

# Leveraging Metadata in NoSQL Storage Systems

Ala' Alkhalidi, Indranil Gupta, Vijayanth Raghavan, Mainak Ghosh

*Department of Computer Science*

*University of Illinois, Urbana Champaign*

*Email: {aalkhal2, indy, vraghvn2, mghosh4}@illinois.edu*

**Abstract**—NoSQL systems have grown in popularity for storing big data because these systems offer high availability, i.e., operations with high throughput and low latency. However, metadata in these systems are handled today in ad-hoc ways. We present Wasef, a system that treats metadata in a NoSQL database system, as first-class citizens. Metadata may include information such as: operational history for portions of a database table (e.g., columns), placement information for ranges of keys, and operational logs for data items (key-value pairs). Wasef allows the NoSQL system to store and query this metadata efficiently. We integrate Wasef into Apache Cassandra, one of the most popular key-value stores. We then implement three important uses cases in Cassandra: dropping columns in a flexible manner, verifying data durability during migrational operations such as node decommissioning, and maintaining data provenance. Our experimental evaluation uses AWS EC2 instances and YCSB workloads. Our results show that Wasef: i) scales well with the size of the data and the metadata; ii) affects throughput minimally by only 9%, and iii) affects operational latencies by only 3%.

## 1. Introduction

With the advent of NoSQL stores, large corpuses of data can now be stored in a highly-available manner. Access to this stored data is typically via CRUD operations, i.e., Create, Read, Update, and Delete. NoSQL storage systems provide high throughput and low latency for such operations.

In NoSQL systems such as Apache Cassandra [1], MongoDB [2], Voldemort [3], Bigtable [4], and DynamoDB [5], data is organized into tables, somewhat akin to tables in relational databases. For instance Cassandra calls these tables as “column families”, while MongoDB calls them as “collections”. Each table consists of a set of rows, where each row is a key-value pair or equivalently a data item. Each row is identified by a unique key. Unlike relational databases, NoSQL systems allow schema-free tables so that a data item could have a variable set of columns (i.e., attributes). Access to these data items is allowed via CRUD operations, either using the primary key or other attributes of the data items. Many applications using NoSQL storage

systems do not require such complex operations as joins, and thus they are not widely implemented.

While NoSQL systems are generally more efficient than relational databases, their ease of management remains about as cumbersome as that in traditional systems. A system administrator has to grapple with system logs by parsing flat files that store these operations. Hence implementing system features that deal with metadata is cumbersome, time-consuming and results in ad-hoc designs. Data provenance, which keeps track of ownership and derivation of data items, is usually not supported. During infrastructure changes such as node decommissioning, the administrator has to verify by hand the durability of the the data being migrated.

We argue that such information needs to be collected, stored, accessed, and updated in a first-class manner. We call this information as *metadata*. For the purposes of NoSQL systems, we define metadata as essential information about a data item, a table, or the entire storage system, but excluding the data stored in the data items themselves. This includes structural metadata that is relevant to the way tables are organized, administrative metadata used to manage system resources, and descriptive data about individual data items.

We present Wasef<sup>1</sup>, a metadata system intended for NoSQL data stores. Wasef functions as a component of the data store and leverages the underlying NoSQL functionalities to deliver its services. Wasef has to address three major challenges. First, it must collect metadata without imposing too much overhead on the foreground CRUD operations arriving from clients. Second, it must allow an administrator or a user to specify (via clean APIs) which metadata is collected and how to use it. Finally, Wasef must scale with: i) the size of the cluster, ii) the size of the data being stored, iii) the rate of incoming CRUD operations, and iv) the size of the metadata itself.

Our work makes the following contributions:

- We present the design and architecture of Wasef, a metadata management system for NoSQL storage systems.
- We implement the W-Cassandra system, a key-value store consisting of Wasef integrated into Apache Cassandra 1.2.
- We implement three important use cases in W-Cassandra:

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

1. Arabic word for “Descriptor.”

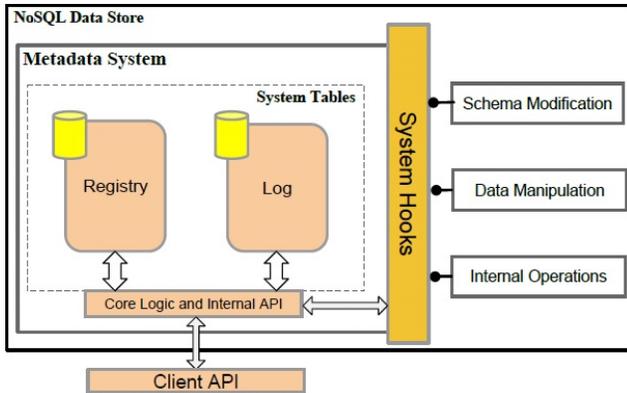


Figure 1: The architecture of Wasef.

- 1) A flexible implementation of the column drop operation, thus addressing a major JIRA (bug) issue in Cassandra 1.2.
  - 2) Durability verification for node decommissioning.
  - 3) A data provenance feature.
- We evaluate W-Cassandra on the AWS cloud [6], and using realistic YCSB [7] workloads. Our evaluation shows that Wasef: i) scales well with the cluster size, data size, operation rate, and metadata size; ii) affects throughput minimally by only 9%, iii) affects operational latencies by only 3%.

The rest of the paper is organized as follows. Section 2 provides an overview of Wasef and its API, and we expand on the details in Section 3. Section 4 describes the design and implementation of the use-case scenarios, and evaluation is presented in Section 5. We describe related work in Section 6, and conclude in Section 7.

## 2. System Design

In this section we first lay out our design principles (Section 2.1). We then describe our architecture (Section 2.2), workflow (Section 2.3) and API (Section 2.4).

### 2.1. Design Principles

Wasef’s design is based on four guiding principles:

- 1) *Modularity and integration with the existing functionality*: The metadata system should integrate with the underlying infrastructure in a modular fashion. It should not affect existing NoSQL system APIs, functionality, or performance.
- 2) *Flexible granularity of collected metadata*: The design should be flexible to collect and store metadata about objects and operations of different kinds and at different granularities (e.g., data items vs. tables). Such metadata includes (but is not fundamentally limited to) the time and description of performed operations, object names, ownership information, and column information.
- 3) *Accessibility of metadata by internal and external clients*: Metadata needs to be accessible by both external clients

(e.g., for data provenance) as well as servers internal to the cluster (e.g., for management operations such as dropping of columns). We provide this via flexible APIs to collect, access, and manipulate metadata.

- 4) *Minimal collection of the metadata*: Due to the enormous size of the data and operations handled by NoSQL data stores, the continuous collection of metadata about every operation might impose a large overhead on the system. To avoid this, Wasef allows the administrator to configure metadata collection for only a selected set of operations.

### 2.2. Architectural Components

Wasef consists of five major components (Figure 1):

- *Registry*: The Registry is a table for registering objects for whom metadata will be collected. Each Registry entry is identified by two attributes: i) name of the target object (e.g., table, row, or cluster node), ii) name of the operation(s) that will trigger metadata collection about the target object (e.g., table truncation, row insertion, or node decommissioning).  
NoSQL systems like Cassandra often offer a type of table called “system tables”. As these tables are persistent and easily accessible at servers, we store the Registry as a system table.
- *Log*: The Log is a table where collected metadata is stored. Unlike a flat file format, a table-formatted storage allows easy querying. Like the Registry, we store the Log as our second system table.
- *Core Logic and Internal API*: The Core Logic embodies the main logic of Wasef. It is implemented as a thin wrapper layer around the Registry and Log. To facilitate efficient metadata operations, it is integrated with the underlying NoSQL system (Section 3). Finally, it exposes an API for internal data store components to exploit Wasef.
- *System Hooks*: The System Hooks component contains implementations dependent on the underlying data store. It monitors data store operations (e.g., schema modification, data manipulation, etc.), and calls the Core Logic to log the metadata into the Log table.
- *Client (External) API*: The client API is a set of functions exposed to external clients (and users) allowing them to register objects and operations for metadata collection.

### 2.3. Operational Workflow

We give an overview of how metadata is registered, logged, and queried.

**2.3.1. Metadata Registration.** Metadata collection starts after a metadata target and operation have been registered in the Registry table. Targets are unique names of data entities, such as tables or data items, or internal components of the data store, such as SSTables or cluster nodes. We use a consistent and unified convention for naming metadata targets. Our Cassandra implementation uses the following naming convention:

```
<Keyspace name>.  
<Column family name>.  
<Comma-separated partitioning keys list>.  
<Dot-separated clustering keys list>.  
<Non-key column name>
```

For example, the target name of a column family (table) called `Teacher` that is part of `School` keyspace is named `School.Teacher`, while the target name for a row (data item for a teacher) with key `John` is `School.Teacher.John`. Partitioning and clustering keys are defined later in Section 3.1.1.

Operations are names of events occurring to targets, which trigger metadata collection. Examples are schema modification, row insertion, and node decommissioning. For instance, the column drop operation can be registered under the name `AlterColumnFamily_Drop`. As per Figure 1, registration can be performed via either the internal or external API.

**2.3.2. Metadata Logging.** System Hooks (Section 2.2) continuously monitor system operations. For each operation, the hook contacts the Core Logic, which checks if the operation and its target are in the Registry. If so, the hook collects the relevant metadata for that operation. For instance, for the `AlterColumnFamily_Drop` operation, the metadata collected is the name of the dropped column, the timestamp, and the database session owner. After that, Core Logic logs this metadata in the Log table.

**2.3.3. Metadata Querying.** The collected metadata accumulates in the Log table. The internal and external clients of the metadata can concurrently query the Log table via the APIs.

## 2.4. Metadata API

Wasef provides an intuitive CRUD-like API to both developers of the NoSQL system (internal API), and developers of client code (external API).

**2.4.1. Internal API.** The internal API provides operations for Registry and Log tables. The Registry API consists of three functions for managing a target-operation pair:

```
Registry.add(target, operation)  
Registry.delete(target, operation)  
Registry.query(target, operation)
```

`Registry.add` registers a target-operation pair by adding a new record into the Registry table. `Registry.delete` removes this entry. `Registry.query` returns a boolean indicating if the target-operation pair exists in the Registry.

The API to the Log consists of the following calls:

```
Log.add(target, operation, timestamp, value)  
Log.delete/Log.query(  
    target, operation, startTime, endTime)
```

`Log.add` inserts a new metadata record into the Log table. It starts by validating the target and operation parameters

against the Registry. Then a record is inserted into the Log containing the target-operation pair (as key), operation timestamp, and the metadata about the target-operation in the value field.

`Log.delete` removes all metadata for the target-operation pair. The `startTime` and `endTime` parameters are optional; if present, Wasef only removes matching records that were timestamped inside that time interval.

Finally, `Log.query` returns all the records identified by the mandatory target parameter. The last three parameters are optional and provide flexibility in querying the Log.

**2.4.2. External API.** Wasef's external API allows external clients to register, unregister, and query metadata:

```
register, unregister, queryAll, queryLatest  
Parameters for all four calls:  
(target, operation)
```

The `register` and `unregister` functions are wrappers around `Register.add` and `Register.delete` respectively from the internal API, and provide identical functionality. For convenience we also provide `queryAll` which retrieves all the records identified by the target-operation pair, as well as `queryLatest` which returns only the last inserted metadata record matching the criteria.

## 3. Implementation of W-Cassandra

We integrated Wasef into Cassandra 1.2. We call the resulting metadata-aware key-value store as W-Cassandra. The code is available for download at:

<http://dprg.cs.uiuc.edu/downloads>.

Before describing W-Cassandra's implementation, we first provide a brief outline of the existing Cassandra system.

### 3.1. Cassandra Background

Cassandra [1] is an open-source NoSQL data store intended for clusters of servers (i.e., nodes). It employs a peer-to-peer design. It can handle massive data sizes and scale out to large clusters. While Cassandra supports geo-distributed datacenters [8], and Wasef applies to those settings just as well, for brevity our discussion in this paper focuses on only a single-datacenter version.

#### 3.1.1. Data Model. Row-oriented Internal Operations:

Data operations such as reading, writing, compaction, and partitioning, are performed at the row-level, i.e., data item-level. Cassandra encodes all row-level operations internally in a unified form called *Row Mutation*. This representation is used for node-specific operations and inter-node communication.

**Hierarchical Organization:** Rows are grouped into column families (i.e., tables) such that the rows within a column family are identified by primary keys of the same type. A keyspace (or schema) is a logical grouping of column families, specifiable by a user. For instance, system tables are typically grouped under the keyspace named `System`.

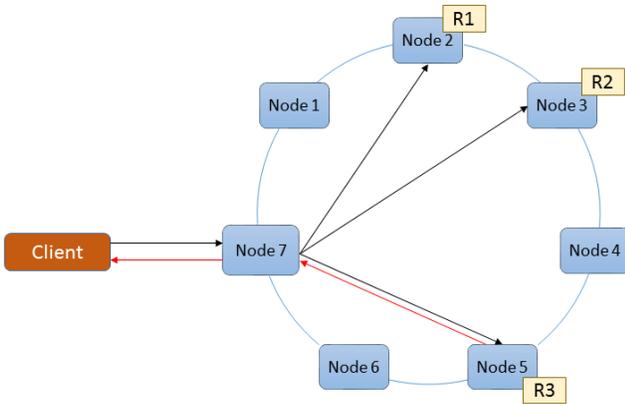


Figure 2: Client write request to Cassandra cluster. For a consistency level of ONE, the coordinator forwards the request to all replicas and waits only for the first reply.

**Partitioning vs. clustering keys:** Within a column family, the primary key of each row is divided into two parts. The first part is the partitioning key, which is used by Cassandra to place the key at a server in the cluster. Replica servers for that key are then assigned clockwise starting from that point in a virtual ring, as shown in Figure 2 (we omit ring details as they are not directly relevant to Wasef). The second part of the primary key is the clustering key, which is used by the Cassandra to cluster the non-key columns into a single storage block associated with the partitioning key.

**Cassandra Query Language (CQL):** Cassandra provides a SQL-like (but not SQL-complete) query language called CQL [9]. CQL commands include data definition queries (e.g., `create table`), data manipulation queries (e.g., `insert` and `select` for rows), and basic authentication queries to create database users and grant them permissions.

**3.1.2. Client Requests.** A client sends its read and write requests to a server called the coordinator. Any server can be used as a coordinator by any client.

**Write requests:** The coordinator forwards the request to all the replica servers of the data (typical replication factor = 3). Each replica writes locally, and acknowledges back to the coordinator. The coordinator then acknowledges back to the client. The client can specify, as a parameter called “consistency level”, how many replicas the coordinator should hear from before acknowledging to the client. Figure 2 shows an example for a consistency level of ONE.

**Read requests:** The coordinator forwards a read to a subset of the replicas, whose cardinality is equal to the client’s specified consistency level (these replicas are often chosen as the ones closest in network round-trip-time from the coordinator). For consistency levels higher than ONE (e.g., QUORUM or ALL), the coordinator receives multiple replica replies – it returns the highest timestamped version to the client.

### 3.2. Wasef Metadata Storage

Wasef collects and stores metadata by using two types of tables, in a way that offers low read latency and flexible querying. While implementing these techniques, we use underlying Cassandra tables. This enables Wasef to inherit Cassandra’s existing functionality such as data compression, caching, fast access, and replication factors.

Concretely, we store all metadata tables as Cassandra’s system tables, and collect them in the `system_metadata` system keyspace. Using system tables provides a read-only protection for the metadata schema, and makes it available immediately after the system is bootstrapped.

As explained earlier in Section 2.2, Wasef stores two system tables: Registry table and Log table. This separation has three advantages: i) because the Registry table is smaller than the Log, it can be queried quickly for each operation during metadata logging (Section 2.3.2); ii) management operations on registry entries are simplified, and iii) we can cleanly group entries within the Log, e.g., based on metadata insertion time or operation type.

Figure 3 illustrates, via an example, the schemas of the Registry and Log. The Registry table consists of two fields: The `target` field stores the name of the metadata target object, and the `operation` field stores the data store operation name which will trigger the metadata collection. The Log table has several fields that describe collected metadata. These include the `target`, the `operation`, and the timestamp of the operation (i.e., `time`). The `client` field reports the ownership information of the metadata target.

The shown primary keys of these two metadata tables are carefully chosen and ordered in order to achieve two goals:

- 1) *Optimizing the storage layout for low read latency:* The `target` key works as the partitioning key for both tables while the clustering keys are joined using a fixed scheme of delimiters. This is shown in Figure 3.B. Grouping the metadata related to one target within the same row orders the fields lexicographically and ensures they reside in the same Cassandra node, which leads to faster reading. As shown in Figure 3.B, all metadata for target name `School.Teacher.John` are grouped in the same row in the Log table. Every column in that row represents one operation. Using this layout, performing a select query that asks about all the operations related to one target is as fast as querying about one operation.
- 2) *Flexible querying of Log table:* In CQL, the `where` clause of the `select` statement filters only based on the table primary key. Thus, including more fields in the primary key increases querying flexibility.

### 3.3. Supported Targets and Operations

Table 1 shows the list of metadata targets and operations currently supported by W-Cassandra: these operations are from among those already present in CQL 1.2 [9]. The list can be extended to other CQL operations by adding

```
create table registry(
  target text,
  operation text,
  primary key( Target, operation) );
```

Partitioning key
Clustering key

```
create table log(
  target text,
  operation text,
  time long,
  client text,
  value text,
  primary key(Target, operation, time, client)
);
```

Partitioning key
Clustering keys

(A)

Registry

|                     |             |            |
|---------------------|-------------|------------|
| School.Teacher      | AlterCF_Add | Truncate   |
|                     | null        | null       |
| School.Teacher.John | Delete_Row  | Update_Row |
|                     | null        | null       |

Log

|                     |  |   |
|---------------------|--|---|
| School.Teacher      | AlterCF_Add-1509051314-admin<br>{col_name:address,col_type:text,<br>compaction_class:<br>SizeTieredCompactionStrategy} | AlterCF_Add-2009051414-admin<br>{col_name:mobile,col_type:text,<br>compression_sstable:<br>DefaultCompressor} |
| School.Teacher.John | Update_Row-1510051314-admin<br>{col_name:address,<br>col_old_val:null,col_new_val:'<br>Urbana,IL',ttl:432000}          | Update_Row-2010051414-admin<br>{col_name:mobile,<br>col_old_val:null,<br>col_new_val:'555555',ttl:432000}     |

(B)

Figure 3: Storage layout of Wasef. (A) Schemas of metadata tables, described in CQL. (B) Example of the internal storage layout of the metadata tables. Note how column names of the Log table are composed of the clustering primary key names.

appropriate per-operation system-dependent hooks to the code.

Metadata reporting starts after the system has been bootstrapped; concretely, after the authentication feature in Cassandra has been started. In order to implement this, we modified the write path so as to propagate the ownership information all the way down to the metadata Log. The ownership information is needed later to implement our data provenance use-case.

We observe that Table 1 is missing the create operation because we require an explicit `Log.add` to start metadata collection about a target. Implicit creates would have been complicated since it would have required adding non-existing objects in the Registry.

### 3.4. Optimizing Metadata Collection

Whenever a new operation arrives into the system, it is first validated against the metadata Registry. This is the sole overhead entailed for operations that do not have a Registry. Next, in case there is a Registry entry that matches, appropriate writes are entered into the Log.

To address the overhead of metadata collection for fine-grained metadata targets such as writes for a data item, we optimize both registry validation and log writing:

**Fast Registry Validation:** We speed up the response time for querying the Registry table in three ways. Each of these leverages underlying Cassandra functionalities.

- 1) **Enabling Dynamic Snitching:** We enable dynamic snitching, which allows the Cassandra coordinator to send read requests to replicas that are closest in network proximity (i.e., round-trip-time from the coordinator) based on history.
- 2) **Setting read consistency level to ANY for the Registry table:** This consistency level is faster than ONE, and it allows the coordinator to acknowledge the client after either storing it locally or receiving the first replica

acknowledgement, whichever occurs earlier. This also reduces the network traffic.

- 3) **Enabling row caching:** When row caching is enabled, Cassandra stores new entries in a cache associated with the destination table. Thus, Cassandra can serve read operations from the cache to shorten the read path.

**Lightweight Log Writing:** We employ two optimization techniques to reduce the Log writing overhead:

- 1) We perform write operations in a background thread. Cassandra uses the SEDA design [10], which divides the work into stages with a separate thread pool per stage. We adopt the same philosophy in order to improve the metadata write efficiency. We do so by injecting write mutations (Section 3.1.1) into the mutation handling stage, in a separate thread.
- 2) We set the write consistency level to ANY, just like for the Registry tables. This improves the time to write an entry into the Log, which is the most common metadata-related operation in W-Cassandra.

### 3.5. Discussion

**Replication, Fault-tolerance, Consistency Levels:** Since Wasef relies on Cassandra’s replication (replication factor = 3), Wasef’s fault-tolerance guarantees for metadata are identical to Cassandra’s fault-tolerance for data. Further, the replication factor is configurable per table in Wasef.

Orthogonal to replication is the issue of consistency level for operations. Wasef writes prefer a consistency level of ANY for speed, but this choice does not affect fault-tolerance: a metadata write is acknowledged only after some server has written it to its commit log, thus making it durable. Consistency level only delays propagation of writes to replicas – Cassandra’s hinted handoff and read repair mechanisms propagate writes eventually (and quickly) to replicas. Finally, the administrator retains the option to

| Target | Identifier                      | Operations             | Metadata  |
|--------|---------------------------------|------------------------|---|
| Schema | Name                            | Alter, Drop            | Old and new names, replication map                                  |
| Table  | Name                            | Alter, Drop, Truncate  | column family name, column names and types, compaction strategy, .. |
| Row    | Partitioning keys               | Insert, Update, Delete | key names, affected columns, TTL and timestamp                      |
| Column | Clustering keys and column name | Insert, Update, Delete | key names, affected columns, TTL and timestamp                      |
| Node   | Node ID                         | nodetool decommission  | Token ranges  |

Table 1: Supported metadata targets and operations in W-Cassandra.

increase Wasef’s consistency level to `Quorum` or `ALL`, depending on the application.

**Wasef vs. Alternative Approaches:** There are two alternative approaches to implementing metadata: i) a stand-alone implementation running on dedicated servers, and ii) application-specific metadata. Firstly, if metadata were stored on dedicated servers, we would have to re-implement querying, replication, fault-tolerance, and load-balancing. Instead Wasef’s integration into Cassandra allows us to directly use these features from Cassandra. Further, dedicated servers would need tuned scaling as datasize and cluster size increases – Wasef scales automatically (as our experiments show). Secondly, implementing metadata in the application would take many man-hours for each application. Our approach is fast yet flexible. If an application does not need metadata, it can disable Wasef.

## 4. Use Cases: Leveraging The Metadata

Wasef opens the door to leverage the power of collected metadata in order to address system pain points and provide new features. To demonstrate this, we implemented three use-case scenarios selected from three different domains.

### 4.1. Enabling Flexible Column Drops

**Problem.** The JIRA issue 3919 [11] reports that in Cassandra 1.2, dropping a column removes its definition from the table schema but keeps the column data stored in the system. This leads to incorrect behaviors, e.g., if another column is added to the table with a name equal to the dropped one, CQL queries referencing the new column will still return data belonging to the old column.

**Metadata Solution.** We present a correct and flexible column drop implementation using Wasef. The intuition behind the flexibility in our approach is as follows – in today’s file systems (e.g., Mac, Windows, Unix) when a user deletes a file or folder, the user has a second chance to retrieve it from a Trash/Deleted Items folder. We implement a similar intuition for columns. Concretely, Figure 4 illustrates our state diagram for our flexible column drop. When a column is dropped the first time, it becomes unavailable for CQL

queries, but the data is retained. Thereafter, either on a user’s explicit second delete command or after a (configurable) timeout, the column data is permanently purged. The user can un-delete the column from the tentatively dropped state.

This state diagram is implemented as follows. When a column is first dropped using the `Alter Table` operation, W-Cassandra inserts a metadata called `AlterColumnFamily_Drop` into the Registry. If the column is subsequently re-added, then this metadata is deleted, leaving the system in the initial state. However, if the column were to be dropped again explicitly or the timeout expires, a metadata log entry that contains the column name and the time of the drop operation would be inserted into the Log table and the column marked for permanent deletion. Before permanent deletion, any attempts to insert new columns with the old name are rejected.

During a CQL select query, W-Cassandra checks both the `AlterColumnFamily_Drop` metadata and the log entry, and excludes the dropped (tentative or permanent) columns from the query’s results.

In fact, Cassandra provided a fix for the column-drop issue in a later release (Cassandra 2.0) by retaining a history of dropped columns. However, this history is maintained in a specialized hash map attached to the column schema, and the hash map is replicated at *all* cluster nodes. This ad-hoc solution took many man-months to implement. Instead, our W-Cassandra approach took only 50 additional lines of code to write, is more flexible and systematic, and leverages underlying NoSQL features such as replication.

### 4.2. Node Decommissioning

**Problem.** Node decommissioning, a utility that is part of Cassandra’s NodeTool, removes a node from the cluster by assigning its token ranges to other nodes and replicating the data accordingly. Verifying that data is not lost during this operation is a major pain point [12] because the administrator needs to manually count the token (i.e., key) ranges assigned to the decommissioned node, and then manually check if each has been reassigned to another node. Under some circumstances, this operation may be quite critical, e.g., when the decommissioned node is the only one holding certain data with a replication factor of 1.

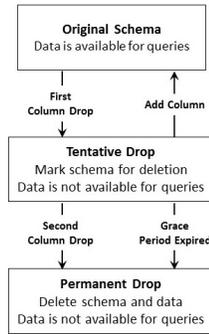


Figure 4: State Machine for Flexible Column Drop in W-Cassandra.

**Metadata Solution.** We implement an automated verification tool for the node decommissioning operation. The tool verifies for system administrators that all the data which existed in a decommissioned node have indeed been safely migrated to other nodes in the cluster.

The basic decommissioning workflow is the same as in Cassandra, and is as follows. The decommissioned node streams its data to other nodes in the cluster and then unbootstraps. The new destinations of the streamed data are calculated as follows. First, the token ranges, at each server, for each non-system table are collected. Second, the partitioner and the replication strategy are used to decide the new replica for each token range. Third, all the collected information is passed to the file streamer to move the blocks (SSTables) to the intended destinations. Finally, the decommissioned node is retired.

We exploit the metadata collected by Wasef to verify node decommissioning as follows. During the decommission command, we use `Log.add` API (Section 2.4) to store the new replica for each token range hosted by the decommissioned node into the metadata Log, by using a metadata tag called `decommission`. Then, when the decommissioned node leaves the cluster, the operation can be verified using the following available command:

```
nodetool decommission -verify <node IP>
```

This command retrieves the metadata records of the decommissioned nodes (using `Log.query` API) and verifies that all the ranges are currently available in the system using the partitioner and the replication strategy of each keyspace. Section 5.6 evaluates the overhead of this operation.

### 4.3. Providing Data Provenance

**Problem.** Data provenance information is essential for users to be able to tell where data comes from, when it was produced, and in general to keep track of it. This is a critical issue for many disciplines that curate and store data, ranging from manufacturing to bio-informatics to astronomy. Unfortunately, modern NoSQL systems largely do not support data provenance collection by default.

**Metadata Solution.** Using the techniques of Section 3.3, we collect the following provenance data about each operation in Cassandra:

- 1) Full name of target of data operation. E.g., dropping a table called `User` located in `Test` keyspace results in logging `Test.User` as the full name.
- 2) Operation name. E.g., `Alter_Add_Column` Family indicates a new column addition to a column family.
- 3) Operation time: timestamp of the operation.
- 4) The authenticated session owner name.
- 5) Results of the operation. E.g., when a new column is added, the name of the column and its attributes are logged. When a column name is modified, its old and new names are logged.

This provides two desirable types of query provenance (discussed later in Section 6): i) *Where Provenance*, which keeps track of the records from which the query results were derived. In the absence of joins in NoSQL systems, Wasef provides Where Provenance through the “full name” of the metadata target. ii) *Why Provenance*, the justification of the results via a listing of operations that produced them in Wasef’s “operation name” field in the Log.

This provenance data can be queried via our external APIs (Section 2.4). Features such as replication, scalability, and accessibility are enabled as usual for this provenance data.

For correctness, we do not automatically garbage-collect old provenance data. However, system administrators may use the delete APIs provided by our system to manually delete old provenance data entries, e.g., based on their timestamps.

## 5. Experimental Evaluation

Our experiments are designed to answer the following questions:

- 1) What is the performance cost of integrating metadata collection and querying into Cassandra? This includes read and write latencies, and the overall throughput, for W-Cassandra.
- 2) How does W-Cassandra scale with cluster size, size of data, size of metadata, and query injection rate?
- 3) How does W-Cassandra perform for the use-case scenarios of Section 4?

We run our experiments on the Amazon Web Services (AWS) EC2 public cloud infrastructure [6]. We inject traces using the YCSB benchmark [7].

### 5.1. Experimental Setup

Our experiments use six AWS EC2 `m1.large` instances, each with 2 virtual CPUs (4 ECUs), 7.5 GB of RAM, and 480 GB of ephemeral disk storage. We run the YCSB client from a separate identical instance. The instances use Ubuntu 12.04 64-bit operating system with swapping turned off, as recommended by Cassandra for production settings. YCSB affords multiple types of workloads: by default we use the

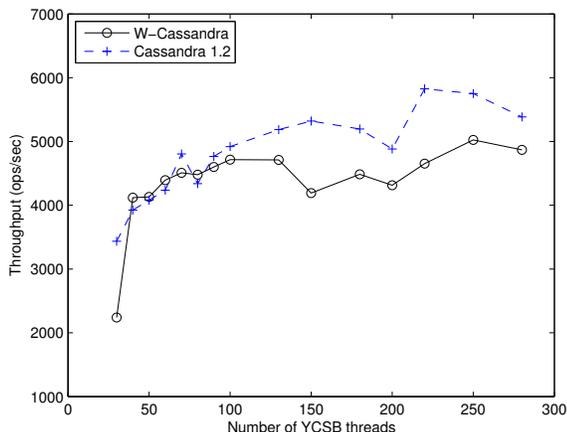


Figure 5: Throughput against number of clients: Throughput comparison between the standard Cassandra 1.2 and W-Cassandra. The experiment uses a cluster of six EC2 instances and a data set size of 12 GB. Each point in the graph is the average ops/s arising from a workload of 1 million YCSB client requests. The average difference between the two lines is 9%.

“Zipfian” workload that uses the Zipf distribution queries to objects based on their popularity.

## 5.2. W-Cassandra Throughput and Latency

We first measure the effect of Wasef on Cassandra’s throughput and latency. We conduct 15 YCSB runs for both W-Cassandra and standard Cassandra (1.2) using a heavy-update workload (50% update and %50 read). The database size is 12 GB, keys are 8 B, and values are 1 KB.

Figure 5 shows the throughput, which increases quickly at first and then saturates at about 300 threads. Over all the data points in the plot, the average performance of W-Cassandra is only 9% worse than Cassandra. This value captures the overhead of Wasef’s metadata collection.

Figures 6 and 7 depict the update and read latencies, respectively, for the two systems. Over all the data points in the plot, the average read latency is only 3% worse in W-Cassandra, and update latency is only 15% worse. The latter is higher because an update engenders additional metadata writes.

## 5.3. Scalability with Cluster Size

We linearly increase the number of nodes in the cluster, while proportionally scaling the data set size and system load. From one run to the next in Figure 8, the cluster size was increased by two nodes, the data set increased by 4 GB, and load increased by 50 YCSB threads.

Figure 8 shows that W-Cassandra retains linear scalability in spite of its metadata overheads (“Scalability difference” line). The percentage overhead in update latency rises slowly with scale because W-Cassandra injects one metadata validation request (i.e., read request) per update. This request has a high probability to be served locally when

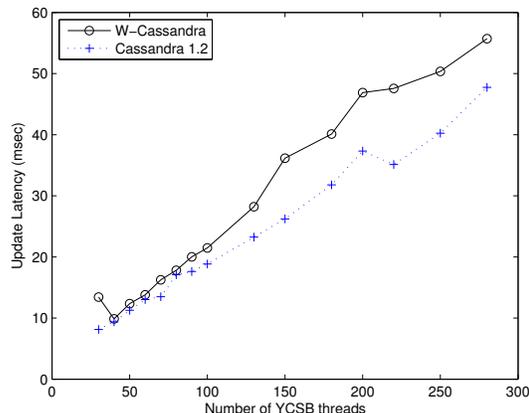


Figure 6: Update latency against number of clients: Update latency comparison between the standard Cassandra 1.2 and W-Cassandra. The experiment uses a cluster of six EC2 instances and a data set size of 12 GB. The workload contains 500K YCSB client requests.

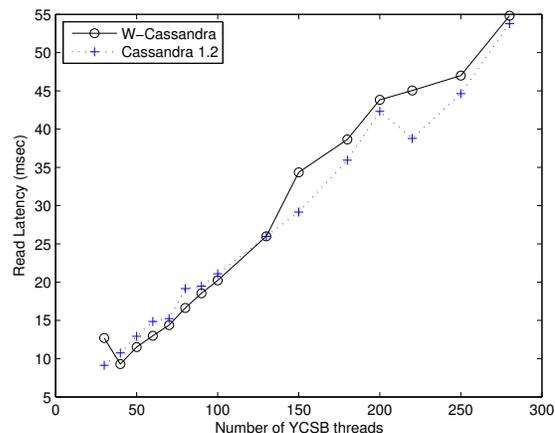


Figure 7: Read latency against number of clients: Read latency comparison between the standard Cassandra 1.2 and W-Cassandra. The experiment uses a cluster of six EC2 instances and a data set size of 12 GB. The workload contains 500K YCSB client requests.

the number of nodes is small. However, when the number of nodes increases the probability to serve the request locally decreases. Yet, the overhead is still small: at 10 nodes, the update latency increase due to W-Cassandra is only about 10%.

## 5.4. Column Drop Feature

We compare the latency for column drop in W-Cassandra with that in Cassandra 1.2 using a data set size of 8 GB. While Cassandra 1.2 implements a crude (and incorrect<sup>2</sup>) version of column drop, we choose to compare against

2. Cassandra 1.2’s column drop deletes the schema definition but retains the data.

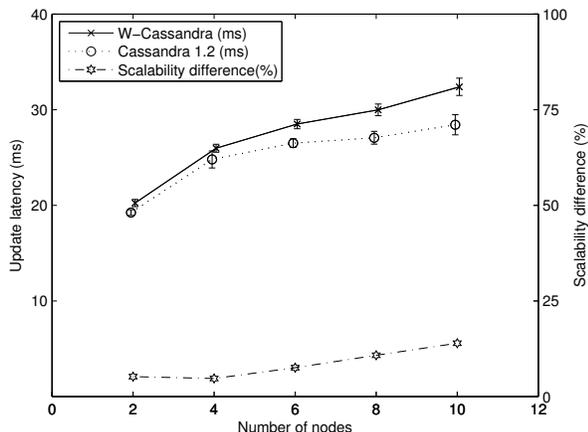


Figure 8: Scalability against cluster size: A scalability comparison between the standard Cassandra 1.2 and W-Cassandra. From one run to the next, the experiment increased the number of cluster size by two nodes, data size by 4 GB, and YCSB load by 50 threads. Datapoints are perturbed horizontally for clarity.

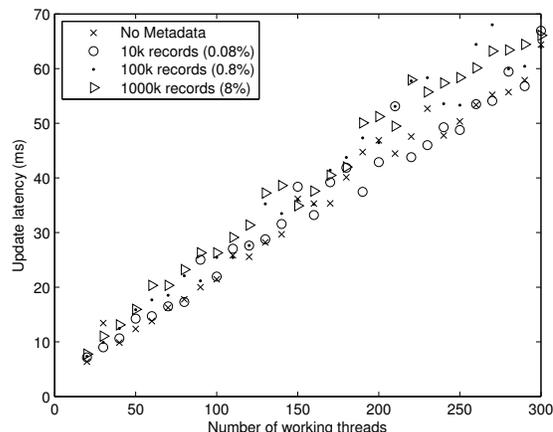


Figure 10: Update latency comparison between different metadata sizes registered in W-Cassandra: The experiment uses six EC2 instances with total data size of 12 GB.

### 5.5. Scalability with Data Size (Collecting Data Provenance)

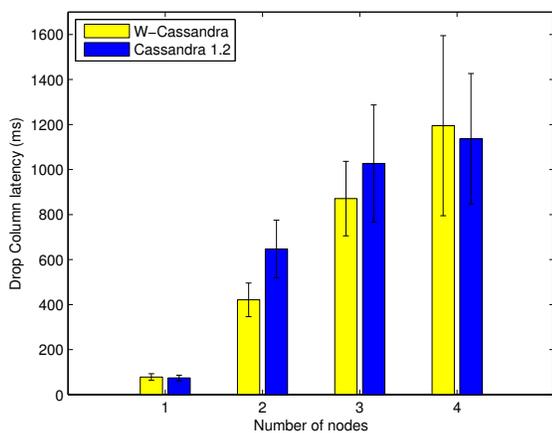


Figure 9: Column drop latency against cluster size: The latency of column drop operation for the standard Cassandra 1.2 compared to W-Cassandra. Each bar in the graph represents the average of 500 drop column operations performed by clients running at a separate machine. The data set size is 8 GB. The overhead at 4 nodes is 5%.

it because: i) Wasef is built into Cassandra 1.2, and ii) comparing against the latest version Cassandra 2.0 would be unfair as this version is faster than 1.2 because of several optimizations that are orthogonal to column dropping.

Since the YCSB benchmark does not offer schema modification tests, we designed a customized test that performs a set of 500 column drop operations. Figure 9 shows that the latency of W-Cassandra hovers at or around Cassandra's column drop latency. Since Cassandra 1.2's implementation is incorrect, this is a positive result.

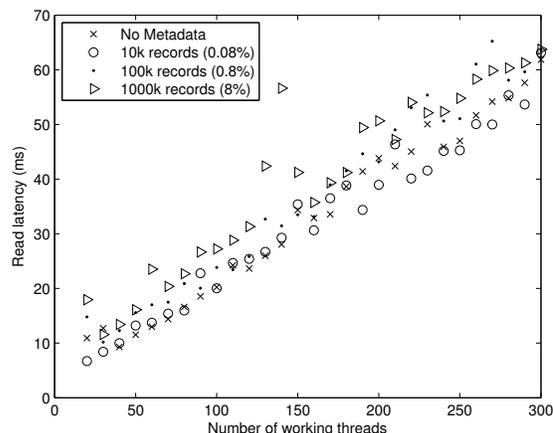


Figure 11: Read latency comparison between different metadata sizes registered in W-Cassandra: The experiment uses six EC2 instances with total data size of 12 GB.

Figures 10 and 11 show that as the metadata size is increased from 0.08% to 8% of data size, the increase in operation latencies, while provenance is being collected, is generally very small. Independent of its size, this metadata is in fact replicated across multiple servers, thus allowing it to scale with data size. Finally, we note that Wasef is memory-bound rather than disk-bound because Cassandra is too. A disk-bound Cassandra would be very slow, and would lead the administrator to add more servers, making it, and thus Wasef, memory-bound again.

### 5.6. Verifying Node Decommissioning

The main overhead faced by the system administrator during node decommissioning is the first stage when token

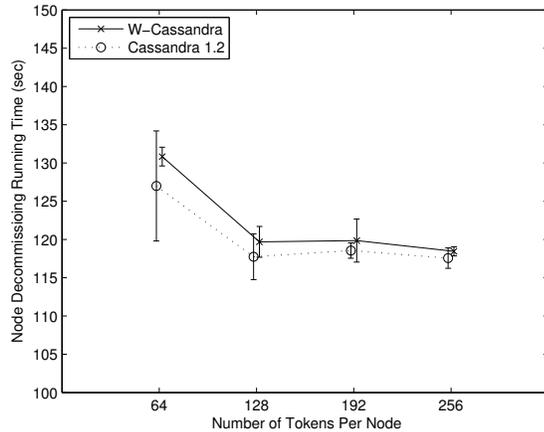


Figure 12: Running time for node decommissioning operation:

The running time of the node decommission operation for standard Cassandra 1.2 compared to W-Cassandra. The data set size is 4 GB. The average difference between the two lines is 1.5%. Datapoints are perturbed horizontally for clarity.

metadata is collected; thereafter the data streaming to other servers is automated. To measure the overhead of the first stage, we vary the number of tokens per node. We use four AWS EC2 instances, and a 4 GB data set size. Figure 12 shows that W-Cassandra is only marginally slower than Cassandra; the average overhead was measured at 1.5%.

## 6. Related Work

**Types of Metadata** The term metadata was coined by Bagley in the context of programming languages [13]. Generally, there are two types of metadata [14]: i) structural metadata that describes database entities (e.g., tables) and their hierarchical relationships, and ii) descriptive metadata is data about the data items in the database. Additionally, NISO also defines administrative metadata [15]. Our Wasef system collects all three types of metadata.

**Database Metadata:** Metadata systems can be implemented internal to a database [16], [17] or externally [18]–[21]. Examples of external metadata systems include those for businesses [21] and for Grids [22]–[25]. Internal metadata systems like [4] are used to collect structural metadata. Trio [17] is a data management system atop a relational database (e.g., Postgres), but it is not a general solution for metadata.

**Metadata in Cloud Data Stores:** Many have argued that metadata should be a feature of cloud data stores [26]–[30]. Client-centric approaches for metadata [26] are too intrusive; we believe that metadata collection should be server-centric.

**Data Provenance:** Provenance information is managed in scientific workflows [31]–[34], monitoring system operations [29], [35]–[37], and database queries [38]–[42]. Query provenance in relational database is of two kinds [43]: i) *Where Provenance* describes source records of a query’s result, and ii) *Why Provenance* justifies query results by its source operations and relations between source records. Wasef provides both these kinds of provenance.

There has been some recent work on provenance in key-value stores [16], [26], [44]. The KVPMC system for Cassandra [44] collects provenance data on request, provides client access, and can store provenance data internally or externally. However, Wasef is preferable for four reasons: i) it is a general solution for any modern NoSQL system, ii) it collects all kinds of metadata, not merely provenance, iii) KVPMC is client-side while Wasef is server-side, and iv) Wasef imposes less overhead and provides good scalability. **Metadata at Cloud Providers:** Cloud datastores support some form of metadata [45]–[47]. AWS S3 [46] provides both system-defined metadata (e.g., object size, creation time) and user-defined metadata. However, these metadata services are inflexible, e.g., metadata size is limited to tens of KBs, and querying is inexpressive.

## 7. Conclusion

We presented a metadata system for NoSQL data stores, called Wasef. We integrated Wasef into Cassandra. We showed how our system, called W-Cassandra, can be used to correctly and flexibly provide features like column drop, node decommissioning, and data provenance. Our experiments showed that our system imposes low overhead on Cassandra throughput of 9% and read latency of 3%. We also showed that our system scales well with cluster size, incoming workload, data size, and metadata size. We believe that Wasef opens the door to treating metadata as first-class citizens in NoSQL systems, and exploring the myriad forms of metadata that abide in this new class of data stores.

## References

- [1] A. Lakshman and P. Malik, “Cassandra: A Decentralized Structured Storage System,” *ACM SIGOPS Operating Systems Review*, vol. 44, no. 2, pp. 35–40, 2010.
- [2] “MongoDB,” <https://www.mongodb.org/>.
- [3] “Project Voldemort,” <http://www.project-voldemort.com/voldemort/>.
- [4] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber, “Bigtable: A Distributed Storage System for Structured Data,” *ACM Transactions on Computer Systems (TOCS)*, vol. 26, no. 2, p. 4, 2008.
- [5] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Voshall, and W. Vogels, “Dynamo: Amazon’s Highly Available Key-value Store,” in *ACM SIGOPS Operating Systems Review*, vol. 41, no. 6. ACM, 2007, pp. 205–220.
- [6] “Amazon Web Services (AWS),” <http://aws.amazon.com/>.
- [7] B. F. Cooper, A. Silberstein, E. Tam, R. Ramakrishnan, and R. Sears, “Benchmarking Cloud Serving Systems with YCSB,” in *The First Symposium on Cloud Computing*. New York, NY, USA: ACM, 2010, pp. 143–154.
- [8] T. Rabl, S. Gómez-Villamor, M. Sadoghi, V. Muntés-Mulero, H.-A. Jacobsen, and S. Mankovskii, “Solving Big Data Challenges for Enterprise Application Performance Management,” *The VLDB Endowment*, vol. 5, no. 12, pp. 1724–1735, 2012.
- [9] “CQL for Cassandra 1.2,” <http://www.datastax.com/documentation/cql/3.0/cql/aboutCQL.html>.
- [10] M. Welsh, D. Culler, and E. Brewer, “SEDA: An Architecture for Well-Conditioned, Scalable Internet Services,” in *Operating Systems Review*, vol. 35, no. 5. ACM, 2001, pp. 230–243.

- [11] "CASSANDRA-3919 JIRA Issue," <https://issues.apache.org/jira/browse/CASSANDRA-3919>.
- [12] "Slightly Remarkable Blog, Removing Nodes from a Cassandra Ring," <http://slightlyremarkable.com/post/57852577144/removing-nodes-from-a-cassandra-ring>.
- [13] P. R. Bagley, "Extension of Programming Language Concepts," DTIC Document, Tech. Rep., 1968.
- [14] F. P. Bretherton and P. T. Singley, "Metadata: A User's View," in *Seventh International Working Conference on Scientific and Statistical Database Management*. IEEE, 1994, pp. 166–174.
- [15] N. Press, "Understanding Metadata," *National Information Standards*, vol. 20, 2004.
- [16] K.-K. Muniswamy-Reddy and M. Seltzer, "Provenance as First Class Cloud Data," *ACM SIGOPS Operating Systems Review*, vol. 43, no. 4, pp. 11–16, 2010.
- [17] C. C. Aggarwal, "Trio: A System for Data Uncertainty and Lineage," in *Managing and Mining Uncertain Data*. Springer, 2009, pp. 1–35.
- [18] E. Deelman, G. Singh, M. P. Atkinson, A. Chervenak, N. P. Chue Hong, C. Kesselman, S. Patil, L. Pearlman, and M.-H. Su, "Grid-Based Metadata Services," in *The 16th International Conference on Scientific and Statistical Database Management*. IEEE, 2004, pp. 393–402.
- [19] G. Singh, S. Bharathi, A. Chervenak, E. Deelman, C. Kesselman, M. Manohar, S. Patil, and L. Pearlman, "A Metadata Catalog Service for Data Intensive Applications," in *ACM/IEEE Conference in Supercomputing*. IEEE, 2003, pp. 33–33.
- [20] M. Balakrishnan, D. Malkhi, T. Wobber, M. Wu, V. Prabhakaran, M. Wei, J. D. Davis, S. Rao, T. Zou, and A. Zuck, "Tango: Distributed Data Structures Over a Shared Log," in *The 24th ACM Symposium on Operating Systems Principles*. ACM, 2013, pp. 325–340.
- [21] "Oracle Metadata Service in Fusion Middleware 11g," <http://www.oracle.com/technetwork/developer-tools/jdev/metadataservices-fmw-11gr1-130345.pdf>.
- [22] M. B. Jones, C. Berkley, J. Bojilova, and M. Schildhauer, "Managing Scientific Metadata," *Internet Computing, IEEE*, vol. 5, no. 5, pp. 59–68, 2001.
- [23] M. Xiong, H. Jin, and S. Wu, "2-Layered Metadata Service Model in Grid Environment," in *Distributed and Parallel Computing*. Springer, 2005, pp. 103–111.
- [24] E. Deelman, G. Singh, M.-H. Su, J. Blythe, Y. Gil, C. Kesselman, G. Mehta, K. Vahi, G. B. Berriman, J. Good *et al.*, "Pegasus: A Framework for Mapping Complex Scientific Workflows Onto Distributed Systems," *Scientific Programming*, vol. 13, no. 3, pp. 219–237, 2005.
- [25] P. Zhao, A. Chen, Y. Liu, L. Di, W. Yang, and P. Li, "Grid Metadata Catalog Service-Based OGC Web Registry Service," in *The 12th International Workshop on Geographic Information Systems*. ACM, 2004, pp. 22–30.
- [26] K.-K. Muniswamy-Reddy, P. Macko, and M. I. Seltzer, "Provenance for the Cloud," in *FAST*, vol. 10, 2010, pp. 15–14.
- [27] M. A. Sakka, B. Defude, and J. Tellez, "Document Provenance in the Cloud: Constraints and Challenges," in *Networked Services and Applications-Engineering, Control and Management*. Springer, 2010, pp. 107–117.
- [28] O. Q. Zhang, M. Kirchberg, R. K. Ko, and B. S. Lee, "How to Track Your Data: The Case for Cloud Computing Provenance," in *The Third International Conference on Cloud Computing Technology and Science (CloudCom)*. IEEE, 2011, pp. 446–453.
- [29] P. Macko, M. Chiarini, M. Seltzer, and S. Harvard, "Collecting Provenance Via the Xen Hypervisor," in *The Third USENIX Workshop on the Theory and Practice of Provenance*, 2011.
- [30] M. Imran and H. Hlavacs, "Provenance in the Cloud: Why and How?" in *The Third International Conference on Cloud Computing, GRIDS, and Virtualization*, 2012, pp. 106–112.
- [31] R. Bose and J. Frew, "Lineage Retrieval for Scientific Data Processing: A Survey," *ACM Computing Surveys (CSUR)*, vol. 37, no. 1, pp. 1–28, 2005.
- [32] Y. L. Simmhan, B. Plale, and D. Gannon, "A Survey of Data Provenance Techniques," *Computer Science Department, Indiana University, Bloomington IN*, vol. 47405, 2005.
- [33] S. B. Davidson, S. C. Boulakia, A. Eyal, B. Ludäscher, T. M. McPhillips, S. Bowers, M. K. Anand, and J. Freire, "Provenance in Scientific Workflow Systems," *IEEE Data Engineering Bulletin*, vol. 30, no. 4, pp. 44–50, 2007.
- [34] C. Goble, "Position Statement: Musings on Provenance, Workflow and (Semantic Web) Annotations for Bioinformatics," in *Workshop on Data Derivation and Provenance, Chicago*, 2002.
- [35] K.-K. Muniswamy-Reddy, D. A. Holland, U. Braun, and M. I. Seltzer, "Provenance-Aware Storage Systems," in *USENIX Annual Technical Conference, General Track*, 2006, pp. 43–56.
- [36] C. Sar and P. Cao, "Lineage File System," <http://crypto.stanford.edu/~cao/lineage>.
- [37] R. K. Ko, P. Jagadpramana, and B. S. Lee, "Flogger: A file-Centric Logger for Monitoring File Access and Transfers Within Cloud Computing Environments," in *The 10th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*. IEEE, 2011, pp. 765–771.
- [38] A. P. Chapman, H. V. Jagadish, and P. Ramanan, "Efficient Provenance Storage," in *International Conference on Management of Data*. ACM, 2008, pp. 993–1006.
- [39] A. Meliou, W. Gatterbauer, and D. Suciu, "Bringing Provenance to Its Full Potential Using Causal Reasoning," *The Third USENIX Workshop on the Theory and Practice of Provenance*, 2011.
- [40] S. Gao and C. Zaniolo, "Provenance Management in Databases Under Schema Evolution," in *The Fourth USENIX Conference on Theory and Practice of Provenance*. USENIX Association, 2012, pp. 11–11.
- [41] P. Buneman, J. Cheney, and E. V. Kostylev, "Hierarchical Models of Provenance," in *The Fourth USENIX Conference on Theory and Practice of Provenance*, 2012, p. 10.
- [42] D. Bhagwat, L. Chiticariu, W.-C. Tan, and G. Vijayvargiya, "An Annotation Management System for Relational Databases," *The VLDB Journal*, vol. 14, no. 4, pp. 373–396, 2005.
- [43] P. Buneman, S. Khanna, and T. Wang-Chiew, "Why and Where: A Characterization of Data Provenance," in *Database Theory-ICDT 2001*. Springer, 2001, pp. 316–330.
- [44] D. Kulkarni, "A Provenance Model for Key-Value Systems," in *The Fifth USENIX Workshop on The Theory and Practice of Provenance*. USENIX, 2013.
- [45] "Google Cloud Data Store," <https://developers.google.com/datastore/>.
- [46] "Amazon Simple Storage Service (S3)," <http://aws.amazon.com/s3/>.
- [47] "Open Stack Software," <https://www.openstack.org>.