# Learn Hadoop in 1 Day

By Krishna Rungta

# Table Of Content

# Chapter 1: Introduction to BIG DATA: What is, Types, Characteristics & Example

In order to understand 'Big Data', you first need to know

## What is Data?

The quantities, characters, or symbols on which operations are performed by a computer, which may be stored and transmitted in the form of electrical signals and recorded on magnetic, optical, or mechanical recording media.

## What is Big Data?

Big Data is also **data** but with a **huge size**. Big Data is a term used to describe a collection of data that is huge in size and yet growing exponentially with time. In short such data is so large and complex that none of the traditional data management tools are able to store it or process it efficiently.

# Examples Of Big Data

Following are some the examples of Big Data-

The **New York Stock Exchange** generates about *one terabyte* of new trade data per day.

**Social Media**

The statistic shows that ***500+terabytes*** of new data get ingested into the databases of social media site **Facebook**, every day. This data is mainly generated in terms of photo and video uploads, message exchanges, putting comments etc.



A single **Jet engine** can generate ***10+terabytes*** of data in ***30 minutes*** of flight time. With many thousand flights per day, generation of data reaches up to many ***Petabytes.***

# Types Of Big Data

BigData' could be found in three forms:

1. **Structured**
2. **Unstructured**
3. **Semi-structured**

## Structured

Any data that can be stored, accessed and processed in the form of fixed format is termed as a 'structured' data. Over the period of time, talent in computer science has achieved greater success in developing techniques for working with such kind of data (where the format is well known in advance) and also deriving value out of it. However, nowadays, we are foreseeing issues when a size of such data grows to a huge extent, typical sizes are being in the rage of multiple zettabytes.

***Do you know? $10^{21}$ bytes*** equal to ***1 zettabyte*** or ***one billion terabytes*** forms ***a zettabyte***.

Looking at these figures one can easily understand why the name BigData is given and imagine the challenges involved in its storage and processing.

***Do you know?*** Data stored in a relational database management system is one example of a **'structured'** data.

**Examples Of Structured Data**

An 'Employee' table in a database is an example of Structured Data

| Employee_ID | Employee_Name | Gender | Department | Salary_In_lacs |
|---|---|---|---|---|
| 2365 | Rajesh Kulkarni | Male | Finance | 650000 |
| 3398 | Pratibha Joshi | Female | Admin | 650000 |
| 7465 | Shushil Roy | Male | Admin | 500000 |
| 7500 | Shubhojit Das | Male | Finance | 500000 |
| 7699 | Priya Sane | Female | Finance | 550000 |

# Unstructured

Any data with unknown form or the structure is classified as unstructured data. In addition to the size being huge, un-structured data poses multiple challenges in terms of its processing for deriving value out of it. A typical example of unstructured data is a heterogeneous data source containing a combination of simple text files, images, videos etc. Now day organizations have wealth of data available with them but unfortunately, they don't know how to derive value out of it since this data is in its raw form or unstructured format.

**Examples Of Un-structured Data**

The output returned by 'Google Search'

# Semi-structured

Semi-structured data can contain both the forms of data. We can see semi-structured data as a structured in form but it is actually not defined with e.g. a table definition in relational DBMS. Example of semi-structured data is a data represented in an XML file.

Examples Of Semi-structured Data

Personal data stored in an XML file-

```
<rec><name>Prashant Rao</name><sex>Male</sex><age>35</age></rec>
<rec><name>Seema R.</name><sex>Female</sex><age>41</age></rec>
<rec><name>Satish Mane</name><sex>Male</sex><age>29</age></rec>
<rec><name>Subrato Roy</name><sex>Male</sex><age>26</age></rec>
<rec><name>Jeremiah J.</name><sex>Male</sex><age>35</age></rec>
```

**Data Growth over the years**

Please note that web application data, which is unstructured, consists of log files, transaction history files etc. OLTP systems are built to work with structured data wherein data is stored in relations (tables).

# Characteristics Of Big Data

*(i) Volume –* The name Big Data itself is related to a size which is enormous. Size of data plays a very crucial role in determining value out of data. Also, whether a particular data can actually be considered as a Big Data or not, is dependent upon the volume of data. Hence, **'Volume'** is one characteristic which needs to be considered while dealing with Big Data.

*(ii) Variety –* The next aspect of Big Data is its **variety**.

Variety refers to heterogeneous sources and the nature of data, both structured and unstructured. During earlier days, spreadsheets and databases were the only sources of data considered by most of the applications. Nowadays, data in the form of emails, photos, videos, monitoring devices, PDFs, audio, etc. are also being considered in the

analysis applications. This variety of unstructured data poses certain issues for storage, mining and analyzing data.

***(iii) Velocity* –** The term **'velocity'** refers to the speed of generation of data. How fast the data is generated and processed to meet the demands, determines real potential in the data.

Big Data Velocity deals with the speed at which data flows in from sources like business processes, application logs, networks, and social media sites, sensors, Mobile devices, etc. The flow of data is massive and continuous.

***(iv) Variability* –** This refers to the inconsistency which can be shown by the data at times, thus hampering the process of being able to handle and manage the data effectively.

## Benefits of Big Data Processing

Ability to process Big Data brings in multiple benefits, such as-

- Businesses can utilize outside intelligence while taking decisions

Access to social data from search engines and sites like facebook, twitter are enabling organizations to fine tune their business strategies.

- Improved customer service

Traditional customer feedback systems are getting replaced by new systems designed with Big Data technologies. In these new systems, Big Data and natural language processing technologies are being used

to read and evaluate consumer responses.

- Early identification of risk to the product/services, if any Better
- operational efficiency

Big Data technologies can be used for creating a staging area or landing zone for new data before identifying what data should be moved to the data warehouse. In addition, such integration of Big Data technologies and data warehouse helps an organization to offload infrequently accessed data.

## Summary

- Big Data is defined as data that is huge in size. Bigdata is a term used to describe a collection of data that is huge in size and yet growing exponentially with time.
- Examples of Big Data generation includes stock exchanges, social media sites, jet engines, etc.
- Big Data could be 1) Structured, 2) Unstructured, 3) Semi- structured
- Volume, Variety, Velocity, and Variability are few Characteristics of Bigdata
- Improved customer service, better operational efficiency, Better Decision Making are few advantages of Bigdata

# Chapter 2: What is Hadoop? Introduction, Architecture, Ecosystem, Components

## What is Hadoop?

Apache Hadoop is an open source software framework used to develop data processing applications which are executed in a distributed computing environment.

Applications built using HADOOP are run on large data sets distributed across clusters of commodity computers. Commodity computers are cheap and widely available. These are mainly useful for achieving greater computational power at low cost.

Similar to data residing in a local file system of a personal computer system, in Hadoop, data resides in a distributed file system which is called as a **Hadoop Distributed File system**. The processing model is based on **'Data Locality'** concept wherein computational logic is sent to cluster nodes(server) containing data. This computational logic is nothing, but a compiled version of a program written in a high-level language such as Java. Such a program, processes data stored in Hadoop HDFS.

***Do you know?*** Computer cluster consists of a set of multiple processing units (storage disk + processor) which are connected to each other and acts as a single system.

# Hadoop EcoSystem and Components

Below diagram shows various components in the Hadoop ecosystem-



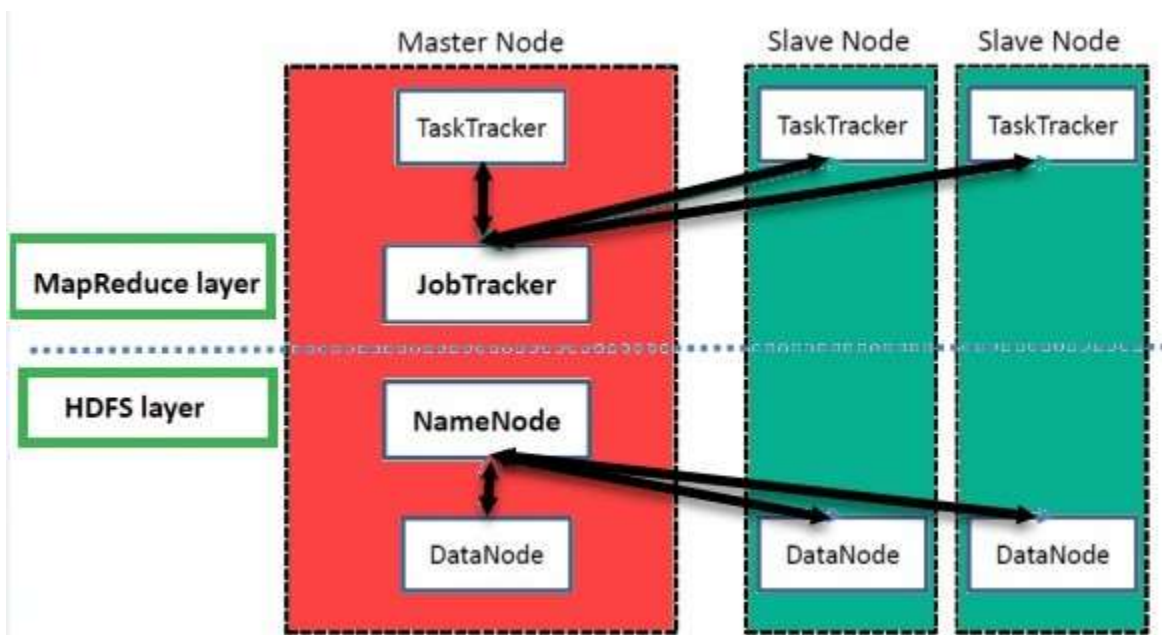Apache Hadoop consists of two sub-projects –

1. **Hadoop MapReduce:** MapReduce is a computational model and software framework for writing applications which are run on Hadoop. These MapReduce programs are capable of processing enormous data in parallel on large clusters of computation nodes.

2. **HDFS** (**Hadoop Distributed File System**): HDFS takes care of the storage part of Hadoop applications. MapReduce applications consume data from HDFS. HDFS creates multiple replicas of data blocks and distributes them on compute nodes in a cluster. This distribution enables reliable and extremely rapid computations.

Although Hadoop is best known for MapReduce and its distributed file system- HDFS, the term is also used for a family of related projects that fall under the umbrella of distributed computing and large-scale data processing. Other Hadoop-related projects at Apache include are **Hive, HBase, Mahout, Sqoop, Flume, and ZooKeeper.**

# Hadoop Architecture



High Level Hadoop Architecture

Hadoop has a Master-Slave Architecture for data storage and distributed data processing using MapReduce and HDFS methods.

**NameNode:**

NameNode represented every files and directory which is used in the namespace

**DataNode:**

DataNode helps you to manage the state of an HDFS node and allows you to interacts with the blocks

**MasterNode:**

The master node allows you to conduct parallel processing of data using Hadoop MapReduce.

**Slave node:**

The slave nodes are the additional machines in the Hadoop cluster which allows you to store data to conduct complex calculations.
Moreover, all the slave node comes with Task Tracker and a DataNode. This allows you to synchronize the processes with the NameNode and Job Tracker respectively.

In Hadoop, master or slave system can be set up in the cloud or on- premise

# Features Of 'Hadoop'

**• Suitable for Big Data Analysis**

As Big Data tends to be distributed and unstructured in nature, HADOOP clusters are best suited for analysis of Big Data. Since it is processing logic (not the actual data) that flows to the computing

nodes, less network bandwidth is consumed. This concept is called as **data locality concept** which helps increase the efficiency of Hadoop based applications.
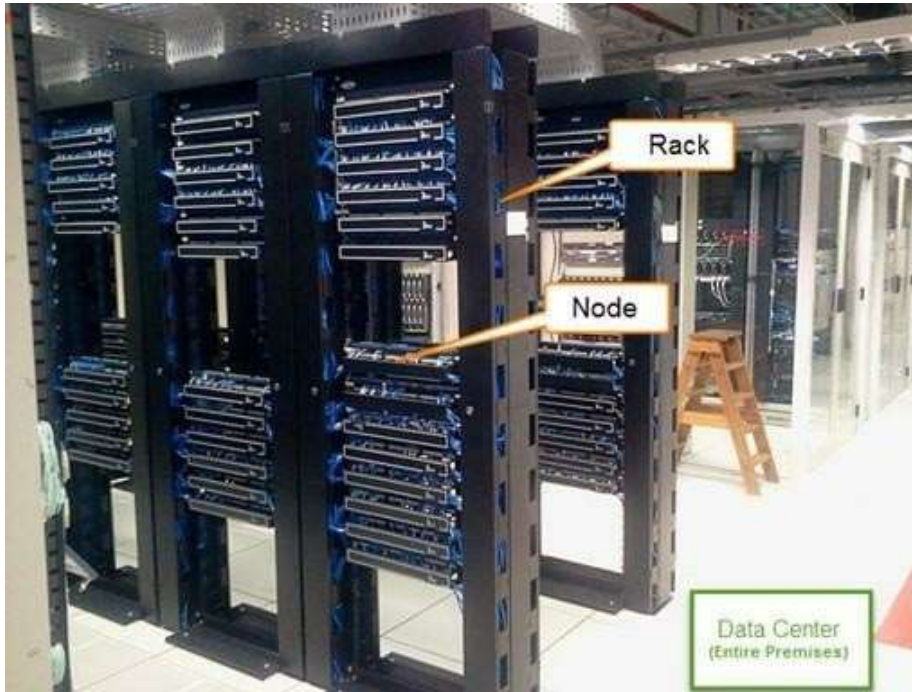
• **Scalability**

HADOOP clusters can easily be scaled to any extent by adding additional cluster nodes and thus allows for the growth of Big Data. Also, scaling does not require modifications to application logic.

• **Fault Tolerance**

HADOOP ecosystem has a provision to replicate the input data on to other cluster nodes. That way, in the event of a cluster node failure, data processing can still proceed by using data stored on another cluster node.

# Network Topology In Hadoop

Topology (Arrangment) of the network, affects the performance of the Hadoop cluster when the size of the Hadoop cluster grows. In addition to the performance, one also needs to care about the high availability and handling of failures. In order to achieve this Hadoop, cluster formation makes use of network topology.

Typically, network bandwidth is an important factor to consider while forming any network. However, as measuring bandwidth could be difficult, in Hadoop, a network is represented as a tree and distance between nodes of this tree (number of hops) is considered as an important factor in the formation of Hadoop cluster. Here, the distance between two nodes is equal to sum of their distance to their closest common ancestor.

Hadoop cluster consists of a data center, the rack and the node which actually executes jobs. Here, data center consists of racks and rack consists of nodes. Network bandwidth available to processes varies depending upon the location of the processes. That is, the bandwidth available becomes lesser as we go away from-

- Processes on the same node
- Different nodes on the same rack
- Nodes on different racks of the same data center
- Nodes in different data centers