



# On the Necessity of Model Checking NoSQL Database Schemas When Building SaaS Applications

Stefanie Scherzinger, Eduardo Cunha de Almeida, Felipe Ickert, Marcos Didonet del Fabro

## ► To cite this version:

Stefanie Scherzinger, Eduardo Cunha de Almeida, Felipe Ickert, Marcos Didonet del Fabro. On the Necessity of Model Checking NoSQL Database Schemas When Building SaaS Applications. International Workshop on Testing the Cloud, Jul 2013, Lugano, Switzerland. ACM, 2013, <10.1145/2489295.2489297>. <hal-01018129>

**HAL Id: hal-01018129**

**<https://hal.inria.fr/hal-01018129>**

Submitted on 3 Jul 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# On the Necessity of Model Checking NoSQL Database Schemas When Building SaaS Applications

Stefanie Scherzinger<sup>1</sup>, Eduardo Cunha de Almeida<sup>2</sup>, Felipe Ickert<sup>2</sup>, and Marcos Didonet Del Fabro<sup>2</sup>

<sup>1</sup>Regensburg University of Applied Sciences, Germany – stefanie.scherzinger@hs-regensburg.de

<sup>2</sup>Federal University of Paraná, Brazil – {eduardo | fvickert | didonet}@inf.ufpr.br

## ABSTRACT

The design of the NoSQL schema has a direct impact on the scalability of web applications. Especially for developers with little experience in NoSQL stores, the risks inherent in poor schema design can be incalculable. Worse yet, the issues will only manifest once the application has been deployed, and the growing user base causes highly concurrent writes. In this paper, we present a model checking approach to reveal scalability bottlenecks in NoSQL schemas. Our approach draws on formal methods from tree automata theory to perform a conservative static analysis on both the schema and the expected write-behavior of users. We demonstrate the impact of schema-inherent bottlenecks for a popular NoSQL store, and show how concurrent writes can ultimately lead to a considerable share of failed transactions.

## Categories and Subject Descriptors

D.2.4 [Software/Program Verification]: [Model checking]; H.2.1 [Logical Design]: [Data models]; H.2.4 [Systems]: [Concurrency, Distributed databases]

## General Terms

Verification

## Keywords

Model Checking, NoSQL Databases

## 1. INTRODUCTION

With Platform-as-a-Service (PaaS) offerings readily available, hosting an application in the cloud has turned into a commodity. The PaaS framework manages the complex deployment of the application across a large-scale hardware infrastructure, and comes with its own hosted databases, either relational or NoSQL.

The learning curves in working with these high-end services seem steep, yet this impression can be treacherous.

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](mailto:Permissions@acm.org).

Copyright is held by the author/owner(s). Publication rights licensed to ACM.

TTC'13, July 15, 2013, Lugano, Switzerland  
ACM 978-1-4503-2162-4/13/07  
<http://dx.doi.org/10.1145/2489295.2489297>

While easy to get started with, NoSQL systems bear their own pitfalls when it comes to writing scalable applications. Unfortunate design choices, made early on, can later manifest in severe performance bottlenecks when the workload increases, and ultimately even prevent scalability. Even worse, these problems arise once the application is up and running and the number of users is growing rapidly.

Such traps are known amongst experienced software engineers. Yet the research community seems to have mainly focused on matters of *physical schema* design, such as investigating data partitioning under transactions [8, 9, 17]. The impact that the *logical schema* design can have on the performance of an application has so far been neglected. The large body of work on relational schema design [1] can only be marginally transferred to NoSQL databases.

In this paper, we point out common pitfalls encountered by novices to NoSQL technologies, particularly when developing against backends related to Megastore [5], (e.g. [3, 4, 12, 14]). These systems co-locate data items that are frequently accessed together, to avoid distributed transactions. Thus, the schema design has an impact on transactional behavior. We stress the need for model checking on the level of the database schema, and make the following contributions.

- We introduce a real-life design challenge to the testing research community. In particular, we draw attention to scalability bottlenecks lurking in the schema design of NoSQL backends.
- We argue that NoSQL stores are still in the process of being better understood by many developers and testers. Most professionals have grown up with relational databases and ACID transactions, and are new to the NoSQL paradigm. This makes poor design choices likely when working with new technologies.
- Our work targets logical schema design, not the physical optimization of NoSQL databases. We are thus catering to the needs of the application developers and testers, rather than those of cloud service providers.
- We provide a method for static, conservative schema analysis that checks for schema-inherent bottlenecks. In doing so, we assume that the schema design and the expected user write behavior are known. Our approach uses methods from tree automata theory.

The remainder of this paper is structured as follows. In Section 2, we discuss a blogging application with alternative schema designs, as well as their impact on scalability. Our model checking approach is presented in Section 3. We experimentally demonstrate the impact of schema design on

application performance in Section 4. Section 5 is dedicated to related work. We then conclude with a summary and an outlook on future work.

## 2. CASE STUDY: MANAGING BLOGS

We model the schema for a blogging site to illustrate the impact of early design decisions on the scalability of a NoSQL-based web application. Our project follows the lead of beacon apps such as the Huffington Post website, which handles 70+ million comments on articles each year. Some ballpark numbers from late 2012 declare that users collectively write up to 25k comments per hour. While an article receives about 40 comments on average, comments are not evenly distributed: The most popular articles count over 4,500 comments, a single video even more than 170k comments [19, 20].

We choose Datastore for persistence, a NoSQL backend provided with Google AppEngine [12, 13, 18] (or an alternative system with a Megastore-style architecture, such as Hypertable [14] or Cassandra [3]). While our contributions are not limited to this particular platform, we stick to Datastore terminology throughout this paper.

Datastore is a key-value store<sup>1</sup> with entities as the basic unit of storage. Entities have a unique *key* and a *value*. Entities are made persistent by calling a simple *put*-operation, and retrieved by key using a *get*-operation. Different from a relational database, the schema is not fixed in advance. Entities need not cohere to any pre-defined structure. This makes Datastore a *schema-less* backend. It is nevertheless plausible to assume that developers maintain some loose notion of a schema, otherwise it becomes too complex to maintain an application. Out of this reason, web applications typically use object mappers [18], which enforce a consistent mapping between the objects in the application code and the persisted entities.

Datastore supports queries in a restricted query language. Due to a weak consistency model, queries may temporarily return stale results. When the application logic requires strong consistency in updating several entities, the entities may be stored as an *entity group*. The system can then physically co-locate these entities, thus allowing for stricter consistency guarantees. Several systems built upon Megastore-like architectures provide this feature (c.f. Section 5).

**Naive design: The article-oriented schema.** Figure 1 shows how entities could be grouped for a blogging application. We store an entity for each of the registered users Alice, Bob, Charlie, Dave, and Eve. Then “users(Alice)” is an entity key. The predicate “users” would then be the *entity kind*, a means of categorizing entities for the purpose of querying. As a key-value store, Datastore maps the key “users(Alice)” to the entity value containing data on Alice’s user account, such as her password and preferences. In the figure, the value is depicted by a document icon, since we do not care for the exact contents.

Users may publish articles. “articles(Alice, a1)” is the key for an article written by Alice. The article entities may be stored with each user, forming a hierarchy of entities, and thus an entity group. Any comments on this article made by other users, such as a comment by user Dave with key “comments(Dave, a1)”, are also stored within this group.

<sup>1</sup>Datastore is actually a *document store*, with structured values. Yet this distinction is of no consequence to our work.

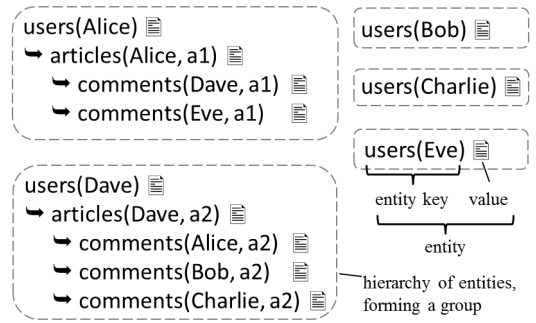


Figure 1: Grouping comments with articles.

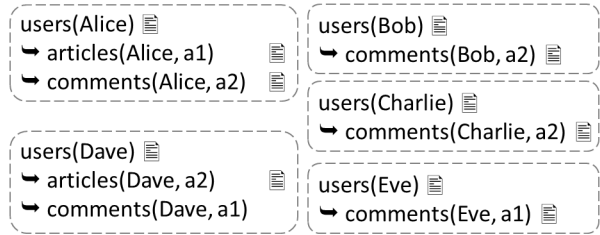


Figure 2: Grouping postings with their authors.

This may seem a natural model, since the conceptual relationships between entities are reflected in the hierarchy. Yes this schema contains a dangerous design issue that functional testing will not reveal: Systems like Datastore are tuned towards supporting massive parallel reads, at the cost of limiting parallel writes. By modeling the relationship between articles and comments in this manner, we force Datastore to co-locate all comments on a particular article in a single group. Yet the write throughput for groups is throttled, causing writes to fail if too many users are writing against a group simultaneously. With Google Datastore, at least one write can be handled against a group per second, while on average, 5–10 concurrent writes can be managed [11]. Due to this limit and an optimistic concurrency control, concurrent writes against the same entity group are likely to fail. Consequently, a heated debate over an article will inevitably lead to a bad user experience.

We take the vivid discussions on this dilemma in developer blogs and forums as a strong indicator that the implications of schema design on the scalability of an application are not easy to grasp for NoSQL novices.

**A safer way: The user-oriented schema.** The safest bet to avoid write contention is to never group any entities together. Yet storing each entity within its very own entity group comes at the cost of losing transactional safety and restricted query capabilities.

An approach that still allows for strong consistency on updates within a group and additional query options is to group all articles and comments with their authors, as depicted in Figure 2. Realistically, a user won’t post several articles or comments per second. Hence, we consider it guaranteed that there will not be any concurrent writes against the same entity group. Thus, the alternative schema is safe from scalability bottlenecks.

### 3. MODEL CHECKING NOSQL SCHEMAS

Our contribution is an analysis that checks whether a schema is guaranteed to be free of concurrent write requests on the granularity of single entity groups.

We make these assumptions: (1) The schema is known, at least on the level of how entity keys are constructed and how entities are assembled inside entity groups. NoSQL systems are per se schema-less. However, any maintainable web application will require some consistent notion of the structure of entity groups. To account for large degrees of freedom, we describe the entity group structure by regular tree grammars, rather than by a relational schema. (2) Since the application code is available, we know the expected write behavior of users. (3) We assume that a single user won't execute any parallel writes against the same entity group. Of course, several users may write concurrently. This is reasonable for interactive web applications, where the application handles user web requests one at-a-time.

**Safe schemas.** Under these assumptions, we can statically check whether a given schema is free of concurrent write requests against single entity groups. We then say the schema is *safe*. Otherwise, we conservatively assume that concurrent writes might occur, making the schema *unsafe*.

EXAMPLE 1. The schema from Figure 1 is unsafe when users write comments. During a heated debate, comments against the same article manifest in writes against the same entity group. In contrast, the schema from Figure 2 is safe since each user writes against a separate entity group.  $\square$

**Outline of our approach.** We encode the schema as a regular tree language. Our formalism closely resembles DTDs, which define the valid structure of XML documents [7]. We then derive a tree automaton  $\mathcal{A}_s$  (the “schema” automaton) which accepts only trees that would be valid in an entity group. We capture what is known about user write behavior in the form of *update expressions*. They are the basis for a further tree automaton  $\mathcal{A}_w$  (the “writes” automaton) which accepts a tree only if all updates in an entity group occur on behalf of the same user. Of course, each tree in this language may be updated by a different user.

Leveraging automata theory, we check whether the language accepted by  $\mathcal{A}_s$  is contained in the language accepted by  $\mathcal{A}_w$ . If this holds, then for any entity group in this schema, we can guarantee that all writes against this group are done by the same user. Containment can be checked in EXPTIME [7] when nondeterministic automata are involved, and is an established procedure in XML processing.

There are two cases in which we classify a schema as unsafe. Either the containment check fails, or the automaton  $\mathcal{A}_w$  is not deterministic to begin with. In either case, it is possible for two users to write against the same entity group.

#### 3.1 Encoding Schemas

An entity group contains a tree of entities, possibly consisting of a single node. We next formalize the schema as a language of unranked trees. We abstract even further and focus on the entity keys, ignoring the values.

**Formalizing the schema.** We encode the tree structures of entity groups by an *extended regular tree grammar* [7]. These are regular grammars where the right-hand-side of a production contains a regular expression, allowing for the specification of languages over unranked, labeled trees.

Let the finite symbol set  $\Sigma$  contain all valid entity keys. This is the grammar alphabet<sup>2</sup>. For each symbol  $\sigma$  in  $\Sigma$  denoting a node label, we define a production of the form  $\sigma ::= \rho$  where  $\rho$  is a regular expression over symbols from  $\Sigma$ . The right-hand-side of the production defines a pattern over the children of this node. We further identify a designated start production for an imaginary root node, defining which entities may be at the top-level in an entity group (a common trick in the literature on XML).

The following example defines the article-oriented schema from Figure 1. Out of convenience, we introduce this shorthand for regular expressions:  $|_{a \in \{a_1, \dots, a_n\}} := a_1 | \dots | a_n$ .

EXAMPLE 2. We start with constructing the alphabet  $\Sigma$ . Let  $U$  be a finite set of user identifiers, and let  $A$  be a finite set of article identifiers. Then the alphabet is defined as

$$\Sigma = \{users(u) \mid u \in U\} \cup \{articles(u, a) \mid u \in U, a \in A\} \\ \cup \{comments(u, a) \mid u \in U, a \in A\}.$$

Next are the grammar productions, defined for each user  $u$  in  $U$  and article identifier  $a$  in  $A$ :

$$root ::= |_{u' \in U} users(u') \quad (1)$$

$$users(u) ::= (|_{a' \in A} articles(u, a'))^* \quad (2)$$

$$articles(u, a) ::= (|_{u' \in U} comments(u', a))^* \quad (3)$$

$$comments(u, a) ::= \epsilon \quad (4)$$

Production (1) specifies that a user entity makes the topmost node in an entity group. (2) states that the user node has article nodes as children. (While a user won't write the same article twice, our grammar would allow this. Yet as long as the grammar is more general than the actual constraints, it is safe to simplify this way.) Production (3) specifies that an article node may have comments as children. The last production declares comments as leaf nodes. All entity groups in Figure 1 correspond to trees defined by this grammar.  $\square$

The grammar in the example above allows only trees that are at most three levels deep. We could have just as well defined recursive productions, so that comments can be posted against comments. In the mere interest of a comprehensible example, we prefer to avoid this additional complexity.

On a second note, the grammar above does not enforce unique entity keys. Yet as long as the grammar is more general than the application, this is permissible in the interest of a more convenient specification. In the worst case, an over-generalization could yield false positives.

**Automata construction.** Now that we have a way to encode the schema, we need a means for checking whether a given tree is part of this language. This is done by hedge automata [7]. We construct a hedge automaton which processes trees bottom-up, from the leaves towards the root. This automaton may even be nondeterministic, since we can always transform it into a deterministic automaton [7].

A hedge automaton  $\mathcal{A} = (\Sigma, Q, Q_f, \Delta)$  has an alphabet  $\Sigma$ , a set of states  $Q$ , final states  $Q_f$ , and state transition rules  $\Delta$ . In our construction, we use the same alphabet as the grammar. For each symbol  $\sigma$  in the alphabet  $\Sigma$ , we define a corresponding state  $q_\sigma$ . For each symbol  $\phi$  occurring on the

<sup>2</sup>In favor of a compact representation, we do not explicitly distinguish between grammar terminals and nonterminals. This follows the custom of DTDs, which are also instances of extended regular tree grammars [7].

right-hand-side of the start production  $root ::= \rho$ , we make the corresponding state  $q_\phi$  final.

For each grammar production  $\sigma ::= \rho$  with  $\sigma \in \Sigma$ , we define a state transition rule  $q_\sigma \leftarrow \sigma(\rho')$ , where  $\rho'$  is the regular expression obtained from  $\rho$  by replacing each symbol  $\tau \in \Sigma$  by state  $q_\tau$ . Transition rules assign states to nodes bottom-up, and read as follows. When the children of a node labeled  $\sigma$  have all been assigned a state, and the sequence of their assigned states read from left to right matches the regular expression  $\rho'$ , then this node is assigned state  $q_\sigma$ .

EXAMPLE 3. Continuing our example, we define state transition rules for each user  $u \in U$  and article identifier  $a \in A$ :

$$\begin{aligned} q_{users(u)} &\leftarrow users(u) ( (|_{a' \in A} q_{articles(u, a')} )^* ) \\ q_{articles(u, a)} &\leftarrow articles(u, a) ( (|_{u' \in U} q_{comments(u', a)} )^* ) \\ q_{comments(u, a)} &\leftarrow comments(u, a) ( \epsilon ) \end{aligned}$$

The final states are defined as  $Q_f = \{q_{users(u)} \mid u \in U\}$ . Figure 3 shows a run of this automaton on a tree from Figure 1. States are assigned to nodes bottom-up. The nodes to be processed next are highlighted in bold. As the root is assigned a final state, the automaton accepts.  $\square$

### 3.2 Encoding User Write Behavior

As the application code is known, we also know the users' update behavior. Key-value stores don't have powerful update languages. Typically, entities are persisted as key-value pairs and can be again retrieved from storage given their key.

**Formalizing updates.** We formalize a write as  $put^u(k, v)$ , persisting the key-value pair  $(k, v)$  on behalf of user  $u$ . This will overwrite any existing entity with the same key. We formalize a delete as  $del^u(k)$ , deleting the entity with key  $k$  from the store on behalf of user  $u$ . Our syntax is adapted from traditional insert-delete-modify transactions [1].

EXAMPLE 4. Users register and create accounts, publish articles, and may comment on articles. Let  $u$  denote a particular user, then we formalize this user's updates as

$$\begin{aligned} put^u( users(u), v ) \\ put^u( articles(u, a), v ) \\ put^u( comments(u, a), v ) \end{aligned}$$

where  $v$  denotes some value to be stored, and  $a$  denotes some article identifier.  $v$  and  $a$  are treated as unbound variables. If users also delete all data they create, we further define

$$\begin{aligned} del^u( users(u) ) \\ del^u( articles(u, a) ) \\ del^u( comments(u, a) ) \end{aligned}$$

where  $a$  denotes some article identifier.  $\square$

**Automata construction.** From the update expressions, we derive an automaton that accepts a tree only if all its nodes are guaranteed to be updated by a single user. These are the trees that denote safe entity groups, since there are no concurrent updates against the same entity group.

To specify this automaton, we again employ alphabet  $\Sigma$ . We define a state  $q_u$  for each user  $u$ , and make all states final. For each update expression of the form  $put^u(k, e)$  or  $del^u(k)$  we then define a state transition rule  $q_u \leftarrow k( q_u^* )$ .

When processing a tree bottom-up, the automaton assigns state  $q_u$  to a leaf if this node could have been written by

user  $u$ . It will assign this state to an inner node only if all of its children have been assigned state  $q_u$ , and the entity represented by the node could have been written by user  $u$ .

EXAMPLE 5. Continuing Example 4, translating the write-operations yields these transition rules for each user  $u$  in  $U$  and each article identifier  $a$  in  $A$ :

$$\begin{aligned} q_u &\leftarrow users(u) ( q_u^* ) \\ q_u &\leftarrow articles(u, a) ( q_u^* ) \\ q_u &\leftarrow comments(u, a) ( q_u^* ) \end{aligned}$$

The deletes do not contribute any additional rules.  $\square$

If the automaton such constructed is deterministic (i.e. it has no ambiguous transitions), it will only accept trees of entities that have been written by the a single user. In case the automaton should be nondeterministic, model checking classifies the schema as unsafe.

EXAMPLE 6. The automaton just constructed is deterministic. It does not accept all trees from Figure 1, meaning this schema is unsafe: Let us pick out the tree from Figure 3(a). The leaf nodes are assigned the states  $q_{Dave}$  and  $q_{Eve}$ . There is no matching state transition rule for assigning their parent node a state. Since the root node cannot be assigned a final state, the tree is not accepted. Conceptually, when the leaf nodes carry different states, they were written by different users – possibly concurrently. The same automaton accepts all trees from Figure 2.  $\square$

### 3.3 Verifying Schema Safety

With hedge automata, we can check whether a particular tree is part of a given language. Even better, we can also make general statements without having to handle the individual trees of a language (of which there may be infinitely many): Let  $L(\mathcal{A}_s)$  be the language defining the schema, and let  $L(\mathcal{A}_w)$  be the language defining the safe writes (as constructed in the previous section). If  $L(\mathcal{A}_s) \subseteq L(\mathcal{A}_w)$ , then all entities in an entity group are guaranteed to be updated by a single user. Thus, the schema is safe under these update transitions. This check can be done on the level of the hedge automata and is a standard procedure described in [7].

EXAMPLE 7. We model check the schema from Figure 1 with the automata  $\mathcal{A}_s$  and  $\mathcal{A}_w$  from Examples 3 and 5 respectively. Then  $L(\mathcal{A}_s) \not\subseteq L(\mathcal{A}_w)$ . For instance, the tree from Figure 3(a) is accepted by  $\mathcal{A}_s$  but not by  $\mathcal{A}_w$ . Hence this schema is unsafe.  $\square$

EXAMPLE 8. Let us describe the schema from Figure 2 by a grammar, again using alphabet  $\Sigma$ . We define productions for each user  $u$  in  $U$  and article identifier  $a$  in  $A$ :

$$\begin{aligned} root ::= & \quad |_{u' \in U} users(u') \\ users(u) ::= & \quad (|_{a' \in A} articles(u, a'))^* \\ & \quad | (|_{a' \in A} comments(u, a'))^* \\ articles(u, a) ::= & \quad \epsilon \\ comments(u, a) ::= & \quad \epsilon \end{aligned}$$

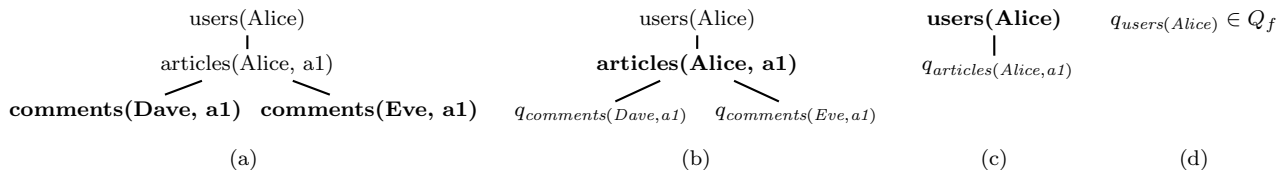


Figure 3: Bottom-up run of the automaton from Example 3.

We then derive the transition rules for automaton  $\mathcal{A}'_s$ :

$$\begin{aligned}
 q_{users(u)} &\leftarrow users(u) \left( (|a' \in A q_{articles(u, a')}|^* \right. \\
 &\quad \left. | (|a' \in A q_{comments(u, a')}|^*)^* \right) \\
 q_{articles(u, a)} &\leftarrow articles(u, a) (\epsilon) \\
 q_{comments(u, a)} &\leftarrow comments(u, a) (\epsilon)
 \end{aligned}$$

As  $L(\mathcal{A}'_s) \subseteq L(\mathcal{A}_w)$ , this schema is safe: For each tree  $t \in L(\mathcal{A}'_s)$  denoting an entity group, since  $t$  is also in  $L(\mathcal{A}_w)$ , this entity group is updated by a single user only.  $\square$

**Discussion.** To the best of our knowledge, this is the first application of tree automata theory for model checking NoSQL schemas. Our approach is conservative in the sense that for all schemas verified as safe, it is guaranteed that no two users will ever write concurrently against the same entity group. However, it may produce false positives. There are schemas where concurrent writes occur without risk, since few users ever write concurrently. For example, sharding global counters effectively distributes the writes over several entity groups, thus increasing write throughput without avoiding concurrent writes [18]. While our decision procedure is conservative, it does detect practical instances of schema-inherent performance bottlenecks.

A known limitation of our method is that the finite set of entity keys must be enumerated in advance. We are currently investigating how we can generalize our approach to work with countably infinite alphabets.

## 4. EXPERIMENTAL EVALUATION

We experimentally tested the scalability of a schema in which comments are stored with articles (“article-oriented”) against a schema where comments are stored with their authors (“user-oriented”). We generate concurrent writes by running a servlet on Google AppEngine [12] which triggers parallel Ajax calls against Datastore [13]. Our application is publicly accessible and may be run from any browser<sup>3</sup>.

In the article-oriented scenario, all comments on an article are stored in the same entity group as the article itself. We created a single article, and ran a series of 20, 100, 500, and 1000 write requests against this single entity group. As Figure 4 shows, this causes between 20% and 35% of the writes to fail. In the user-oriented scenario, all comments are stored in the same entity group as the entity representing their author. Since all writes go against separate entity groups, we achieve a success rate of 100% (see Figure 4).

These observations confirm that the logical schema design directly impacts the scalability of software-as-a-service applications backed by NoSQL data stores. Interestingly, the effect can be confirmed for as few as 20 concurrent writes.

<sup>3</sup><http://vanquish-ufpr-bottleneck.appspot.com>

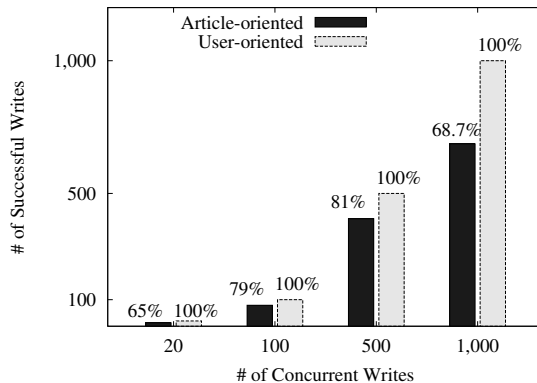


Figure 4: Successful writes per schema.

## 5. RELATED WORK

Schema-inherent bottlenecks may be found in several systems that share a software stack similar to Megastore [5] such as Hypertable [14], HBase [4], Google Datastore [13], BigTable [6], and Cassandra [3]. In this family of systems, transactions are executed on only a subset of columns at row level. Groups of columns may be accessed together, but through a locality constraint indicated by the developers, which is the very principle underlying entity groups.

On the necessity of testing cloud database schemas, the investigation naturally heads towards the model transformation area where test oracles may explore variations of the database schema for validating scalability. Different transformations may be explored, including database fragmentation in order to scale out the database design [16]. The decomposition of a relation into fragments allows transactions and queries to execute concurrently, and is a technique commonly used by high-throughput main memory databases [15].

A possible fragmentation approach based on XML Schema is done by breaking ER-graphs into alternative hierarchies (trees) [10]. Due to the XML nature, this approach accepts multi-valued attributes that may be used to mimic *schema-less* databases. Unlike our approach, the fragmentation customizes the schema towards specific applications without taking performance into account.

In Schism [8], a workload-aware fragmentation approach improves the scalability of distributed databases. The approach creates a graph, where nodes are tuples and edges connecting the nodes are transactions that update them. The graph is fragmented to balance the transactional workload and to minimize costly cross-fragment transactions. While Schism delivers scalability to different types of applications that require distributed databases, its focus is on relational rather than NoSQL databases.

Two further approaches build fragments for transactional

multi key accesses, which may be used in cloud databases. G-Store [9] builds fragments for applications that require multi access to a group of keys. A fragment is built for each non-disjoint group of keys on-demand, where each fragment is kept in a single node (i.e., co-located data collection). Similarly, the One Hop Replication [17] builds fragments for social network databases to keep transactions close to each other, in a cluster of nodes (one-hop away in the best case). Interestingly, these systems can also suffer from concurrency contention when a group of keys becomes popular, such as in our case study.

We refer to [2] for a comprehensive survey on transactions on co-located data in NoSQL data stores.

## 6. SUMMARY AND OUTLOOK

In this paper, we have discussed crucial design decisions in crafting NoSQL schemas for software-as-a-service applications. We have shown that the logical NoSQL schema has a vital impact on the overall scalability of an application. The schema is effectively responsible for physically balancing concurrent write requests. This relationship between schema and application performance is easy to miss for programmers who are new to this domain.

To address this knowledge gap, we perform model checking so that application testers may discover schema-inherent bottlenecks. Our approach requires extracting an abstraction of the schema from a cloud database or the application code, as well as the write requests executed on behalf of users. We draw on formal methods from tree automata theory to perform a conservative static schema analysis. This reveals whether entity groups are guaranteed to be free of concurrent write requests. We show that our approach can detect performance contention that manifests in high rates of failed transactions. Based on this feedback, we were able to manually devise an alternative schema that scales up to highly concurrent write requests on Google Datastore.

Next, we plan to investigate the automated recommendation of alternative schemas. Ideally, such a *schema advisor* for NoSQL databases would suggest alternative designs which show better performance under workloads with concurrent writes.

## 7. ACKNOWLEDGEMENTS

This work was supported by the CNPQ grant 481715/2011-8, and Fundação Araucária grants 17591 and 20874.

## 8. REFERENCES

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [2] D. Agrawal, S. Das, and A. El Abbadi. *Data Management in the Cloud: Challenges and Opportunities*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2012.
- [3] Apache Cassandra. Available on: <http://cassandra.apache.org/>.
- [4] Apache HBase. Available on: <http://hbase.apache.org/>.
- [5] J. Baker, C. Bond, J. C. Corbett, J. Furman, et al. “Megastore: Providing Scalable, Highly Available Storage for Interactive Services”. In *Proc. CIDR*, pages 223–234, 2011.
- [6] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, et al. “Bigtable: A distributed storage system for structured data”. In *Proc. OSDI*, pages 15–15, 2006.
- [7] H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. *Tree Automata Techniques and Applications*. 2007.
- [8] C. Curino, Y. Zhang, E. P. C. Jones, and S. Madden. “Schism: a Workload-Driven Approach to Database Replication and Partitioning”. *PVLDB*, 3(1):48–57, 2010.
- [9] S. Das, D. Agrawal, and A. El Abbadi. “G-Store: A scalable data store for transactional multi key access in the cloud”. In *Proc. SoCC*, pages 163–174, 2010.
- [10] R. Elmasri, Y.-C. Wu, B. Hojabri, C. Li, and J. Fu. “Conceptual Modeling for Customized XML Schemas”. In *Proc. ER*, pages 429–443, 2002.
- [11] A. Fuller and M. Wilder. “More 9s Please: Under The Covers of the High Replication Datastore”. Google I/O Conference, presentation available on: <http://www.google.com/events/io/2011/>, May 2011.
- [12] Google AppEngine. Available on: <https://developers.google.com/appengine/>, June 2013.
- [13] Google AppEngine Datastore. Available on: <https://developers.google.com/appengine/docs/java/datastore/>, June 2013.
- [14] Hypertable Inc. Available on: <http://hypertable.org/>.
- [15] E. P. C. Jones, D. J. Abadi, and S. Madden. “Low overhead concurrency control for partitioned main memory databases”. In *Proc. SIGMOD*, pages 603–614, 2010.
- [16] M. T. Özsu and P. Valduriez. *Principles of Distributed Database Systems, Third Edition*. Springer, 2011.
- [17] J. M. Pujol, G. Siganos, V. Erramilli, and P. Rodriguez. “Scaling online social networks without pains”. In *Proc. NetDB*, 2009.
- [18] D. Sanderson. *Programming Google App Engine*. O’Reilly, 2013.
- [19] J. Sonderman. “How the Huffington Post handles 70+ million comments a year”. Article available on: <http://www.poynter.org>, October 2012.
- [20] The Huffington Post. Available on: <http://www.huffingtonpost.com/theblog/index/>, June 2013.