , OPS2')
TRL'' (OPS , read(ID,X) , OPS') TRL'''
```

where the variables **TRL** denote *lists of transactions*, and the variables **OPS** denote parts of a single transaction (sets of operations). In particular, in the client’s original transaction list, there are two transactions (**OPS1 , write(ID1,X,V1) , OPS1'**) and (**OPS2 , write(ID2,X,V2) , OPS2'**) both writing **X** before a transaction (**OPS , read(ID,X) , OPS'**) reads **X**. It is clear that if **read(ID,X)** reads the value written by **write(ID1,X,V1)**, then the client has not read its own (latest) writes. The second equation takes care of the case when the read operation reads **null** even though it itself had an earlier transaction writing **X**. If none of these equations apply, the **owise** equation returns **false**:

```
op tooOldRead : Results TrList -> Bool .
eq tooOldRead((RS (ID |-> version(X,V1,TS,MD) , R) RS2),
  (TRL (OPS1 , write(ID1,X,V1) , OPS1')
  TRL' (OPS2 , write(ID2,X,V2) , OPS2')
  TRL'' (OPS , read(ID,X) , OPS') TRL''')) = true .
eq tooOldRead((RS1 (ID |-> version(X,null,TS,MD) , R) RS2),
  (TRL (OPS1 , write(ID1,X,V1) , OPS2)
  TRL' (OPS , read(ID,X) , OPS') TRL''')) = true .
eq tooOldRead(RS, TRL) = false [owise] .
```

Our analysis shows that only RAMP with one-phase writes does not satisfy the property. The counterexample obtained by analyzing the initial state where a client has transactions [*write(x); write(y)*] [*read(x); read(y)*] shows that the read operations both return **null**. The reason is that one-phase writes do not forbid the client to start the subsequent read-only transaction before its last write-only transaction has committed (by **commit** messages) on the partitions.

Summary. We have analyzed all four key properties of our seven versions of RAMP. Our results agree with the proved and conjectured results in [2]: All versions satisfy the properties, except that

- RAMP without 2PC does not satisfy read atomicity, “companions present” and read-your-writes; and
- RAMP with one-phase writes does not satisfy read-your-writes.

5. RELATED WORK

Despite the importance of distributed data stores for cloud computing, we are not aware of much work on formalizing and verifying such systems. The recent paper [16] describes experiences at Amazon Web Services using TLA+ to formalize and model check Amazon’s scalable high-performance replicated NoSQL data store DynamoDB. The authors report that TLA+ model checking “[found] bugs in system designs that cannot be found through any other techniques we know of” and that “formal methods are routinely applied to the design of complex real-world software, including public cloud services,” at Amazon. Maude was used in [8, 9] to define a formal model of Google’s widely-replicated data store Megastore, which ensures serializability for certain classes of transactions, and to develop an extension of Megastore. These models were simulated for QoS estimation and model checked for functional correctness. In a similar vein, [13] presents a formal model of the popular distributed key-value store Cassandra, and formally analyzes different consistency properties using Maude. The authors in [3] reduce the problem of verifying *eventual consistency* of optimistic data replication systems in large-scale networks to reachability and model checking problems. The main difference between all the above-mentioned work and this paper is that we formalize and analyze a different consistency property (read atomic isolation) than the eventual consistency and serializability properties analyzed in the related work, and that we do so for several alternative designs of RAMP, a system that, to the best of our knowledge, has not been formally modeled and analyzed before.

6. CONCLUDING REMARKS

Today’s distributed storage systems are seeing a convergence of traditional “strong consistency” databases and the new generation of “fast but eventually consistent” NoSQL databases. In this paper, we have analyzed a promising bridge system, called RAMP [2], which seeks to provide strong consistency for client transactions while still offering fast performance. While the original RAMP paper included hand proofs only for the basic RAMP algorithms, we adopted a model checking approach where we: (1) first developed fully-executable formal models of many RAMP variants, extensions, and optimizations using Maude; (2) formally analyzed these models to confirm the original correctness properties (within some constraints); (3) used our models to evaluate RAMP’s extensions and optimizations (which the original RAMP paper did not do); and (4) used our models to evaluate the critical dependence of RAMP’s correctness properties on its building blocks. In addition, we have presented a general technique in Maude for model checking systems for all possible initial configurations up to certain bounds. Symmetry reduction should be added to this method in future work to reduce the number of generated initial states.

Based on our experience, we believe that the formal modeling and model checking approach can be powerful if used by developers of cloud storage systems before, or even alongside, their implementation. This approach allows developers to go beyond mere hand proofs or implementation-based simulations. Once the formal model is written, correctness

properties can be automatically analyzed. Splicing system features in and out becomes easy, and this can be used to isolate the effect of each constituent building block, as well as to evaluate the effect of further optimizations.

7. REFERENCES

- [1] P. Bailis, A. Davidson, A. Fekete, A. Ghodsi, J. M. Hellerstein, and I. Stoica. Highly available transactions: Virtues and limitations. *PVLDB*, 7(3), 2013.
- [2] P. Bailis, A. Fekete, J. M. Hellerstein, A. Ghodsi, and I. Stoica. Scalable atomic visibility with RAMP transactions. In *Proc. SIGMOD’14*. ACM, 2014.
- [3] A. Bouajjani, C. Enea, and J. Hamza. Verifying eventual consistency of optimistic replication systems. In *Proc. POPL’14*. ACM, 2014.
- [4] E. A. Brewer. Towards robust distributed systems (abstract). In *PODC*, page 7, 2000.
- [5] Cassandra. <http://cassandra.apache.org>.
- [6] F. Chang et al. Bigtable: A distributed storage system for structured data. *ACM Trans. Comput. Syst.*, 26(2), 2008.
- [7] M. Clavel, F. Durán, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, and C. L. Talcott. *All About Maude*, volume 4350 of *LNCS*. Springer, 2007.
- [8] J. Grov and P. C. Ölveczky. Formal modeling and analysis of Google’s Megastore in Real-Time Maude. In *Specification, Algebra, and Software*, volume 8373 of *LNCS*. Springer, 2014.
- [9] J. Grov and P. C. Ölveczky. Increasing consistency in multi-site data stores: Megastore-CGC and its formal analysis. In *SEFM*, volume 8702 of *LNCS*. Springer, 2014.
- [10] E. Hewitt. *Cassandra: The Definitive Guide*. O’Reilly Media, 2010.
- [11] A. Lakshman and P. Malik. Cassandra: a decentralized structured storage system. *ACM SIGOPS Operating Systems Review*, 44(2):35–40, 2010.
- [12] S. Liu et al. Formal modeling and analysis of RAMP transaction systems. <https://sites.google.com/site/siliunobi/ramp>, 2015.
- [13] S. Liu, M. R. Rahman, S. Skeirik, I. Gupta, and J. Meseguer. Formal modeling and analysis of Cassandra in Maude. In *Proc. ICFEM’14*, volume 8829 of *LNCS*. Springer, 2014.
- [14] N. Lynch and S. Gilbert. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *ACM SIGACT News*, 33(2):51–59, 2002.
- [15] Market research media, NoSQL market forecast 2013-2018. <http://www.marketresearchmedia.com/?p=568>.
- [16] C. Newcombe, T. Rath, F. Zhang, B. Munteanu, M. Brooker, and M. Deardeuff. How Amazon Web Services uses formal methods. *Communications of the ACM*, 58(4):66–73, 2015.
- [17] Basho Riak. <http://basho.com/riak/>.