

# Survey on Architecture for Implementing Big Data in Cloud

K. Uma\*, A. Gowri Shankar, V. Rajesh Kumar

Department of IT, VIT University, Vellore – 632 014, India

\*Corresponding author: E-Mail: drumakphd@gmail.com

## ABSTRACT

The rapid advancement of the Internet of Things and Electronic Commerce that entered the era of big data. The characteristics, such as great amount and heterogeneity of big data bring the dare to the storage and analytics. As there are many architecture in the field for implementing big data in cloud computing literature. This paper open a comparison of universal storage architecture, smash and Hadoop YARN with Apache Spark for big data in cloud environment. In recent times, big data has become a trendy research topic and brought about a scope of new challenges that must be tackled to sustain many commercial and research demands. Tackling these big data issues requires capabilities not characteristically found in common Cloud platforms. This includes a distributed file system for capturing and storing data; a high performance computing engine able to process such huge quantities of data; a reliable database system able to optimize the indexing and querying of the data, and geospatial capabilities to imagine the resultant analyzed data. With the high-scalability cloud technologies, Hadoop and Spark, the proposed system architecture is first implemented successfully and resourcefully. Experimental results illustrate the effectiveness and efficiency of the proposed system services via an advanced web technology. In addition, some experimental results signify that the computing ability of Spark is better than that of Hadoop.

**KEY WORDS:** Cloud, big data, data analytics, traffic analysis, data model, Apache Spark, NoSQL, Hadoop YARN, HBase.

## 1. INTRODUCTION

Today in this world there is plenty of information to be stored in every field: Social media, Ecommerce websites, Banking, Communication etc. To stock up these data and to utilize them properly we go for big data storage technology but yet there are some troubles when accessing these information or data in a shared environment, so we introduced information shared over the internet using cloud so that we can effortlessly share and retrieve data from any part of the world. In this paper we have discussed the various architectures for implementing big data in cloud. The presented storage architectures can support multiple data models, including all kinds of relational data and non-relational heterogeneous data, by isolating nodes in cloud storage centre into several clusters, each of which stores data with special model such as key value model and document model (Atzori, 2010). Furthermore, the architecture provides users with unified storage interface and query interface. A astounding amount of data across all walks of life is being generated: social media data, scientific data, data collected by government, by industry and across all research endeavours. This is often referred to as the data deluge and big data is now a term in common parlance. In order to process, analyse and derive knowledge from this ever-growing amount of data, large-scale and above all scalable infrastructures are vital. For many domains, the challenges in developing such infrastructures are one of scale, e.g. the volume of data is the challenge that must be tackled; for other domains it is the speed (velocity) at which data itself is produced that is the issue to be tackled; for other domains it is the accuracy, authenticity and provenance of data that is essential (veracity). The transport arena is one example that has much to benefit from big data capabilities in allowing to process voluminous amounts of data that is created in real time and in vast quantities. Some disciplines demand all of these big data capabilities. Transport and traffic flows more generally are areas that demand the ability to cope with large-scale data sets that are produced in near real-time and have significant demands on the accuracy and validation of information.

## 2. RELATED WORKS

For amorphous data set NoSQL data stand is proposed by Sun (2010). The node parting for normalisation is planned by (Dobre nad Xhafa, 2014). The data which is used within the universal structural design need not just before be a structured data or uniform data however may existing data the design uses data analysis methods to make-up them with the aid of NoSQL record (Kortuem, 2010). Straight forward understanding of layers by in-between complete layers into modules. The separation of notes are done by, how much can a node store that is the space in them which was projected by Jlain Zou. Apache Spark is more faster than Hadoop MapReduce which was used to compare these architectures which was projected by (Brugmann, 2014). RDDs allow upturn of failed nodes by re-computation by means of the DAG ancestor is proposed by Sinnott (2015). The utilize of Geo based software makes it extraordinary for traffic based function. As there are numerous data to accumulate in cloud HDFS is proposed by Gereon Frahling (2006). NoSQL to hoard big data since NoSQL provides an instrument for storage and recovery of data that is improved than the tabular relations old in relational databases. Purpose can be performed by Spark in three modes: Spark Standalone, YARN, and mesos.

**Comparative Analysis:**

**A Universal Storage Architecture for Big Data in Cloud Environment:** This architecture can hold up different data models, it includes every type of relational data and non-relational assorted data called NoSQL data (Welbourne, 2009). The architecture is based on isolating nodes in cloud cargo space centre into a number of clusters, all of which supplies data through special model such as rate model in addition to document model.

**Data analysis layer:** The cloud surroundings nodes are alienated into altered types of clusters in the middle of which the strongest solitary is considers while statistics analysis layer and to facilitate, data store section and data enquiry module are incorporated.

The enormous varied data that the client submitted towards the data layup module is normalized resting on a frequent basis as well as then is stored within the cluster (QIN Xiong-Pai, 2012).

**Data storage:** The construction uses diverse clusters to amass different information models by isolating the notes of the cloud storeroom. Assume with the purpose of the data core needs shore up  $n$  kinds of fact models and afterwards the cloud nodes will be at odds into  $n+1$  cluster (Anderson, 2010).

**Disadvantage of Universal Architecture:**

- The construction only considers the cluster through strongest computing supremacy is only certain so previous clusters with fewer computation powers are not engaged into account.
- The heterogeneous records given by the users are normalized and then stored hooked on unlike notes or clusters, the time occupied to do this work is added and the supply is also high (Kranz, 2010).

**Smash Architecture for Big Data in Cloud Environment:** SMASH architecture was urbanized to attempt issues like fetching the statistics, pointed the data, storing the facts and validating the information in big data transportation. The architecture largely concentrates on circulated data storage other than to facilitate it, it contains Hadoop Distributed File System (HDFS). It uses Apache Spark in its place of Hadoop MapReduce module. This architecture is chiefly for traffic statistics so the architecture consists of machinery like GeoMesa as well as GeoServer are extra to the stack (Hall, 1980).

**Performance of Data loading:** SMASH contains Apache Spark for collecting the actual time data so that it knows how to be stored in a circulated fashion and to decrease the load of the structure. It used Resilient Distributed Dataset (RDD) as an alternative of Hadoop MapReduce technology (Jianting Zhang, 2014). A judgment chart is given for both Hadoop MapReduce plus Spark technology on definite benchmarks on which we can see the act of Spark more than Hadoop MR.

**Performance of Data Processing:** When comparing the performance with HDFS, HDFS stores facts on the HDFS system whereas Apache Spark stores statistics inside RDD memory rather than in the recollection disk. Only when the remembrance exceeds its restriction it spills data on the material disk by this Spark is a good deal faster than MapReduce by dropping the amount of I/O operations (Yin, 2013).

**Performance of Data Backup:** RDDs agree to mending of failed nodes by re-computation with the DAG predecessor.

**Disadvantage of SMASH Architecture:**

- The exercise of spark in SMASH has exacting memory requirements, therefore in a Spark group the memory must be at slightest large data we need to process, as the data has to robust into the memory for additional operation.
- Block sizes are more vital when there is tiny dataset to route then there will be more recollection wastage.

**Novel architecture for big data in cloud environment:** Novel architecture was urbanized to solve issues similar to long time dispensation, inaccurate results, inadequacy of the structure etc. It uses spark requires a bunch manager and a circulated storage system. For cluster director it uses Hadoop YARN method and for distributed storage space HDFS is used (Angles, 2008).

**Performance of Data loading:** The information is loaded in the shape of clusters plus the clusters are divorced into notes by analyzing the data here in it. NoSQL data base is old in big data as there is together structured and unstructured data here in it. And in spread storage Hadoop is used (Yan, 2015).

**Performance of Data Processing:** After loading of assets into function the Spark practice can be executed by in between the application into collect mode and consumer mode. As a replacement for distributing JOB to every node like MapReduce it uses IN recollection of the apache spark to accomplish high-speed since the sparks IN reminiscence is 100 times quicker for several applications (Sinnott, 2015).

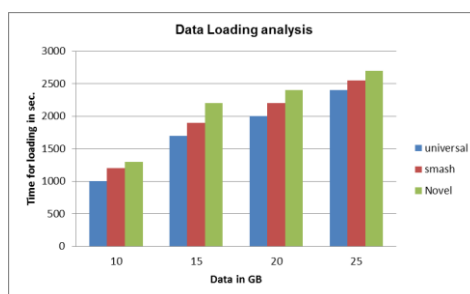
**Performance of Data Backup:** As the structural design uses only NoSQL and Spark technique for dispensation the data there is no extraordinary process like RDD for data backup, as it supplies in cloud data is forever secure but, in holder if there is any thrashing in data then it will be a trouble meant for the application (Chang, 2008).

**Disadvantage of Novel Architecture:** There is no chief disadvantage in this architecture; only thing is it suits only for traffic linked applications since it's based lying on monitoring the traffic plus coverage the data to the users. One more thing is that the scheme should have endorsement option like former architectures.

**Comparison analysis results:** In this paper to show which of the architecture is best among the three architectures, some of the data are collected and compared with respect to loading of data, backing up the data and their query processing time. Values are tabulated in a table.1 for comparative graphs of data verses time to load them is charted and presented in fig.1 and for checking query processing speed of each architecture values are tabulated in table.2 and graph is drawn for the values in table for querying the data fig.2 and for data backup values are tabulated in table.3 and for that fig.3, graph is drawn.

**Table.1. Data in GB verses time for loading the data**

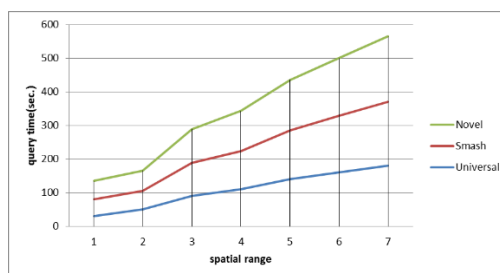
Data in GB	Universal Time in sec.	Smash Time in sec.	Novel
10	1000	1200	1300
15	1700	1900	2200
20	2000	2200	2400
25	2400	2550	2700
30	2650	2800	3000



**Figure.1. Comparison of architectures based on their loading capacity of data**

**Table.2. Query verses Time for processing the Query**

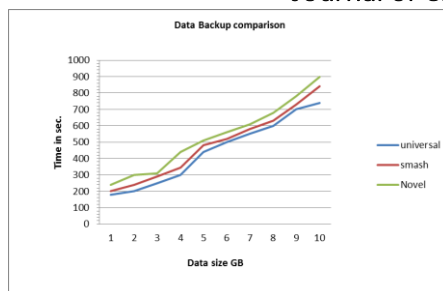
Spatial Range	Universal Query time in sec	Smash Query time in sec	Novel Query time in sec
1	30	50	55
2	50	56	60
3	90	98	101
4	110	114	119
5	140	146	150
6	160	168	172
7	180	190	195



**Figure.2. comparing the architecture by the speed of their querying time**

**Table.3. Data verses Time for backup**

Data size GB	Universal	Smash	Novel
1	180	200	240
2	200	240	300
3	250	290	310
4	300	345	440
5	440	480	510
6	500	520	560
7	550	580	610
8	600	630	680
9	700	730	780
10	740	840	900



**Figure.3. Comparison of architectures based on their capacity to backup data**

#### 4. CONCLUSION

In this paper we have compared three architectures and as of effect, when seeing the graphs it is clearly concluded that universal design is best as the loading time less and the backup point is also little bit higher than the other architectures so if it is ecommerce sites then the best to use is universal architecture or any other data dispensation application then universal architecture is best. In case if you want to go for any traffic manage or any data monitoring purpose then SMASH architecture is just right as it includes more rapid processing of data than novel architecture and the processing time for queries is lesser than novel architecture and loading time for higher data is also less compared to novel so smash is best when implementing big data on cloud.

#### REFERENCES

- Anderson C, Lehnardt J, Slater N, Couch DB, the Definitive Guide, the Definitive Guide, O'Reilly Media, 2010.
- Angles R, Gutierrez C, Survey of graph database models, ACM Computing Surveys (CSUR), 40 (1), 2008, 1.
- Atzori L, Iera A, Morabito G, The internet of things - A survey, Computer Networks, 54 (15), 2010, 2787-2805.
- Baoyun W, Review on internet of things, Journal of Electronic Measurement and Instrument, 23 (12), 2009, 1-7.
- Brugmann J, Schreckenber M and Luther W, A verifiable simulation model for real-world microscopic traffic simulations, Simulation Modelling Practice and Theory, 48, 2014, 58-92.
- Chang F, Dean J, Ghemawat S, Bigtable, A distributed storage system for structured data, ACM Transactions on Computer Systems (TOCS), 26 (2), 2008, 4.
- Clement S and Anderson J, Traffic signal timing determination, the Cabal model, 1997.
- Gereon Frahling and Christian Sohler, A fast k-means implementation using coresets, In *Proceedings of the 22nd ACM Symposium on Computational Geometry, Sedona, Arizona, USA, June 5-7, 2006*, 135-143.
- Intelligent services for big data science and its Applications by Dobre C and Xhafa F, 37 (0), 2014, 267-281.
- Hall M and Willumsen L, Saturn-a simulation assignment model for the evaluation of traffic management schemes, Traffic Engineering & Control, 21 (4), 1980.
- Jianting Zhang S, You and Gruenwald L, High-performance spatial query processing on big taxi trip data using gpgpus, In Big Data (BigData Congress), 2014 IEEE International Congress, 2014, 72-79.
- Kortuem G, Kawsar F, Fitton D, Smart objects as building blocks for the internet of things, Internet Computing, IEEE, 14 (1), 2010, 44-51.
- Kranz M, Holleis P, Schmidt A, Embedded interaction, interacting with the internet of things, Internet Computing, IEEE, 14 (2), 2010, 46-53.
- Monteil J, Nantes A, Billot R, Sau J, Microscopic cooperative traffic flow, calibration and simulation based on a next generation simulation dataset, IET Intelligent Transport Systems, 8 (6), 2014, 519-525.
- QIN Xiong-Pai, Wang Hui-Ju, DU Xiao-Yong, WANG Shan, Big Data Analysis, Competition and Symbiosis of RDBMS and MapReduce, 23 (1), 2012, 32-45.
- Sinnott R.O, Morandini L & Wu S, Smash, A Cloud-based Architecture for Big Data Processing and Visualization of Traffic Data, International Conference on Data Science and Data Intensive Systems, IEEE, 2015, 53-60.
- Sun Q, Liu J, Li S, Internet of Things, Summarize on Concepts, Architecture and Key Technology Problem, Journal of Beijing University of Posts and Telecommunications, 3, 2010, 003.

Welbourne E, Battle L, Cole G, Building the internet of things using RFID, the RFID ecosystem experience, Internet Computing, IEEE, 13 (3), 2009, 48-55.

Yan Y.Z, Liu R.H, Yang C.T & Chen S.T, Cloud City Traffic State Assessment System Using a Novel Architecture of Big Data, International Conference on Cloud Computing and Big Data (CCBD) IEEE, 2015, 252-259.

Yin D and Qiu T.Z, Compatibility analysis of macroscopic and microscopic traffic simulation modeling, Canadian Journal of Civil Engineering, 40 (7), 2013, 613-622.