



## Big Data Analytics: A Survey

<sup>1</sup>K.Saraswathi, and <sup>2</sup>S.Logeshwaran

<sup>1</sup>Assistant Professor, Department of Electronics and Instrumentation Engineering, SCSVMV, Kanchipuram, India

<sup>2</sup>UG Scholar, Department of Computer science and Engineering, PERI Institute of Technology, Chennai, India

### ABSTRACT

It's veritably delicate to storing, managing and recycling huge quantum of data. The term 'Big Data' describes colorful ways and technologies to store, distribute, manage and dissect huge quantum of data with different structures. Big data consists of structured, unshaped or semi-structured data so there's problems do regarding incapacity of conventional data operation styles. To reuse these huge quantities of data in an affordable and effective way, community is used. Big Data is a data which is in large quantum and having complexity in it and this complexity bear new armature, ways, algorithms, and analytics to manage it and prize knowledge from it. Hadoop is a frame for recycling large quantum of data and provides better storehouse capacity for large datasets and performs resemblant processing of big data that gives better computational power to all the tasks. It works in batch processing mode and Hadoop is the core platform for structuring Big Data, it also solves the problem of making it useful for analytics purposes. In this paper, we give a brief overview of big data operation involving hadoop and highlight exploration sweets and the challenges to big data.

Indicator Terms Big Data, Hadoop, Map Reduce, HDFS, Hadoop Component.

### 1. INTRODUCTION:

#### 1.1. Big Data: Definition

Big data is a term used to describe the exponential growth and vacuity of data, having structured, unshaped and semi-structured data, whose size ( volume), complexity (variability), and rate of growth ( haste) make them delicate or indeed insolvable to be managed and anatomized using conventional software tools and technologies. When the quantum of data to be increases than the time to produce results is also increased. Recaptured data from big data is still a complex and time consuming approach. Big data provides tremendous openings for enterprise information operation and decision timber. In the recent study big data isn't only limited to business requirements but also helps in exploration and scientific issues.

The Big Data problem is characterized by the 3V features

Volume-a huge quantum of data, Volume of big data can be measured in terms or several megabytes, gigabytes, terabytes or petabytes.

Haste- a high data ingestion rate or the speed with which the data can be anatomized. Variety-a blend of structured data, semi-structured data, and unshaped data.

These 3V features gives a challenge to data processing systems since these systems cannot moreover gauge to the huge data volume in a cost-effective way or fail to handle data with variety of types. The results to the Big Data problem are largely grounded on the Map Reduce framework (9) and its open source perpetration Hadoop. Although Hadoop handles the data volume challenge successfully. Hadoop is the open source software innovated by Apache and its Linux grounded software. It's used by notorious websites like Google, Yahoo, Facebook, Amazon and numerous further. Hadoop is a frame for recycling large quantum of data and provides better storehouse capacity for large datasets and performs resemblant processing of big data that gives better computational power to all the tasks. It works in batch processing mode and having two major factors HDFS (Hadoop Distributed Train System) (12) for huge data storehouse and Map Reduce for recycling huge quantum of datasets. When the data size is increased it produce problems to being algorithms to manage that so then main problem is to store and reuse that huge quantum of data and this problem is break by hadoop because it store and process huge quantum of data in lower time.

*Hadoop:*

Hadoop is an open-source software framework used for distributed storage and processing of big data using the Map Reduce programming model.

Modules present in Hadoop are designed with a fundamental assumption that hardware failures are common occurrences and should be automatically handled by the framework. The core of hadoop consists of two parts the storage part and processing part.

a) Storage part Storage part of hadoop is HDFS (Hadoop distributed train system) which stores huge quantum of data with high degree of outturn and this huge data is stored in form of clusters.

b) Processing part of Hadoop is MapReduce which is a software frame which process large quantum of data in the form of clusters.

Hadoop distribute clusters to the knot so that they reuse parallel and this approach also takes advantage of data position This allows the dataset to be reused briskly and more efficiently which make it a more conventional supercomputer armature which work on a resembling train system where calculation and data are distributed via high- speed networking.

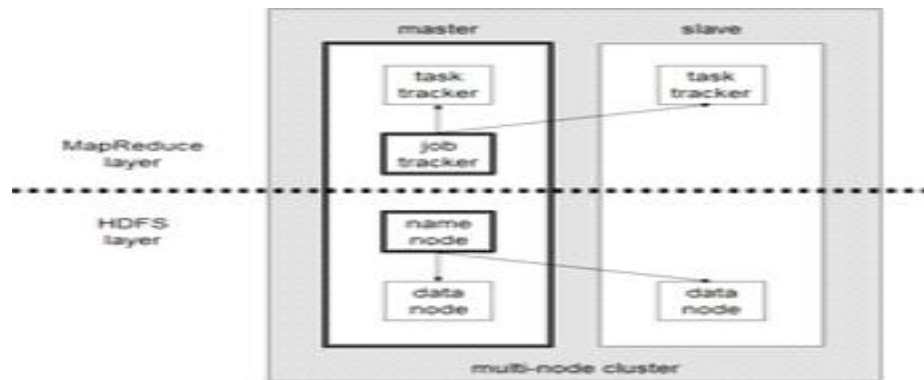


Fig.1.1.Hadop architecture

As depicted in Fig. 1.1, a tiny Hadoop cluster with a single master node and many slave nodes. A Task Tracker, Job Tracker, Name Node, and Data Node make up the master node [14], whilst the slave or worker node serves as both a Data Node and a Task Tracker.

1.2. Hadoop Distributed File System (HDFS): Hadoop Distributed File System (HDFS) is the storage component that stores a large number of structured, unstructured, and seminar-structured data. HDFS is a file system that is based on Java. HDFS is a dependable and easy-to-manage file system. Fast availability, load balancing, security, flexible access, fault tolerance, easy management, and high data throughputs are just a few of the benefits. It allows data to be processed in parallel. The architecture of HDFS is master/slave. [23]

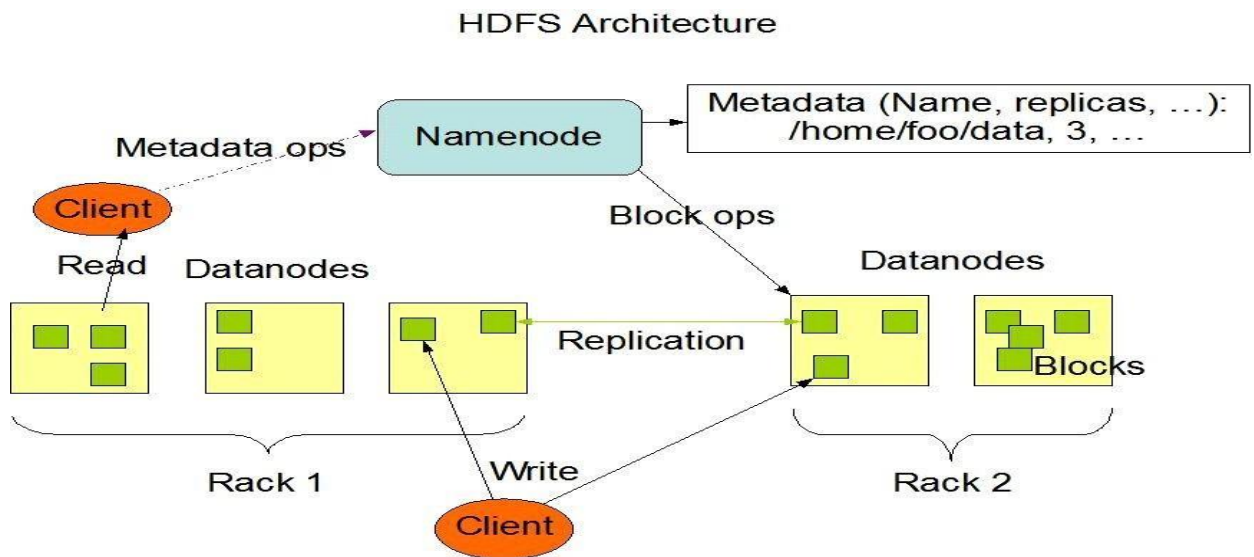


Fig. 1.2. HDFS Architecture

Affair. The reduce job takes the affair of the chart task as the input and combines them to lower set of tuples (reduces the large dataset into a lower one) grounded on the metamorphoses and color full logic. The advantage of Map Reduce is that it's easy to gauge data processing over multiple computing bumps.



Simon Fong, Raymond Wong, and Athanasios V. Vasilakis [25]	2	0	1	5	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Yanhao Huang and Xiaomin Zhou [22]	2	0	1	5	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Marcu Vazcani, Peter Hunter, and Rod Hesse [21]	2	0	1	5	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Alan Evans, In-Veget Joseph, et al. [20]	2	0	1	5	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Syed Akhter Hossain [29]	2	0	1	5	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Xiao Wen Chen AND Xiaoping Lin [30]	2	0	1	4	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Mattias Beck, Zhou Kanwei, Li Shuai, Li Nifeng [32]	2	0	1	4	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Suman Arora, Dr. Madhu Goel [7]	2	0	1	4	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Chang Lin, Junyi Chen, Chi Yang, Xujin Zhang, and Ramesh Babarao Kotagiri [16]	2	0	1	4	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Shifeng Fang, Li Da, Yao, Yinqiang Zhu, Jueheng Han, Huan Pei, Jianwei Yan, and Zhibin Liu [17]	2	0	1	4	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Daichi Takashi, Hiroki Nishiyama, Naoki Kato, and Ryo Mizu [18]	2	0	1	4	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Andrea Marzoni, Arianna Degliani, Riccardo Bellazzi, Paolo Garbelli [27]	2	0	1	3	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Kongqi Qin, and Xiaoyu Zhou [4]	2	0	1	3	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Balkeesh Varma [6]	2	0	1	3	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Daniel Waméke [15]	2	0	1	1	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Jasmin Kemovic, Denis Masic [13]	2	0	1	0	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Mengjie Zhou, Hongbin and Ming Zhou [14]	2	0	1	0	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Leonardo Nemeque Barros, Roberto André Nogueira, et al. [5]	2	0	1	0	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
El Sharbat, Xu Yu, Jiao Feng, Li Guozhen, PEI Anping [12]	2	0	0	9	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Hui Fang, Ming Yang, Ruiqing Yang [11]	2	0	0	7	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Seemal Mehta, Vijay K. Viswanath [10]	2	0	0	6	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
John H. Pan, Chang F. Qiu, and May D. Wang [9]	2	0	0	4	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.
Sushant Goel, Hemant Khand, David Tandi [8]	2	0	0	3	Research on the design of a distributed system for the management of a large number of users in a cloud environment.	Research on the design of a distributed system for the management of a large number of users in a cloud environment.

### 3. CHALLENGES:

Because big data is such a large volume of data, it poses a number of obstacles in terms of management, storage, scheduling, security, and processing. For starters, effective online analysis necessitates data preparation, efficiently distributed storage, and search, all of which necessitate excellent data mining tools. The efficient processing of a large data stream is a major challenge that requires a variety of programming models. Second, the scheduling technique should be intelligent enough to respond in real time to a changing environment. Third, data integration necessitates the development of new protocols and interfaces capable of managing structured, semi-structured, and unstructured data. The fourth point is user interaction and visualization. There are numerous research problems in big data visualization, therefore more efficient real-time visualization techniques are required.

Furthermore, security and privacy are important considerations.

---

## CONCLUSION:

In recent years, a survey of various big data techniques has been presented. It has been discovered that Big Data solutions are mostly built on the Map Reduce architecture and its open source implementation Hadoop. Hadoop successfully solves the data volume difficulty. For task scheduling in Hadoop, big data management encompasses a variety of tools, strategies, and algorithms. This paper is beneficial to a newcomer who want to pursue a profession in the field of big data.

## FUTURE DIRECTION:

This work could be expanded by creating a new job scheduling algorithm that takes into account all of the variables that influence performance. Second, the user profile (similar people) and usage profile (invoked services) should be taken into consideration, as well as some related collaborative filtering strategies to combine with our service selection methodology.

---

## REFERENCES:

- [1] Sreedhar C.N, Kasiviswanath, P. Chenna Reddy, "A Survey on Big Data Management and Job Scheduling" *International Journal of Computer Applications* (0975 – 8887) Volume 130 – No.13, November 2015.
- [2] Dawei Jiang, Sai Wu, Gang Chen, Beng Chin Ooi and Kian-Lee Tan, Jun Xu, "epiC: an extensible and scalable system for processing Big Data" *The VLDB Journal* (2016) 25:3–26 DOI 10.1007/s00778-015-0393-2.
- [3] Marcos D. Assunção, Rodrigo N. Calheiros, Silvia Bianchi, Marco A.S. Netto, Rajkumar Buyya, "Big Data computing and clouds: Trends and future directions" *J. Parallel Distrib. Comput.* 79–80 (2015) 3–15.
- [4] Xiongpai Qin, and Xiaoyun Zhou, "A Survey on Benchmarks for Big Data and Some More Considerations" H. Yin et al. (Eds.): *IDEAL 2013*, LNCS 8206, pp. 619–627, 2013. Springer-Verlag Berlin Heidelberg 2013.
- [5] Leonardo Neumeyer, Bruce Robbins, Anish Nair, Anand Kesari, "S4: Distributed Stream Computing Platform" 2010 IEEE International Conference on Data Mining Workshops.
- [6] Rakesh Varma, "Survey on MapReduce and Scheduling Algorithms in Hadoop" *International Journal of Science and Research (IJSR)* ISSN (Online): 2319-7064 Index Copernicus Value (2013): 6.14 | Impact Factor (2013): 4.438.
- [7] Suman Storage and Dr. Madhu Goel, "Survey Paper on Scheduling in Hadoop" Volume 4, Issue 5, May 2014 ISSN: 2277 128X *International Journal of Advanced Research in Computer Science and Software Engineering*.
- [8] Sushant Goel, Hema Sharda and David Tanid, "Distributed Scheduler for High Performance Data-Centric Systems" b7803-76CI-XIO11B17.00 0 2003 IEEE.
- [9] John. H. Phan, Chang. F. Quo, and May D. Wang, "Comparative Study of Microarray Data for Cancer Research" proceedings of the 26th Annual International Conference of IEEE EMBS San Francisco, CA, USA \* September 1-5, 2004.
- [10] Seema Metikurke and Vijay K. Vaishnavi, "Grid-Enabled Automatic Web Page Classification" 2006 IEEE International Conference on Fuzzy Systems Sheraton Vancouver Wall Centre Hotel, Vancouver, BC, Canada July 16-21, 2006.
- [11] Hui Fang, Ming Yang and Ruqing Yang, "Ground Texture Matching based Global Localization for Intelligent Vehicles in Urban Environment" Proceedings of the 2007 IEEE Intelligent Vehicles Symposium Istanbul, Turkey, June 13-15, 2007.
- [12] BI Shuoben, XU Yin, JIAO Feng, Lü Guonian, PEI Anping, "Study on Data Mining in First Period of Jiangzhai Site Based on the Association Algorithms" 2009 International Conference on Artificial Intelligence and Computational Intelligence, 2009 IEEE DOI 10.1109/AICI.2009.
- [13] Jasmin Azemovic, Denis Music, "Comparative analysis of efficient methods for storing unstructured data into database with accent on performance" 2010, IEEE 2nd International Conference on Education Technology and Computer (ICETC).
- [14] Daniel Warneke, "Exploiting Dynamic Resource Allocation for Efficient Parallel Data Processing in the Cloud" *IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS*, VOL. 22, NO. 6, JUNE 2011.
- [15] Chang Liu, Jinjun Chen, Chi Yang, Rajiv Ranjan and Ramamohanarao Kotagiri, "Authorized Public Auditing of Dynamic Big Data Storage on Cloud with Efficient Verifiable Fine-Grained Updates" *IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS*, VOL. 25, NO. 9, SEPTEMBER 2014.
- [16] Shifeng Fang, Li DaXu, Yunqiang Zhu, Jiaerheng Ahati, Huan Pei, Jianwu Yan, and Zhihui Liu, "An Integrated System for Regional Environmental Monitoring and Management Based on Internet of Things" *IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS*, VOL. 10, NO. 2, MAY 2014.