



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Hadoop – Big Data

An Important technology in IT Sector

Written By,
IIT Bombay Alumni Foundation's
Embedded Technosolutions





Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Hadoop - Big Data Overview

“90% of the world’s data was generated in the last few years.”

Due to the advent of new technologies, devices, and communication means like social networking sites, the amount of data produced by mankind is growing rapidly every year. The amount of data produced by us from the beginning of time till 2003 was 5 billion gigabytes. If you pile up the data in the form of disks it may fill an entire football field. The same amount was created in every two days in 2011, and in every ten minutes in 2013. This rate is still growing enormously. Though all this information produced is meaningful and can be useful when processed, it is being neglected.

What is Big Data?

Big Data is a collection of large datasets that cannot be processed using traditional computing techniques. It is not a single technique or a tool, rather it involves many areas of business and technology.

What Comes Under Big Data?

Big data involves the data produced by different devices and applications. Given below are some of the fields that come under the umbrella of Big Data.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



- **Black Box Data** : It is a component of helicopter, airplanes, and jets, etc. It captures voices of the flight crew, recordings of microphones and earphones, and the performance information of the aircraft.
- **Social Media Data** : Social media such as Facebook and Twitter hold information and the views posted by millions of people across the globe.
- **Stock Exchange Data** : The stock exchange data holds information about the 'buy' and 'sell' decisions made on a share of different companies made by the customers.
- **Power Grid Data** : The power grid data holds information consumed by a particular node with respect to a base station.
- **Transport Data** : Transport data includes model, capacity, distance and availability of a vehicle.
- **Search Engine Data** : Search engines retrieve lots of data from different databases.





Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Thus Big Data includes huge volume, high velocity, and extensible variety of data. The data in it will be of three types.

- **Structured data** : Relational data.
- **Semi Structured data** : XML data.
- **Unstructured data** : Word, PDF, Text, Media Logs.

Benefits of Big Data

- Using the information kept in the social network like Facebook, the marketing agencies are learning about the response for their campaigns, promotions, and other advertising mediums.
- Using the information in the social media like preferences and product perception of their consumers, product companies and retail organizations are planning their production.
- Using the data regarding the previous medical history of patients, hospitals are providing better and quick service.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Big Data Technologies

Big data technologies are important in providing more accurate analysis, which may lead to more concrete decision-making resulting in greater operational efficiencies, cost reductions, and reduced risks for the business.

To harness the power of big data, you would require an infrastructure that can manage and process huge volumes of structured and unstructured data in realtime and can protect data privacy and security.

There are various technologies in the market from different vendors including Amazon, IBM, Microsoft, etc., to handle big data. While looking into the technologies that handle big data, we examine the following two classes of technology:

Operational Big Data

These include systems like MongoDB that provide operational capabilities for real-time, interactive workloads where data is primarily captured and stored.

NoSQL Big Data systems are designed to take advantage of new cloud computing architectures that have emerged over the past decade to allow massive computations to be run inexpensively and efficiently. This makes operational big data workloads much easier to manage, cheaper, and faster to implement.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Some NoSQL systems can provide insights into patterns and trends based on real-time data with minimal coding and without the need for data scientists and additional infrastructure.

Analytical Big Data

These includes systems like Massively Parallel Processing (MPP) database systems and MapReduce that provide analytical capabilities for retrospective and complex analysis that may touch most or all of the data.

MapReduce provides a new method of analyzing data that is complementary to the capabilities provided by SQL, and a system based on MapReduce that can be scaled up from single servers to thousands of high and low end machines.

These two classes of technology are complementary and frequently deployed together.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Operational vs. Analytical Systems

	Operational	Analytical
Latency	1 ms - 100 ms	1 min - 100 min
Concurrency	1000 - 100,000	1 - 10
Access Pattern	Writes and Reads	Reads
Queries	Selective	Unselective
Data Scope	Operational	Retrospective
End User	Customer	Data Scientist
Technology	NoSQL	MapReduce, MPP Database



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Big Data Challenges

The major challenges associated with big data are as follows:

- Capturing data
- Curation
- Storage
- Searching
- Sharing
- Transfer
- Analysis
- Presentation

To fulfill the above challenges, organizations normally take the help of enterprise servers.

Hadoop - Big Data Solutions

Traditional Enterprise Approach

In this approach, an enterprise will have a computer to store and process big data. For storage purpose, the programmers will take the help of their choice of database vendors such as Oracle, IBM, etc. In this approach, the user interacts with the application, which in turn handles the part of data storage and analysis.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

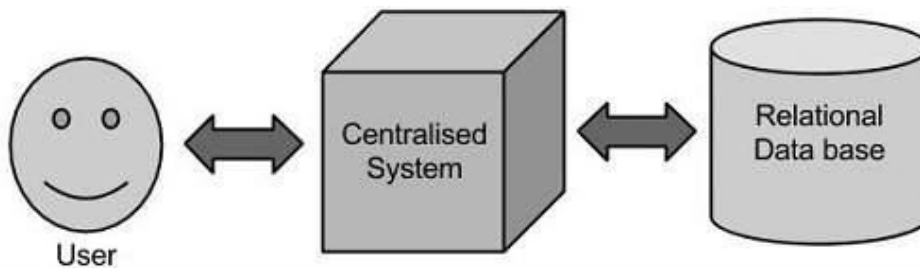
Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Limitation

This approach works fine with those applications that process less voluminous data that can be accommodated by standard database servers, or up to the limit of the processor that is processing the data. But when it comes to dealing with huge amounts of scalable data, it is a hectic task to process such data through a single database bottleneck.

Google's Solution

Google solved this problem using an algorithm called MapReduce. This algorithm divides the task into small parts and assigns them to many computers, and collects the results from them which when integrated, form the result dataset.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

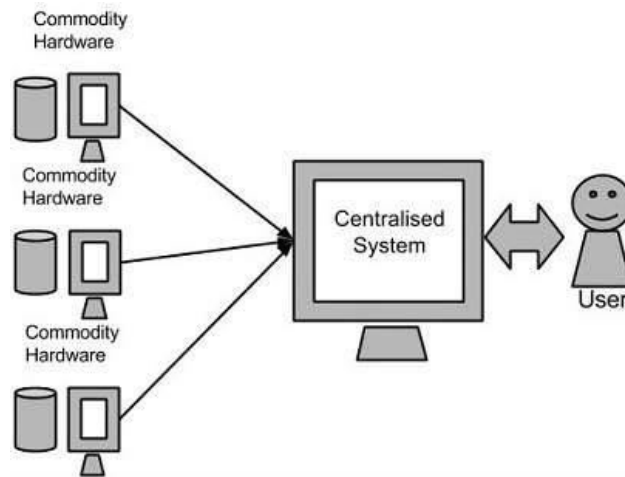
Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Hadoop

Using the solution provided by Google, Doug Cutting and his team developed an Open Source Project called HADOOP.

Hadoop runs applications using the MapReduce algorithm, where the data is processed in parallel with others. In short, Hadoop is used to develop applications that could perform complete statistical analysis on huge amounts of data.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

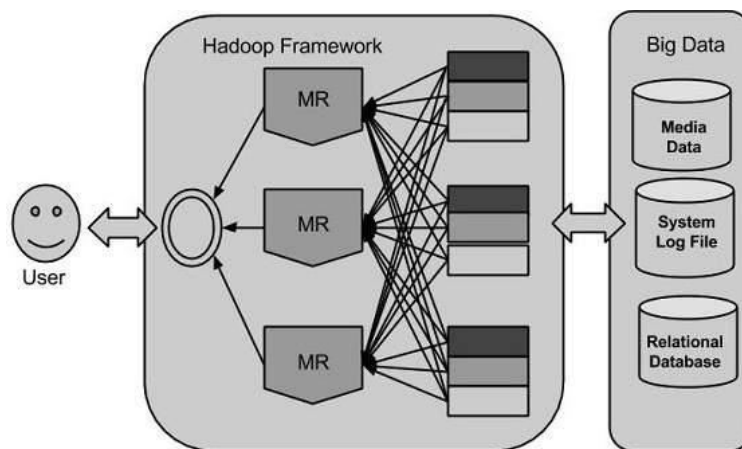
Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



Hadoop - Introduction to Hadoop

Hadoop is an Apache open source framework written in java that allows distributed processing of large datasets across clusters of computers using simple programming models. The Hadoop framework application works in an environment that provides distributed storage and computation across clusters of computers. Hadoop is designed to scale up from single server to thousands of machines, each offering local computation and storage.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

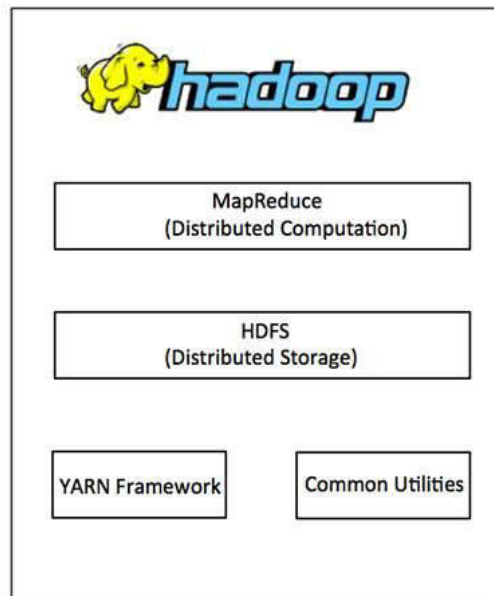
www.embeddedtechnosolutions.com



Hadoop Architecture

At its core, Hadoop has two major layers namely:

- Processing/Computation layer (MapReduce), and
- Storage layer (Hadoop Distributed File System).



MapReduce

MapReduce is a parallel programming model for writing distributed applications devised at Google for efficient processing of large amounts of data (multiterabyte data-sets), on large



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



clusters (thousands of nodes) of commodity hardware in a reliable, fault-tolerant manner. The MapReduce program runs on Hadoop which is an Apache open-source framework.

Hadoop Distributed File System

The Hadoop Distributed File System (HDFS) is based on the Google File System (GFS) and provides a distributed file system that is designed to run on commodity hardware. It has many similarities with existing distributed file systems. However, the differences from other distributed file systems are significant. It is highly fault-tolerant and is designed to be deployed on low-cost hardware. It provides high throughput access to application data and is suitable for applications having large datasets.

Apart from the above-mentioned two core components, Hadoop framework also includes the following two modules:

- **Hadoop Common** : These are Java libraries and utilities required by other Hadoop modules.
- **Hadoop YARN** : This is a framework for job scheduling and cluster resource management.

How Does Hadoop Work?

It is quite expensive to build bigger servers with heavy configurations that handle large scale processing, but as an alternative, you can tie together many commodity computers with single-CPU, as a single functional distributed system and practically, the clustered machines can read the dataset in parallel and provide a much higher throughput. Moreover, it is cheaper



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



than one high-end server. So this is the first motivational factor behind using Hadoop that it runs across clustered and low-cost machines.

Hadoop runs code across a cluster of computers. This process includes the following core tasks that Hadoop performs:

- Data is initially divided into directories and files. Files are divided into uniform sized blocks of 128M and 64M (preferably 128M).
- These files are then distributed across various cluster nodes for further processing.
- HDFS, being on top of the local file system, supervises the processing.
- Blocks are replicated for handling hardware failure.
- Checking that the code was executed successfully.
- Performing the sort that takes place between the map and reduce stages.
- Sending the sorted data to a certain computer.
- Writing the debugging logs for each job.

Hadoop Operation Modes

Once you have downloaded Hadoop, you can operate your Hadoop cluster in one of the three supported modes:

- **Local/Standalone Mode** : After downloading Hadoop in your system, by default, it is configured in a standalone mode and can be run as a single java process.



Embedded Technosolutions

Venture of IIT Bombay & VJTI Alumni

Embedded Systems | Software | Mechanical | Automation

Trainings & Jobs

100% Placement Assistance

Contact : 8828222688 / 9224301650

www.embeddedtechnosolutions.com



-
- **Pseudo Distributed Mode** : It is a distributed simulation on single machine. Each Hadoop daemon such as hdfs, yarn, MapReduce etc., will run as a separate java process. This mode is useful for development.
 - **Fully Distributed Mode** : This mode is fully distributed with minimum two or more machines as a cluster. We will come across this mode in detail in the coming chapters.

Advantages of Hadoop

- Hadoop framework allows the user to quickly write and test distributed systems. It is efficient, and it automatic distributes the data and work across the machines and in turn, utilizes the underlying parallelism of the CPU cores.
- Hadoop does not rely on hardware to provide fault-tolerance and high availability (FTHA), rather Hadoop library itself has been designed to detect and handle failures at the application layer.
- Servers can be added or removed from the cluster dynamically and Hadoop continues to operate without interruption.
- Another big advantage of Hadoop is that apart from being open source, it is compatible on all the platforms since it is Java based.