

Ledger as an auditing solution for Cloud Database

Akash Malla
Computer Science
and Engineering Dept.
Santa Clara University
amalla1@scu.edu

Sharon Subathran
Computer Science
and Engineering Dept.
Santa Clara University
ssubathran@scu.edu

Neha Soma
Software
Engineering Dept.
Santa Clara University
nsoma@scu.edu

Sneha Mule
Computer Science
and Engineering Dept.
Santa Clara University
smule@scu.edu

<https://github.com/SomaNeha/COEN241CloudChainTeamProject>

Abstract—NoSQL databases, being non-relational, distributed and horizontally scalable are able to satisfy most of the needs of the present day applications. They provide high performance and scalability even when handling large volumes of unstructured data. This paper aims to implement a system that captures and aggregates data from server logs, providing a platform from which data can be easily reviewed and analyzed. In this project, we capture all long-running database queries that are requested by applications deployed on AWS EC2 instance. We maintain a list of all such query requests and call it a ledger. This ledger can be used to carry out performance analysis of such queries.

Keywords: Distributed Ledger, Digital Signature, Virtual Machine (VM), NoSQL (MongoDB)

I. INTRODUCTION

As data keeps increasing in size and has an unstructured form, there is a dire urgency to store big data efficiently and in a secure manner. Traditional RDBMS databases have limitations in terms of non-linear query execution time and static relational schema model. As a workaround to increase performance, a primary and secondary index could be created on selected table columns in a schema in order to retrieve data faster. The addition of indexes has a cost to it as it requires more space, possibly terabytes of data, which in turn degrades performance drastically. The cost of hardware to store data and maintenance to support are major drawbacks. Although RDBMS has a fair share of benefits that include robustness, simplicity, flexibility, and performance, it is not the right solution to handle the current scale of data warehousing, data inflow from cloud applications, social media, and IoT devices. As a result of problems from RDBMS, NoSQL database has been developed to be schemaless and provide better performance for big data storage. NoSQL (MongoDB) databases are capable to scale horizontally to ensure high availability and scalability. It supports dynamic schema to store semi-structured or unstructured data.[3].

The main reason why many organizations still use RDBMS is solely due to inadequate security in place of NoSQL databases. In 2013, Cloud Security Alliance (CSA) alluded that NoSQL databases have poor authentication and no protection for data integrity. A common solution to protect sensitive data is through data encryption[4]. MongoDB is a popular choice to store application logs thanks to its schema-free

The following Table-I, gives the comparison between Distributed Ledger based CryptDB and other databases

TABLE I
COMPARISON BETWEEN VARIOUS DATABASES

Name	Description	Pros	Cons
RDBMS	Structured	Complies all ACID properties and has some security	Slow to run query for big data. Requires creation of index on data to improve performance
NoSQL	Schemaless, Scalable, unstructured database	Good performance to run queries to retrieve big data compared to RDBMS	Lacks proper security. Eventual consistency can let update, insert and deletion of records take time
CryptDB	Schemaless, Scalable, unstructured database	Provides strong security as data is encrypted	Eventual consistency can let update, insert and deletion of records take time.
Distributed Ledger (Our proposed solution)	Schemaless, Scalable, unstructured database that uses distributed ledger for validation of queries	Improves current performance by monitoring active MongoDB operations	Additional Space for ledger and requires daemon process to run in order to capture current operations for an interval

design making it flexible enough to use for storing data whose schema tends to change from time-to-time.

Servers generate a large number of events/logs (i.e. logging,) that contain information about their operation including errors, warnings, and users behaviour. Clients extract data from these database logs to carry out analysis on the performance. By default, this data is stored in plain text log files which are time-consuming, difficult to review, reference, and visualize without an efficient system for aggregating and storing this data. This is the system that we implement- a way to eliminate the heavy lifting of formatting and capture all active operation-related data in a single collection, making it easy to reference and analyze for a client to be able to query. The project was

implemented in public cloud, it could implement as a database component in private cloud models.

A. Motivation/Background

The use of Blockchain technology[2] was first introduced by Bitcoin. The chronic feature of this technology enhances trust through transparency and traceability within any transaction of data and financial resources. It has a set of benefits: decentralized control, there is no central authority that owns or controls the network thus avoiding a single point of failure; immutability, where written data is tamper-resistant forever; and the ability to create and transfer assets on the network, without reliance on a central entity. Decentralized Ledger is a key feature of this technology and Bitcoin meant for it to behave as a distributed transaction management system to avoid spending the coin multiple times. So there is a list which consists of various transactions, said to be a ledger or block. Each block/ledger is generated by proof-of-work and appended to the existing blockchain. The concept of distributed ledger provides an incentive platform for distributed database since it exists across a decentralized network of machines that allows transparent transaction mechanisms in various businesses [Fig 1].

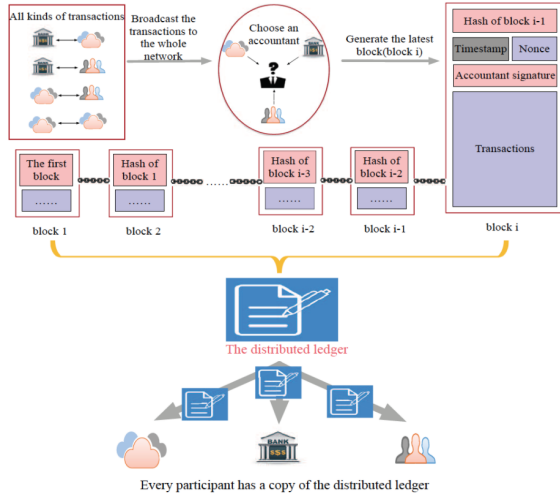


Fig. 1. Demonstrates blockchain for financial institutions

As VM management system plays a central role in a distributed ledger for cloud database operations that involve important information, a protection mechanism is built. This mechanism guarantees only authorized users can access corresponding data on the ledger. We also evaluate the proposed system from performance and security perspectives to demonstrate its usefulness in the real world environment.

II. PROBLEM ANALYSIS

Virtual Machine refers to the software implementation of any computing device or machine which is used to execute the program as physical machines. When a user tries to work on virtual machines they acquire the resources of the

virtualization installed in the remote machines which are been accessed by related protocols so as to behave like a local machine. We used public cloud in our project as private clouds often cost much more than public clouds. In order to implement and work on shared systems so that the resources and workload are been shared among the systems.

There are many ways to isolate workloads; such as are separate machines, physically separate machines, virtualization, containerization (OS-level virtualization). It is said that virtualization can be used on physically separated machines, or by running containers inside hardware virtualization which says it is one of the strongest isolation mechanisms in the cloud environment with inbuilt cloud security concerns to protect from data breaches. A strong form of isolation is to not connect machines at all through any network. This has its advantages for security. It means that the vulnerabilities in a computer, or an operating system's networking stack, or the OS low-level services, cannot be exploited over a network. In virtualization. For more security concerns instead of using virtualization, we used VM instead Docker- container to avoid a large number of data leakages that damage the organizations.

As we have log and Operations in the database so that we can have a history of queries running and keep track of them by time. The log is based on changes that are made to data in the tables that you track. Also, the log consists of, entries are which chronologically ordered and show changes that are made to the fields on the specified tables.

In our project we would like to have the collection of all logs when and what changes with respective time, date and updated queries that are made in the ledger. As Servers generate a large number of events/logs (i.e. logging,) that contain information about their operation including errors, warnings, and users behavior. Clients extract data from these database logs to carry out analysis on the performance. By default, this data is stored in plain text log files which are time-consuming, difficult to review, reference, and visualize without an efficient system for aggregating and storing this data. This is the system that we implement a way to eliminate the heavy lifting of formatting and capture all active operation-related data in a single collection, making it easy to reference and analyze for a client to be able to query.

A. Contributions

Our contributions in this work are summarized as follows:

- We propose a solution that employs the concept of Distributed Ledger Technology (DLT) to enhance the security and consistency of data.
- This paper focuses typically to help, many clients who have to extract data from database logs to analyze performance issues.
- Reviewing and formatting log data is time-consuming, if not, more time consuming than analyzing formatted data using a visualization utility or machine learning algorithm
- We wanted to eliminate the heavy lifting of formatting and capture all active operation-related data in a single collection for any client to be able to query.

III. PROPOSED SOLUTION

We started off the initial phase of our project researching blockchain concepts and wanted to implement them in such a way so as to improve the security and consistency of data in a NoSQL DB. With further research, we felt that we really could not make relevant contributions to existing solutions. So with the knowledge that we had compiled from our research on NoSQL databases and the concept of ledgers used in blockchains, our new aim was to provide a concise and clean platform to audit logs in a NoSQL database. In a real-world scenario, there will be multiple applications that contact the database to either make changes or retrieve some data. Servers generate a large number of events/logs (i.e. logging,) that contain information about their operation including errors, warnings, and users behavior. Typically, many clients have to extract data from database logs to analyze performance issues. Reviewing and formatting log data is time-consuming if not, more time consuming than analyzing formatted data using a visualization utility or machine learning algorithm. We wanted to eliminate the heavy lifting of formatting and capture all active operation-related data in a single collection for any client to be able to query. This is the problem that we aim to solve with this solution. The core component of our solution is to capture the data related to longer running queries. We stored data in the ledger in the MongoDB. MongoDB is installed in Amazon AWS EC2 instance (VM)[15]. we propose a solution with a ledger that captures all operations related to queries which are taking more time to run in MongoDB. For a simulation of a number of queries generated, in this solution, we used to stress test which inserts, update data for every second. We wrote a python script which can fetch important data into ledger about MongoDB long-running queries. we used mongo compass tool to visualize data in the ledger.

The output of above queries will provide operation IDs of the current running or queued operations in mongoDB. The output is analyzed and appended to ledger if the operation ID does not exist. Ledger is a collection in mongoDB that stores operation ID, timestamp, and a status flag as a document in the ledger collection. By default Status Flag is set to FALSE. The daemon process interval is set to 1 minute, so we collect and maintain the output of the run by updating the ledger every time. Executed operations can be determined by comparing the operation IDs that are present in the ledger against the Daemon process's result. If they are not present, it implies that the operation has been successfully executed and therefore the status flag for those operations is set to "TRUE". The ledger will hold operations for an interval of one hour.

IV. IMPLEMENTATION

This section discusses how ledger as a concept could be integrated into any NoSQL database. Ledger, as previously mentioned, can be thought of data that consists of a list of transactions in cryptocurrencies such as bitcoin. Similarly, keeping track of certain transactions that occur in NoSQL database (for example MongoDB)[11] provides great insight

for improving performance and analyzing what types of operations are performed. A ledger could easily answer questions like what and why is the database performing the way it is. Just like how customers currently create and maintain their own tables or collections to hold their data, we propose to create and maintain a ledger that would hold database related statistics added continuously. MongoDB is not only the most popular NoSQL databases in the current market but also the fastest growing database today. That is why MongoDB has been chosen as the database to do a test and present the model. As you can see in [Fig 2], one box consists of multiple collections that hold customer data and the other box consists of the ledger.

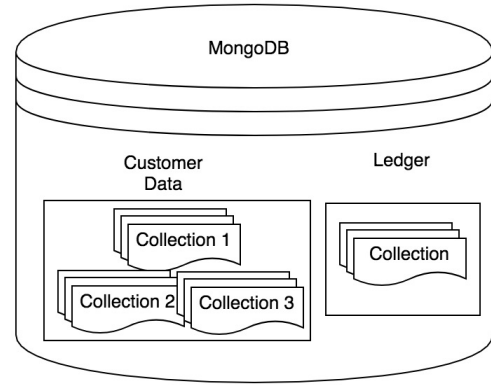


Fig. 2. Customer Data will consist of multiple collections and ledger is a single collection

A. Ledger and its parameters

Ledger is a single collection that is appended with data as a JSON object which is a document in MongoDB for a set interval. An operation related details are added to the ledger if and only if the operation is active, which means it is currently running, and if the operation is running longer than 50ms. The reason for maintaining a separate collection for auditing purposes provide benefits such as:

- Being able to query the ledger using MongoDB query syntax.
- The flexibility of adding or removing parameters or key-value pairs easily in the ledger per each customers requirement.
- Each document in the ledger is constructed by fetching real-time database statistics using commands such as current and server status. In other words, it is a central location to find all necessary auditing details.

The parameters or fields chosen to be added to the ledger comprises of the utmost important data that can generally help for analysis purpose. Below is a list of parameters and why they would be useful to keep track of [Fig 3]

- Operation ID: It is a unique id that each active MongoDB operation is associated with.

very quickly. But with further analysis, it was found that this was due to the inherent property of the working of the load balancing tool and our proposed system suffers no separate performance impact.

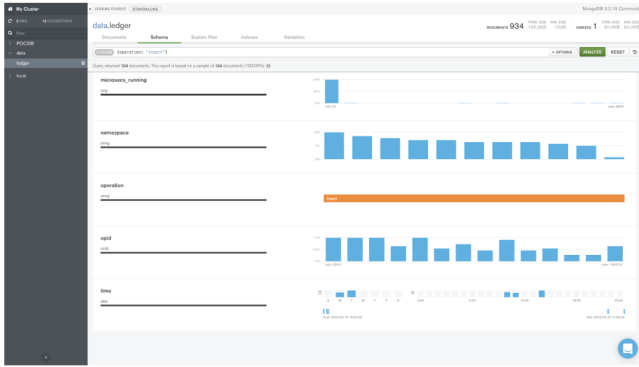


Fig. 6. Visualization of data

V. CONCLUSION

In our project, we tried to implement a solution which is to enhance the data storage about the queries which takes a longer time to execute in a systematic format. We used a ledger to store MongoDB log data. The customer can choose an appropriate time interval to insert into ledger (For example every second, 30 seconds, a minute or more). The customer has the flexibility to keep or remove ledger data depending on how much space available in their cloud environment. It is a concise version of typical NoSQL log entries. Ledger stores data related to the operation and can be useful for data analysis purpose. By using this solution in the data analysis, the service provider can take action in advance to prevent long-running query delay. The service provider can use machine learning algorithms to find common patterns or performance of each operation type which create long-running queries and take precautionary action to avoid delay in the queries. It can give the ability to monitor data in heavy traffic (For example Thanksgiving weekend) and Low Traffic time during the course of a year. This solution is advisable assuming space is not a constraint for a client. In a cloud environment, space constraint can be resolved by simply adding more machines (horizontally Scaling), however, if adding more machines leads to going over a clients budget, then this solution may not be advisable in that scenario.

VI. FUTURE WORK

This solution can be useful for the data analysis purpose. Using this data service provider can Use machine learning algorithm to analysis common pattern to find common patterns or performance of each operation type which create long-running queries and take precautionary action to avoid delay in the queries. It can be used for visualization of clustering data, error detection, graphical representation.

The solution of maintaining system logs in the form of a ledger can be extended to other NoSQL databases like Cassandra and Redis with a prediction of similar results.

VII. RELATED WORK

In this section,[5] we briefly review existing works on using Crypt-NoSQL, the first prototype to support execution of query over encrypted data on NoSQL databases with high performance. CryptDB utilized a proxy server to translate and rephrase the client's query before it goes into the database, making this query executable over encrypted contents. It applies several encryption schemes that are called SQL-aware to data, which means the execution of some specific queries over the encrypted contents will become possible by using their encryption schemes. However, CryptDB has many limitations like it require more storage space. CryptDB cannot perform server-side computations on values encrypted for different principals because the ciphertexts are encrypted with different keys. And it only has the worst security guarantee, and it requires too much space.

This paper is an approach which[8] Corda R3 enables the technologies by using distributed ledger or a blockchain technology. Corda platform in AWS which uses a shared smart contract to encapsulate the business logic of a transaction between organizations. This is all achieved with cryptographic technologies that provides an immutable record on a shared ledger system where data is ensured by validated users, which deals to enable strongest privacy assurance in the industry.

REFERENCES

- [1] A. Brinckman et al, "A Comparative Evaluation of Blockchain Systems for Application Sharing Using Containers," 2017 IEEE 13th International Conference on e-Science (e-Science), Auckland, 2017, pp. 490-497.
- [2] Z. Zheng, S. Xie, H. Dai, X. Chen and H. Wang, "An Overview of Blockchain Technology: Architecture, Consensus, and Future Trends," 2017 IEEE International Congress on Big Data (BigData Congress), Honolulu, HI, 2017, pp. 557-564.
- [3] S. Chickerur, A. Goudar and A. Kinnerkar, "Comparison of Relational Database with Document-Oriented Database (MongoDB) for Big Data Applications," 2015 8th International Conference on Advanced Software Engineering and Its Applications (ASEA), Jeju, 2015, pp. 41-47.
- [4] M. H. Shih and J. M. Chang, "Design and analysis of high performance crypt-NoSQL," 2017 IEEE Conference on Dependable and Secure Computing, Taipei, 2017, pp. 52-59.
- [5] M. Ahmadian, "Secure query processing in cloud NoSQL," 2017 IEEE International Conference on Consumer Electronics (ICCE), Las Vegas, NV, 2017, pp. 90-93. doi: 10.1109/ICCE.2017.7889242
- [6] L. Xu, L. Chen, Z. Gao, Y. Lu and W. Shi, "CoC: Secure Supply Chain Management System Based on Public Ledger," 2017 26th International Conference on Computer Communication and Networks (ICCCN), Vancouver, BC, 2017, pp. 1-6.
- [7] J. Garay, A. Kiayias, and N. Leonardos, *The bitcoin backbone protocol: Analysis and applications*, in Annual International Conference on the Theory and Applications of Cryptographic Techniques. Springer, 2015, pp. 281310.
- [8] <https://www.r3.com/blog/2017/12/05/r3s-corda-becomes-one-of-the-first-dlt-platforms-available-on-aws-marketplace/>
- [9] Building Enterprise-Grade Blockchain Databases with MongoDB white paper, 2017. <https://www.mongodb.com/white-papers>
- [10] A. Stanciu, "Blockchain Based Distributed Control System for Edge Computing," 2017 21st International Conference on Control Systems and Computer Science (CSCS), Bucharest, 2017, pp. 667-671.
- [11] MongoDB Documentaion <https://docs.mongodb.com/manual/tutorial/rotate-log-files/>
- [12] MongoDB Compass <https://www.mongodb.com/products/compass>

- [13] PyMongo Documentation <https://api.mongodb.com/python/current/>
- [14] Apache Jmeter <https://jmeter.apache.org/>
- [15] <https://docs.aws.amazon.com/quickstart/latest/mongodb/architecture.html>