



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

## A Study on New Challenges of Big Data and Apache Spark

Ravula Kartheek<sup>1</sup>, Koteswara Rao Pilly<sup>2</sup>, K. Pardhasaradhi<sup>3</sup>, A.V.Ramana<sup>4</sup>

Assistant Professor, Dept. of Computer Science and Engineering, Rise Krishna Sai Gandhi Group of Institutions  
Ongole, India

Assistant Professor, Dept. of MCA, Rise Krishna Sai Gandhi Group of Institutions, Ongole, India

Assistant Professor, Dept. of Computer Science and Engineering, Rise Krishna Sai Prakasam Group of Institutions  
Ongole, India

Assistant Professor, Dept. of Computer Science and Engineering, Rise Krishna Sai Prakasam Group of Institutions  
Ongole, India

**ABSTRACT:** In this paper we clearly shows the challenges between the working procedure of Hadoop and Apache Spark. It shows the detailed analysis of big data and Apache Spark. And also it will explains the benefits of big data and big data technologies, operational big data and analytical big data, challenges and some more topics will be explained clearly. This paper mainly shows a brief idea on Hadoop and Apache Spark. And it will explain about the data storing procedure in olden days and at present.

**KEYWORDS:** Cloud, Apache, Analytics, Hadoop, Big Data, Spark

### I. INTRODUCTION

Due to the appearance of recent technologies, gadgets, and communication means like social networking sites, the amount of records produced via mankind is developing swiftly every 12 months. The amount of data produced by us from the beginning of time till 2003 becomes five billion gigabytes. If you pile up the data in the shape of disks it can fill a whole football subject. The equal quantity was created in every two days in 2011 and in every 10 minutes in 2013. This price continues to be growing highly. Although all this data produced is significant and can be beneficial whilst processed, its miles being ignored. Ninety percent of the world's records become generated inside the last few years. Big Data means sincerely a huge facts, it is a collection of huge information sets that cannot be processed using traditional computing techniques. Large statistics is not merely a facts, as a substitute it has become an entire situation, which entails various equipment, techniques and frameworks. Apache Spark runs on Hadoop, Mesos, standalone, or inside the cloud. It may access various information resources along with HDFS, Cassandra, HBase, and S3. You could run Spark the use of its standalone cluster mode, on EC2, on Hadoop YARN, or on Apache Mesos. Get admission to facts in HDFS, Cassandra, HBase, Hive, Tachyon, and any Hadoop facts supply.

This paper was clearly explains what comes under big data, benefits, big data technologies, big data challenges. And after that here we will explained about the traditional approaches and limitations of big data. In last we will explained about the Apache Spark.

### WHAT IS BIG DATA?

Big data means absolutely a large records, its miles a set of huge datasets that cannot be processed the use of conventional computing techniques. Big data aren't always simply a statistics, as a substitute it has come to be an entire subject, which entails numerous equipment, techniques and frameworks.

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

## WHAT COMES UNDER THE BIG DATA?

Large information involves the facts produced by means of special gadgets and applications. Given underneath are a number of the fields that come under the umbrella of massive data.

**\*Black box data\***: It is part of helicopter, airplanes, and jets, and so forth. It captures voices of the flight team, recordings of microphones and earphones, and the overall performance records of the plane.

**\*Social Media information\***: Social media such as fb and Twitter maintain records and the perspectives published by using tens of millions of people throughout the globe.

**\*Stock exchange data\*** : The stock change facts holds data approximately the 'buy' and 'sell' decisions made on a proportion of various companies made by using the customers.

**\*Power Grid information \***: The electricity grid records holds statistics fed on via a selected node with recognize to a base station.

**\*Transport data\***: delivery statistics consists of version, capability, distance and availability of a automobile.

**\*search engine data\***: engines like Google retrieve masses of records from special databases.



Fig1: what comes under the Big Data?

Therefore massive information consists of huge quantity, high pace, and extensible type of facts. The records in it'll be of 3 types.

**\*Structured data\***: Relational information.

**\*Semi structured data\***: XML statistics.

**\*Unstructured data\***: word, PDF, text, Media Logs.

## II. BENEFITS OF BIG DATA

Big data are definitely crucial to our existence and it's emerging as one of the maximum crucial technologies in contemporary international. Observe are just few advantages that are very an awful lot recognized to absolutely everyone:

The usage of the facts saved in the social network like face book, the advertising agencies are studying about the response for his or her campaigns, promotions, and other marketing mediums.

The usage of the statistics within the social media like possibilities and product belief of their consumers, product agencies and retail corporations are making plans their manufacturing.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

The usage of the information regarding the preceding clinical records of sufferers, hospitals are presenting higher and short provider.

### III. BIG DATA TECHNOLOGY

Big data technologies are essential in offering extra correct analysis, which can also lead to extra concrete choice-making resulting in greater operational efficiencies, value reductions, and decreased risks for the enterprise.

To harness the power of huge statistics, you'll require an infrastructure that can control and process huge volumes of established and unstructured facts in real-time and can protect statistics privacy and protection.

There are various technologies within the market from extraordinary carriers which includes Amazon, IBM, Microsoft, etc., to address big data. While searching into the technologies that take care of huge statistics, we examine the subsequent classes of technology:

#### a. OPERATIONAL BIG DATA

This include systems like MongoDB that provide operational abilities for actual-time, interactive workloads wherein statistics is generally captured and stored.

NoSQL large records systems are designed to take gain of new cloud computing architectures that have emerged over the last decade to permit huge computations to be run inexpensively and efficaciously. This makes operational huge information workloads an awful lot less difficult to manipulate, less expensive, and faster to enforce.

Some NoSQL systems can offer insights into styles and traits based on real-time facts with minimum coding and without the want for records scientists and additional infrastructure.

#### b. ANALYTICAL BIG DATA

This consists of systems like massively Parallel Processing (MPP) database systems and Map Reduce that offer analytical capabilities for retrospective and complicated evaluation that may contact maximum or all the facts.

Map Reduce offers a brand new approach of studying statistics this is complementary to the capabilities supplied by means of square, and a machine primarily based on Map Reduce that may be scaled up from single servers to thousands of excessive and coffee give up machines.

Those two instructions of era are complementary and frequently deployed collectively.

#### c. OPERATIONAL VERSUS ANALYTICAL STRUCTURES

	Operational	Analytical
<b>Latency</b>	1 ms - 100 ms	1 min - 100 min
<b>Concurrency</b>	1000- 100,000	1 - 10
<b>Access Pattern</b>	Writes and Reads	Reads
<b>Queries</b>	Selective	unselective
<b>End User</b>	Customer	Data Scientist
<b>Data Scope</b>	Operational	Retrospective
<b>Technology</b>	NoSQL	Map Reduce, MPP Database

Table 1: Operational Versus Analytical Structures

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

## IV. BIG DATA CHALLENGES

The predominant demanding situations related to big data are as follows:

- Capturing information
- Curation
- Storage
- Searching
- Sharing
- Transfer
- Evaluation or Analysis
- Presentation

To meet the above demanding situations, businesses normally take the help of organization servers.

## V. CONVENTIONAL METHOD

In this approach, an organization will have a computer device to store and technique large data facts. Right here information can be stored in an RDBMS like Oracle Database, MS SQL Server or DB2 and sophisticated software's may be written to interact with the database, process the required information and present it to the customers for evaluation cause.

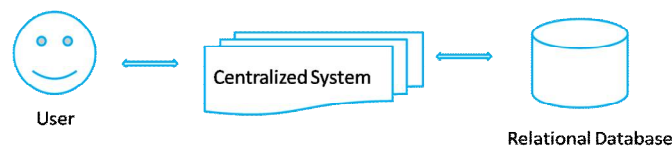


Fig2: Conventional Method

## VI. LIMITATIONS

This technique works nicely where we've got much less extent of facts that can be accommodated with the aid of preferred database servers, or up to the limit of the processor that's processing the statistics. But on the subject of coping with massive quantities of records, it is truly a tedious venture to procedure such records thru a conventional database server.

## VII. GOOGLE'S SOLUTION

Google solved this hassle the usage of a set of rules known as Map Reduce. This set of rules divides the challenge into small components and assigns the ones components to many computers related over the network, and collects the consequences to form the final result dataset.

## VIII. HADOOP

Doug reducing, Mike Cafarella and group took the solution furnished by way of Google and commenced an Open supply task known as HADOOP in 2005 and Doug named it after his son's toy elephant. Now Apache Hadoop is a registered trademark of the Apache software program basis.

Hadoop runs packages the usage of the Map Reduce algorithm, wherein the records is processed in parallel on one of a kind CPU nodes. In brief, Hadoop framework is capable enough to develop applications able to running on clusters of computers and they may carry out complete statistical analysis for large amounts of statistics.

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

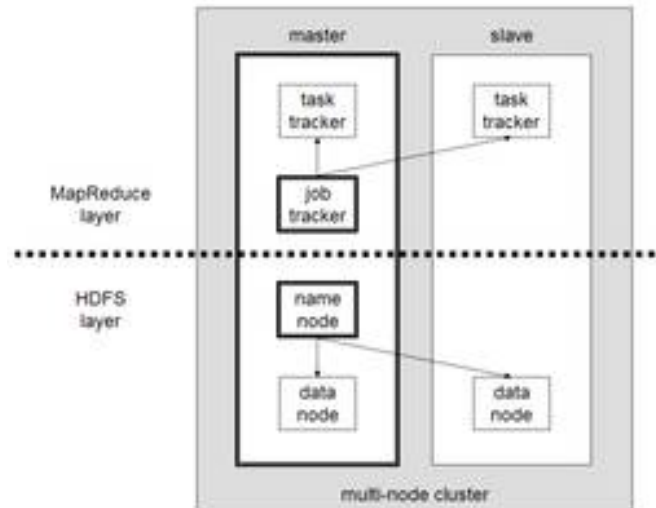


Fig3: Hadoop Framework

## IX. SECURITY ISSUES IN BIG DATA

The assurance of data is another huge concern, and one that augmentations in the association of Big Data. Not at all like standard security strategy, is security in enormous data on a very basic level in the sort of how to process data mining without revealing fragile information of customers. Moreover, current advancements of security protection are essentially in perspective of static data set, while data is reliably vigilantly changed, including data outline, assortment of characteristic and development of new data. Along these lines, it is a test to execute effective security protection in this brain boggling circumstance. Additionally, authentic and authoritative issues moreover require thought. For electronic prosperity records, there are strict laws regulating what ought to and can't be conceivable. For other data, controls, particularly in the US, are less solid. Managing security is effectively both a particular and a sociological issue, which must be tended to commonly from both perspectives to comprehend the assurance of gigantic data. Learning driven security relies on upon immense data examination. By keeping data in one spot, it happens to be a goal for attackers to hurt the affiliation. It obliged that immense data stores are properly controlled. To ensure affirmation a cryptographically secure correspondence framework must be executed.

Controls should use standard of decreased advantages, especially for get to rights, except for a head who have approval data to physical get to. For suitable get to controls, they should be unendingly watched and traded as change specialists affiliation parts so agents don't add up to radical rights that could be mishandled. Other security techniques are relied upon to get and analyze framework development, for instance, metadata, package catch, stream and log information. Affiliations should guarantee interests in security things using agile progressions based examination not static supplies. Another issue is associated with orchestrating consistence of data security laws. Affiliations need to consider legal fanning for securing data. In any case, gigantic data has security purposes of intrigue. Right when affiliations arrange data, they control data as showed by dictated by the directions, for instance, constraining store periods.

This licenses relationship to pick data that has neither little regard nor any ought to be kept so it is not any more open for burglary. Another preferred standpoint is gigantic data can be burrowed for perils, for instance, affirmation of malware, irregularities, or phishing. The for the most part less sorted out and easygoing nature of various Big Data strategies is their quality, be that as it may it furthermore speaks to an issue: if the data included is tricky for reasons of security, attempt security, or authoritative need, then using such philosophies may address a honest to goodness security crack. Database organization structures reinforce security techniques that are genuinely granular, guaranteeing data at both coarse and fine grain level from wrong get to.

Enormous Data programming generally has no such secures. Wanders that fuse any sensitive data in Big Data operations must ensure that the data itself is secure, and that similar data security approaches that apply to the data when it exists in databases or records are in like manner approved in the Big Data association. Powerlessness to do in that capacity can have real negative results.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

## X. APACHE SPARK

Getting a decent hang on the tremendous information, it should be quite quick. For instance, in the event that we discuss the information being created by Wal-Mart stores everywhere throughout the world, that would be a large number of offer sections each hour, correct? In this way, it would be truly awful to have information researchers to give bits of knowledge on deals amid a specific time of day just if the calculation takes a day. Additionally, the information researchers ought to have the capacity to process it in aggregate without a moment's delay. Henceforth, it is required of Spark to be accessible on groups, as opposed to demanding the prerequisite of a solitary machine. Along these lines, Spark gloats a world record in expansive scale sorting by Data bricks. That was made conceivable due to the two components talked about above. Spark stores information sets in memory, which makes it a 100x speedier than Hadoop, which does as such on plate.

Additionally, permitting client projects to load information to a bunch's memory and permitting rehashed questioning, it is a structure appropriate to machine learning calculations. Parts of Spark:

1. **Versatile Conveyed Datasets and the Spark Center:** The Spark Center is the establishment and gives fundamental I/O functionalities, errand dispatching and booking. RDDs are fundamentally an accumulation of divided information. These are for the most part made by referencing datasets in stockpiles, for example, Cassandra, HBase et al., or by applying changes, for example, delineate, and channel and so on existing RDDs.

2. **Spark SQL:** Spark SQL, a segment on the Center, presents another information deliberation called Data Frame, for giving backing to organized information. It gives a dialect to control Data Frames in Java, Python or Scala.

3. **Spark Spilling:** Spark gushing lays on the Center also, and levera on top of the Center which is ended up being ten times quicker than Hadoop's circle based Apache Mahout because of the appropriated memory-based Spark engineering. It executes normal calculations to streamline vast scale machine learning pipelines, as strategic or straight relapse, choice trees or k-implies grouping.

4. **MLlib Machine Learning Library:** This is a machine learning system on top of the Center which is ended up being ten times speedier than Hadoop's plate based Apache Mahout because of the circulated memory-based Spark engineering. It actualizes normal calculations to improve substantial scale machine learning pipelines, as strategic or straight relapse, choice trees or k-implies bunching.

5. **GraphX:** It is a diagram handling system on the Center, and gives a Programming interface to chart calculation that can demonstrate the Pregel deliberation, giving an upgraded runtime.

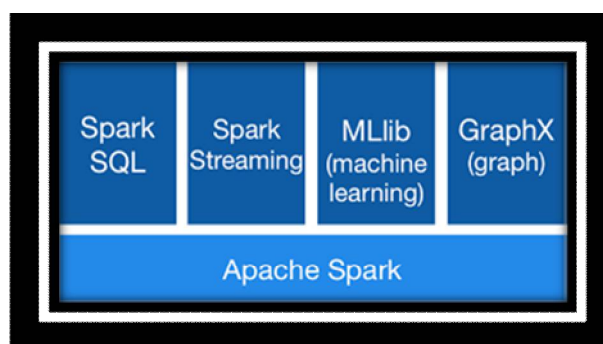


Fig. 4 Apache Spark



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

## XI. WORLD RECORD SET BY APACHE SPARK

This is an examination between Hadoop Outline and Apache Spark for sorting information and setting world record:

	Hadoop MR Record	Spark Record	Spark 1 PB
Data Size	102.5 TB	100 TB	1000 TB
Elapsed Time	72 mins	23 mins	234 mins
Cluster disk throughput	3150 GB/s	618 GB/s	570 GB/s
Sort Rate	1.42TB/min	4.27TB/min	4.27TB/min

Table 2 correlation between hadoop map reduce and apache spark

## XII. CONCLUSION

The paper finishes up with the recommendation that Enormous Information is a blasting field at the present minute, and the tremendous measure of information that gets created each minute requires an exceptionally powerful administration and investigation framework that can manage the greatness. Moreover, this paper looks to legitimize the qualities of Apache Spark and its remaining as exceptionally productive programming relating to the present situation of Enormous Information. In October 2014, Data bricks appreciated the Sort Benchmark and set a different universe record for sorting 100 terabytes (TB) of data, or 1 trillion 100-byte records. The gathering used Apache Spark on 207 EC2 virtual machines and sorted 100 TB of data in 23 minutes. In correlation, the past world record set by Hadoop Map Reduce used 2100 machines as a part of a private server farm and took 72 minutes. This area tied with a UCSD inquires about gathering developing superior systems.

## REFERENCES

1. Sanjay Rathee, "Big Data and Hadoop with components like Flume, Pig, Hive and Jaql", International Conference on Cloud, Big Data and Trust 2013, Vol. 15, November 2013
2. E. Begoli and J. Horey, "Design Principles for Effective Knowledge Discovery from Big Data", Software Architecture (WICSA) and European Conference on Software Architecture (ECSA) Joint Working IEEE/IFIP Conference on, Helsinki, August 2012
3. BhartiThakur, ManishMann, "Data Mining for Big Data: A Review", International Journal of Advanced Research in Computer Science and Software Engineering, Vol. 4, May 2015
4. Wei Fan, Albert Weifet, "Mining Big Data: Current Status, and Forecast to the Future", SIGKDD Explorations, Vol. 14, February 2012
5. Hsinchun Chen, Roger H.L. Chiang, Veda C. Storey, "Business Intelligence and Analytics: From Big Data to Big Impact", MIS Quarterly Special Issue: Business Intelligence Research, Vol. 36, December 2012
6. A. Vailaya, "What's All the Buzz Around 'Big Data'?", IEEE Women in Engineering Magazine, December 2012, pp. 24-31
7. Harshwardhan S. Bhosale, Prof. Devendra P. Gaddekar, "A review, paper on big data and hadoop", International Journal of Scientific Research and Publications, Vol. 4, October 2014
8. [https://www.tutorialspoint.com/hadoop/hadoop\\_big\\_data\\_solutions.htm](https://www.tutorialspoint.com/hadoop/hadoop_big_data_solutions.htm)
9. S. Singh and N. Singh, "Big Data Analytics", 2012 International Conference on Communication, Information & Computing Technology Mumbai India, IEEE, October 2011.
10. Jeffhurlblog.com "three-vs.-of-big-data-as-applied-conferences", July 7, 2012.

## BIOGRAPHY



**R. Kartheek** has received his B.Tech in Information Technology and M.Tech degree in Computer Science and Engineering from JNTU Kakinada in 2012 and JNTU Kakinada in 2015 respectively. He is dedicated to teaching field from the past 1 year. His research areas included Computer Networks, Data mining, Wireless Networks, oops etc. At present he is working as Assistant Professor in RISE Krishna Sai Gandhi Group of Institutions: Ongole, Andhra Pradesh, India.



ISSN(Online): 2320-9801  
ISSN (Print): 2320-9798

## International Journal of Innovative Research in Computer and Communication Engineering

*(An ISO 3297: 2007 Certified Organization)*

**Vol. 4, Issue 11, November 2016**



**Koteswara Rao Pilly**, Presently Working as an “Assistant Professor in MCA Department” in Rise Krishna Sai Gandhi Group of Institutions, Ongole, Prakasam District, A.P, India. Affiliated to Jawaharlal Nehru Technological University, Kakinada. Approved by AICTE, New Delhi. His MCA completed at Rajeev Gandhi Degree & PG College, East Godavari District, A.P, and India. His M.Tech completed in Rao & Naidu Engineering College, Ongole. His research interests are network security, Computer Networks etc.



**K. Pardhasaradhi**, Presently Working as an “Assistant Professor in CSE Department” in Rise Krishna Sai Prakasam Group of Institutions, Ongole, Prakasam District, A.P, India. Affiliated to Jawaharlal Nehru Technological University, Kakinada. Approved by AICTE, New Delhi. His B.Tech completed at Krishna Chaitanya Institute of Technology & Sciences, Prakasam District, A.P, and India. His M.Tech completed in Vasireddy Venkatadri Institute of Technology, Nambur, Guntur. His research interests are network security, Computer Networks etc.



**A.V..RAMANA** Presently Working as an “Assistant Professor in CSE Department” in Rise Krishna Sai Prakasam Group of Institutions, Ongole, Prakasam District, A.P, India. Affiliated to Jawaharlal Nehru Technological University, Kakinada. Approved by AICTE, New Delhi. His M.Sc(CS) at S.V.P.G College, Kadapa District, A.P, and India. His M.Tech completed in ANU,Guntur. His research interests are network security, Computer Networks etc.