

Prototype Big Data Archive in a Public Cloud

Group 56: Pathfinder of Big Data

Zhi Jiang, Isaac T Chan, Zhaoheng Wang

CS 461: Senior Capstone Fall 2016

Oregon State University

Abstract

OSU campuses generate data constantly from multiples sources, including computer labs, wireless usage, student devices, and many others. This quantity of data, also known as big data, can effectively represent all kinds of behaviors of students for information technology. For example, analysis can be run to determine common student behaviors in order to allocate OSU resources more effectively. Currently, the data is very difficult to manage because it is collected from multiple sources and is impossible to analyze. The data is neither stored in the same formats nor in the same locations, meaning it is inaccessible and useful information is unable to be extracted. Our goal for this project is to unify and organize the data onto the consistent cloud platform of Amazon Web Services, which additionally provides utilities to manage and analyze. To achieve this, we plan to have a working prototype at the Engineering Expo that demonstrates the value of analyzing OSU big data and how the cost-to-value of our Amazon cloud solution compares to locally-hosted hardware. Our prototype will allow OSU big data to be analyzed and eventually it can be scaled to analyze all the data that OSU collects.



CONTENTS

1	Isaac T Chan	3
1.1	Methods to measure performance metrics of database functionality	3
1.2	Methods of database security	3
1.3	Methods of user interaction with the system	4
2	Zhi Jiang	5
3	Zhaoheng Wang	5
3.1	The storage way for dealing with processed data	5
3.2	Programming language for achieving Database functionality	6
3.3	The visualization tool use to display the data	7

1 ISAAC T CHAN

1.1 Methods to measure performance metrics of database functionality

The following are the three best options for benchmarking and measuring performance of our implemented NoSQL database: YCSB, AWS Cloudwatch, and TPC-H. YCSB, or Yahoo! Cloud Serving Benchmark, is an open source framework for evaluating and comparing the performance of multiple types of data-serving systems[1]. AWS Cloudwatch is a built-in utility of AWS and can be used to collect and track metrics. TPC-H benchmark consists of a suite of business oriented ad-hoc queries and concurrent data modifications[2]. With the benchmarking and performance measurement utility, we hope to obtain a baseline for our database performance and examine how various data and query loads compare to the baseline. We will be evaluating operation speed of operations such as database inserts, updates, and reads.

	Inserts	Updates	Reads	Visualization	Extra Notes
YCSB	Yes	Yes	Yes	Raw data, can be plotted	Open-source utility means we can customize tests to fit our use-case
AWS Cloudwatch	Yes	Yes	Yes	On AWS UI and also provides raw data	Built-in utility eliminates complexity of implementation
TPC-H	Yes	Yes	Yes	Raw data	Lack of customizable tests

YCSB is a very customizable, open-source utility that can produce relevant and informational metrics for our database. It has a wide user-base and should be easy to implement.

AWS Cloudwatch is a built-in tool that can deliver relevant metrics and should work well with our database on AWS. It also has customizable metrics, which we would implement using AWS CLI (command line interface).

Finally, TPC-H is an enterprise-grade option, mostly used for server-production companies to measure how their products compare to alternatives. There is a lack of customization and a lack of a community for troubleshooting. Documentation is minimal. Implementing TPC-H is likely to be a challenge.

We will select AWS Cloudwatch as the best option. Cloudwatch can also use customized metrics tests which is important in order to know metrics to fit our use case. Although YCSB may be more easily customizable, the lack of installation makes Cloudwatch the best option.

1.2 Methods of database security

Traditionally, NoSQL databases have minimal security. In most implementations, the only security is to allow access to the database via trusted machines. However, relying solely on the network is almost certainly an invitation for a breach to sensitive information. NoSQL databases also cannot use external encryption tools, such as LDAP or Kerberos. Our best options for database security is AWSs user authentication policies, encrypting sensitive data fields, and using sufficient input validation to avoid injection attacks. With methods of database security, we hope to improve security of our database by restricting access and preventing malicious utilization of the data. We will evaluate database security by ensuring only approved user access, whether or not sensitive data is encrypted, and resistance to injection attacks.

	User restricted access	Sensitive data encryption	Resistance to injection attacks
AWS user authentication	Yes	No	Yes
Sensitive data encryption	No	Yes	No
Input validation	No	No	Yes

AWS offers a useful utility for user authentication, where users of the system can be granted different levels of access. Naturally, with authentication there is a built-in resistance to injection attacks because hopefully authenticated users are less prone to malicious intent.

There will be sensitive data in our database, most notably user-identification information, such as student IDs. This will be encrypted prior to given to us, and will most likely remain encrypted as the database is implemented and real data is inserted.

Finally, input validation is a minimal concern if we decide to implement NoSQL using AWS DynamoDB. DynamoDB does not support multiple actions with a single command, removing the risk of injection attacks. If we choose to implement the database using another tool, then there will need to be test cases to identify and ignore injection attacks.

As shown in the comparison table above, no one method can cover the security of the table. We will need to implement all three methods. Due to the lack of external tool support by NoSQL databases, we must resort to using modular security methods.

1.3 Methods of user interaction with the system

There will only be one type of user that interacts with the system - OSU staff that are able to manage and analyze the data. There are many analysis tools that can be used in conjunction with AWS. Most notably, Amazon Machine Learning (AML), Amazon EMR, and Amazon QuickSight. AML is a service that provides easy-to-use machine learning technology. Amazon EMR is a comprehensive utility that allows users to interact with databases, data warehouses, and customize their analysis. Amazon QuickSight is a fast business intelligence service that allows users to visualize data and provide responsive analysis. Our goal is to shape our implementation database to work with a utility that allows users to interact with our system. The ideal utility would provide tools to manage the data, provide analysis, machine learning support, and visualization.

	Data Management	Analysis	Machine Learning	Visualization
AML	No	No	Yes	No
EMR	Yes	Yes	Yes	No
QuickSight	No	Yes	No	Yes

AML is a very specific utility that only really offers machine learning analysis in a simple interface. It does contain much else within the tool, but the user can visualize the results from machine learning models with Amazon Cloudwatch.

Amazon EMR is a comprehensive tool that can manage data and provide analysis, through conventional queries or machine learning technology. It is more difficult to use, and does not provide native visualization.

Lastly, Amazon QuickSight is a business intelligence tool that can deliver fast analysis and boasts a very attractive visualization tool. However, more in-depth analysis and methods of data management are unavailable.

Amazon EMR is the technology that we will primarily be considering. After the conclusion of our project, maintaining the database, as well as analysis, is critical to the survival of our prototype. Our chosen technology must support all of these, and Amazon EMR is the most comprehensive tool. However, in the end, all of these technologies can be used in conjunction with our final big data prototype, but it is important to have a primary tool in order to ensure database maintenance and analysis is continuable and extensible after the Expo.

2 ZHI JIANG

3 ZHAOHENG WANG

3.1 The storage way for dealing with processed data

The goal for this part is to figure which is the optimal option of storage way for dealing with processed data.our product requires to use the NoSQL database however it may not be the optimal option for storage. As a result, we will first compare other storage way with database to check whether database is the optimal choice. If the database is the optimal choice, we will compare the Sql with NoSQL databases. After that, it is necessary to figure out which type of SQL or NoSQL is the best option.

There are many ways for storing data such as Database, Files, Cookies and other ways. we would like to start with cookie storage. Usually, cookies are used to store tiny bits of data. These data are very small only below 4 kilobytes per domain[3]. Besides. the cookies are pass the data by request which is not suitable for our product. Therefore, the cookie storing is not a good choice for our product. Another option is using Files such as XML to store the data. The data will be very large quality in our product thus there will be a lots of XML files for storing. Because of this, the management will be complex by using XML files. So the file storing is also not the suitable option for storage. The next option is using database. The database is used to store the data. The data could be managing, retrieving and organizing in the database[4].The database gives promote accessibility for the data and the data could be easily accessed by using query. Besides, it makes the data more security and reducing the cost for data insert, storing[5]. As a result, the database is the optimal choice for our product.

The database could divide into two types in general: the Relationship database and Non-relational database.The Relational database is usually represent the data base on the table however the Non-relational database use dynamic schema for data. When dealing with multiple type of data, the Relationship database is not the optimal choice because it is not the optimal choice if the data is storing in hierarchical. However, the Non-relational database is the optimal choice for dealing large quantity of data which stores in hierarchical. For scalability, the Relationship database could increase its scalability by promoting the power of hardware however the Non- relational database will increase its scalability

by reducing the load. After comparing with the Relationship database and Non-relational database, we find that the Non-relational database Is the optimal choice for our product because the Non-relational database is suitable for dealing with large quantity of data which stores in hierarchical[6].

Tool	Platform	Security for the data	Speed for loading data	Features	cost for the tool
SimpleDB	Amazon web services	Yes	Fast	High availability and simple to use	Low
MongoDB	Cross-platform	Yes	Fast	Document database, high performance and high availability	Low
DynamoDB	Amazon web services	Yes	Fast	Support key-value model,High availability, free-text search,flexible database schema	25GB for free and the cost is low
Cloud Datastore	Google cloud platform	Does not mention on the product page	Fast	High availability and high scalability	Low

The Non-relational database we want to evaluating are: SimpleDB, MongoDB, DynamoDB, Cloud Datastore. we generate several criteria for evaluating these NoSQL databases. These criteria will evaluating the NoSQL database in various features such as cost,platform support, security,loading speed and many other features.

SimpleDB is offered by Amazon web services. It is security for the data and the speed for loading is fast. Besides, the cost for it is low.

MongoDB is supported by cross-platform. It is also security for the data and the speed for loading is fast. It has the high availability and high performance. The cost for it is also low.

DynamoDB is supported by Amazon web services. It is more security for the data than others and the speed for loading is fast. It has many features for example, the free-text search will make the information searching more convenient[7]. Besides, it has high availability and flexible database schema.

Cloud Datastore is supported by Google cloud platform.it is fast and it has high availability and high scalability.

We will select the DynamoDB for storing because our client requires to use the Amazon web services as the platform. Besides, the DynamoDB has more features which is suitable for our product. For example, the free-text search allows to search the information easier and flexible schema make the schema easy to development. Furthermore, the cost for DynamoDB is very low.

3.2 Programming language for achieving Database functionality

There are many options for programming language such as java, python, php. Each programming language will have different features. Our goal for this part is to figure out the suitable language for our product and avoid using many

different languages. Because the more languages we use the more mistake we will might have.we generate several criteria for evaluating different language such as APIs, testability,security and tool support .

language	API for inserting	API for updating	API for listing table	Testability	Security	DynamoDB support
Java	Yes	Yes	Yes	Testable	Secure	Yes
php	Yes	Yes	Yes	Testable	Normal	Yes
python	No	No	No	Testable	Secure	Yes

Java is the optimal choice for achieving Database functionality. Here are the reason as following. The document of Amazon DynamoDB provide the APIs for basic functionality in java such as inserting data, updating data and listing table. These APIs make the database functionality achieving more easily. Besides, Java is testable and more secure than the php. Another point is most of tool for our product is using Java.If we use the language other than java, the process will become complex because we also need to figure out the translation way between java and that language. This will also make the test process become much more complex. therefore , java is the suitable language for achieving Database functionality.

3.3 The visualization tool use to display the data

There are various visualization tools could be using to make the data visualization. Each visualization tool has different features. Our goal for this part is to figure out the optimal choice of visualization tool for displaying the data. The visualization tool we choose to evaluate are Tableau,QuickSight and FusionChart. we generate several criteria for evaluating them such as platform, speed, features and cost.

Visualization Tool	Platform	Speed	Features	Cost
Tableau	Tableau Online	Fast	Allows cross database, beautiful design	Low
QuickSight	Amazon Web Services	Fast	High accessibility, Get answer fast, Easy share business insight, Smart Visualizations	Low
FusionChart	No platform	Fast	Controllable for chart	Normal

Tableau could generate the graph fast and the cost of it is not high. The QuickSight is supported by Amazon web services. It has high accessibility which allows access data from multiple source. Besides, it could share the business insight in security way.Another important feature for QuickSight is it could generate visualizations very fast for very large quantity of data.Furthermore, the cost of it is low[8]. FusionChart does not require the platform to support and it could generate the graph fast. The cost of Fusionchart is expensive than Tableau and QuickSight.

Comparing with Tableau and FusionChart, the QuickSight is the optimal choice for our product. The QuickSight is supported by Amazon web services which fits the requirement of platform. Besides, it has strong accessibility which allows to communicate with other data service on Amazon web services such as Amazon DynamoDB easily.Furthermore, it has fast speed for generate the large data set which fits our purpose. Therefore, the QuickSight will be the optimal option as the visualization tool use to display the data.

REFERENCES

- [1] G. Kamat, "YCSB, the Open Standard for NoSQL Benchmarking, Joins Cloudera Labs", *Cloudera*, 2015. [Accessed: 14- Nov- 2016].
- [2] "TPC-H - Homepage", *Tpc.org*, 2016. [Online]. Available: <http://www.tpc.org/tpch/>. [Accessed: 14- Nov- 2016].
- [3] "Why you probably shouldn't use cookies to store session data", *Wonko.com*, 2016. [Online]. Available: <http://wonko.com/post/why-you-probably-shouldnt-use-cookies-to-store-session-data>. [Accessed: 14- Nov- 2016].
- [4] "Benefits of Using a Database", *Work.chron.com*, 2016. [Online]. Available: <http://work.chron.com/benefits-using-database-3792.html>. [Accessed: 14- Nov- 2016].
- [5] "Advantages of Database", *Ecomputernotes.com*, 2016. [Online]. Available: <http://ecomputernotes.com/fundamental/what-is-a-database/advantages-of-database>. [Accessed: 14- Nov- 2016].
- [6] "SQL vs NoSQL Database Differences Explained with few Example DB", *Thegeekstuff.com*, 2016. [Online]. Available: http://www.thegeekstuff.com/2014/01/sql-vs-nosql-db/?utm_source=tuicool. [Accessed: 14- Nov- 2016].
- [7] "Amazon DynamoDB Product Details", *Amazon Web Services, Inc.*, 2016. [Online]. Available: https://aws.amazon.com/dynamodb/details/?nc1=h_ls. [Accessed: 14- Nov- 2016].
- [8] "Amazon QuickSight Business Intelligence Software", *Amazon Web Services, Inc.*, 2016. [Online]. Available: <https://aws.amazon.com/quicksight/>. [Accessed: 14- Nov- 2016].