

Unit 05 Introduction to Big Data.

Big Data.

Big data refers to large and complex datasets that can not be efficiently managed, processed, and analyzed using traditional data processing techniques.

Types of big data.

- (1) Structured.
- (2) Unstructured.
- (3) Semi-structured.

5 V's of Big data.

Volume: volume refers to the vast amount of data generated and collected. It includes both structured and unstructured data. With the growth of the technology and the internet, data volumes have significantly increased, leading to storage and processing challenges.

Velocity.

Velocity refers to the speed at which data is generated, collected and processed. It refers to the real time or near real time nature of data streams, where the data needs to be analyzed and acted upon quickly. Examples include streaming of data from social media, sensors, financial transactions and more.

Variety

Variety refers to the diversity of the data types and sources. It encompasses structured, unstructured, and semi-structured data. The variety of data presents challenges in terms of data integration, management, and analysis.

Veraelity

Veraelity related to the quality and reliability of data. It refers to the trustworthiness, accuracy and consistency of data.

Value

Value represents the ultimate goal of big data analytics. It refers to the ability to extract meaningful insights, patterns, and value from big data.

challenges

- ① Lack of understanding of big data.
- ② Dealing with data growth.
- ③ Confusion while big data tool selection.
- ④ Generating insights in timely manner.
- ⑤ Recruiting and retaining big data talent.
- ⑥ Integrating disparate data sources.
- ⑦ Securing big data.
- ⑧ Organizational resistance.

Commonly used tools for Big data.

- ① Hadoop
- ② Map Reduce programming.
- ③ HDFS.
- ④ Spark.
- ⑤ Apache Hive.
- ⑥ Apache Pig.
- ⑦ Apache Kafka.
- ⑧ HBase.
- ⑨ NoSQL Database: MongoDB, HBase, Cassandra etc.

Map Reduce.

Map Reduce is a software framework and programming models used for processing huge amount of data. Map Reduce program work in two phase, namely map and reduce.

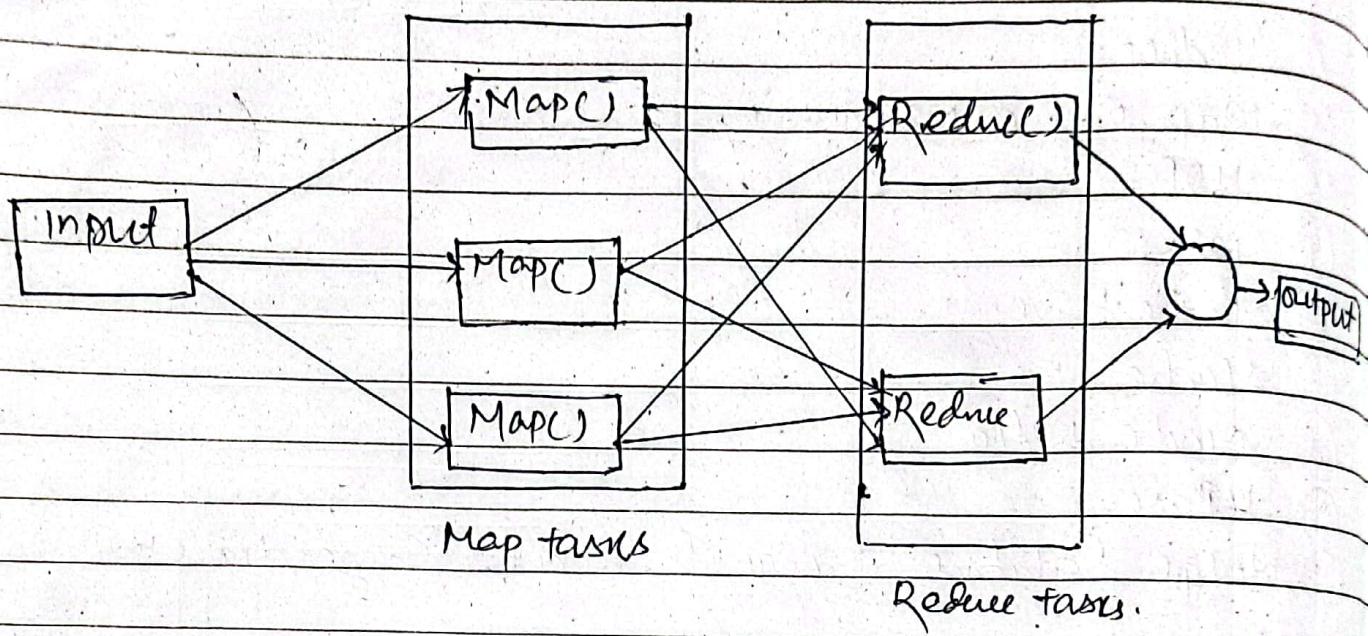
Map

It's tasks deal with splitting and mapping of data. It takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key-value pairs).

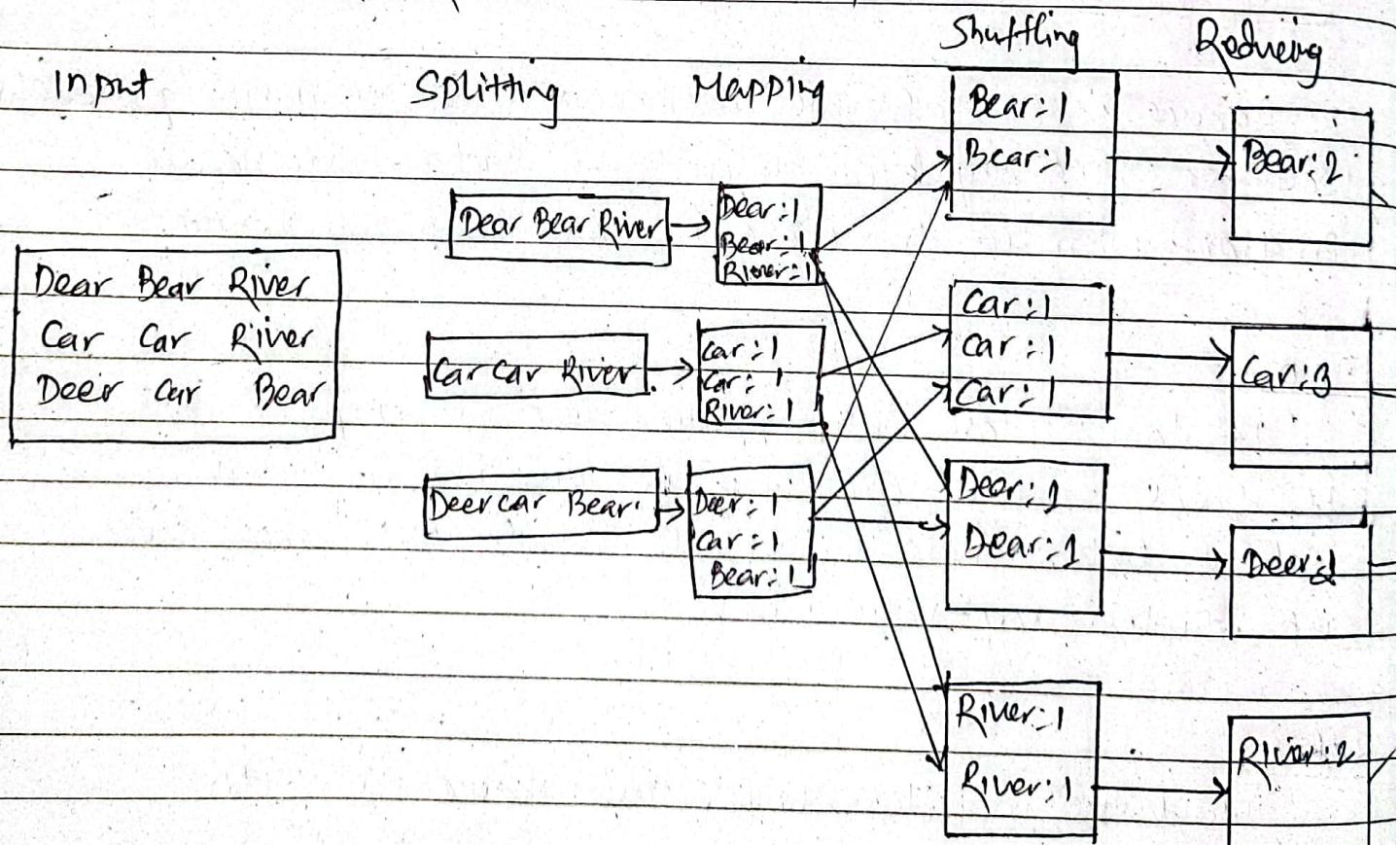
Reduce

It's tasks includes shuffle and reduce the data. The reduce tasks takes the output from the map as an input and combines those data tuples (key-value pairs) into a smaller set of tuples.

Input



Example.: Map Reduce word count process.



Apache Hadoop is an open source framework that is used to efficiently store and process large datasets ranging from in size from gigabyte to petabyte of data.

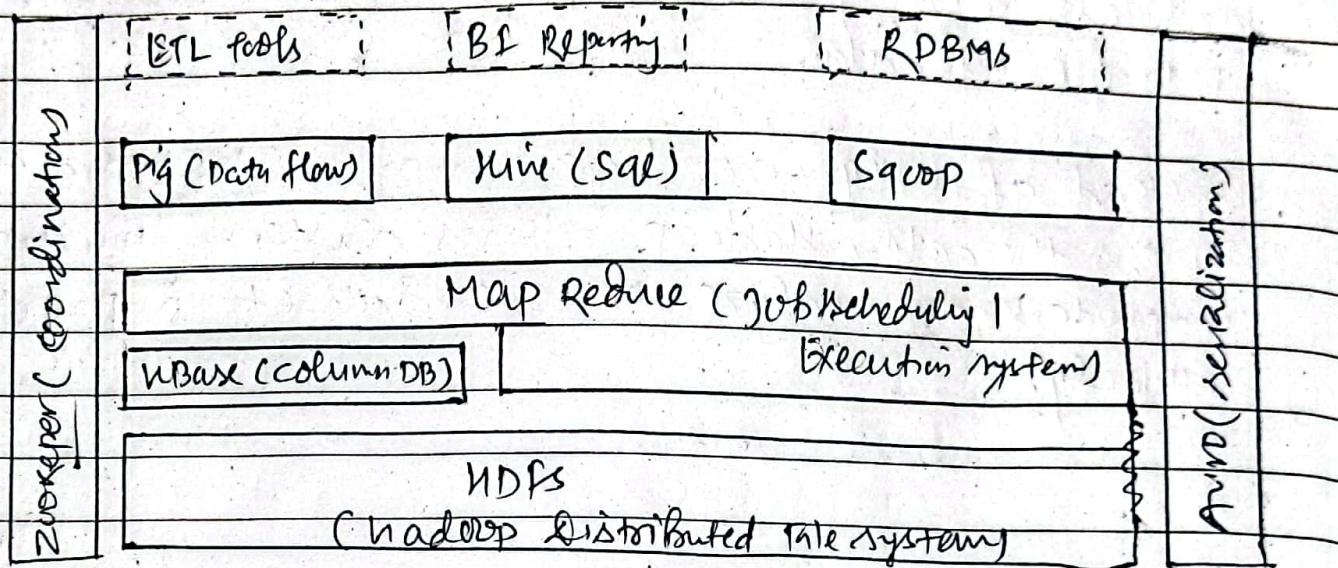
Instead of using one large computer to store and process the data, Hadoop allows clustering multiple computers to analyze massive datasets in parallel quickly.

Final Result:

Bear, 2
Car, 3
Deer, 2
River, 2



Hadoop Ecosystem



Hadoop Ecosystem is shown in figure above.

The Hadoop Ecosystem consists of following components

- ① Hadoop Common.
- ② Hadoop Distributed file system(HDFS)
- ③ YARN.
- ④ Hadoop Map Reduce.

Hadoop Common.

It consists of the common utilities that support other Hadoop modules.

- It is also known as Hadoop core.
- ~~It is also known.~~
- It provides file system and operating system level abstractions.

It contains Java archive files and scripts needed to start Hadoop.

YARN

- It provides the solution to the scalability problems of classical MapReduce.
- It is framework for job scheduling and cluster resource management.
- It enables Hadoop to support more varied processing approach and broader array of applications.

MapReduce

It is method of distributing computations across multiple nodes.

Each node processes the data that is stored at that node.

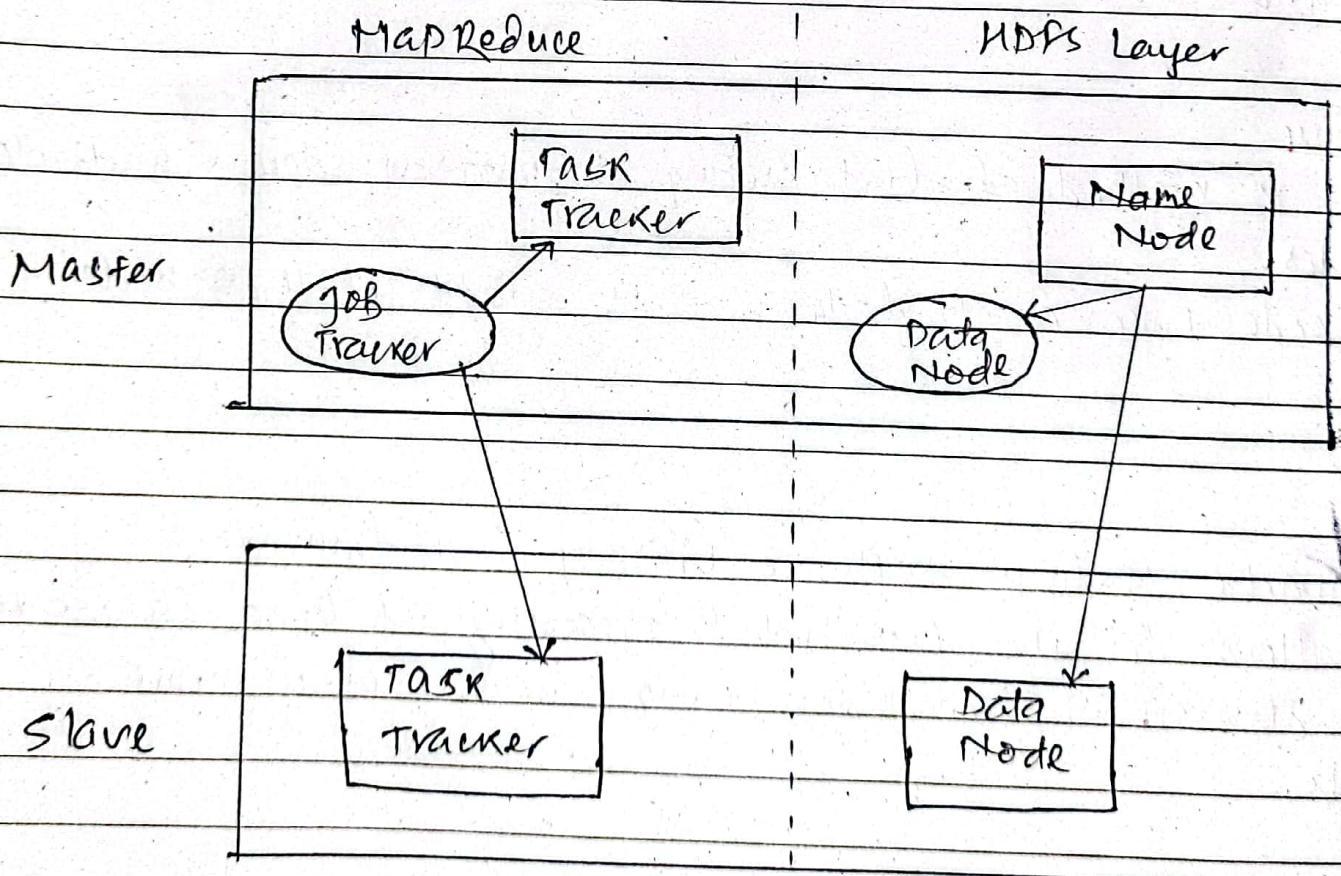
Note

The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models.

Hadoop Architecture

The hadoop architecture is a package of the file system, MapReduce engine and HDFS (Hadoop Distributed File System). The MapReduce engine can be MapReduce/MR1 or YARN/MR2.

A hadoop cluster consists of a single master and multiple slave nodes. The master node includes job tracker, task tracker, NameNode and DataNode whereas the slave node includes DataNode and Task tracker.

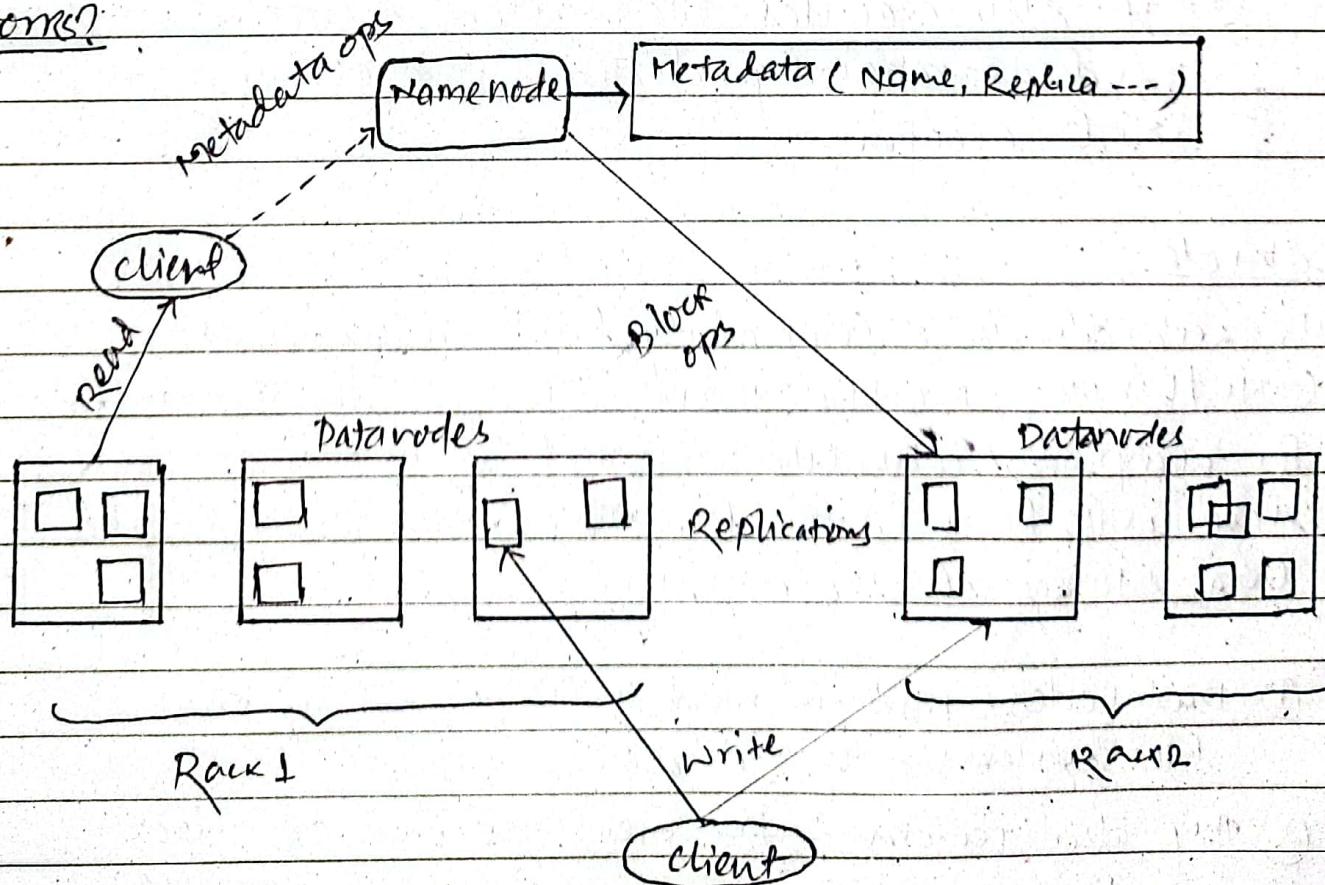


HDFS

The Hadoop distributed file system (HDFS) is the primary data storage system used by Hadoop applications. HDFS employs a Namenode and Datanode architecture to implement a distributed filesystem that provides high performance access to data across highly scalable Hadoop clusters.

Hadoop itself is an open source distributed processing framework that manages data processing and storage for big data applications. HDFS is a key part of the many Hadoop ecosystem technologies. It provides a reliable means for managing pools of big-data and supporting related big data analytics applications.

WORKS?



HDFS follows the master-slave architecture and it has the following elements.

NameNode

The NameNode is commodity hardware that contains the GNU/Linux operating system and the NameNode software. It is a software that can be run on commodity hardware.

The system having NameNode acts as the master server and does the following tasks.

- ① Manages the file system namespace.
- ② Regulates client access to file.
- ③ It also executes file system operations such as renaming, closing, and opening files and directories.

DataNode

The DataNode is a commodity hardware having the GNU/Linux operating system and DataNode software.

For every node (commodity hardware | system) in a cluster, there will be a DataNode. These nodes manage the data storage of their system.

- ① DataNodes perform read-write operations on the file systems as per client request.
- ② They also perform such operations such as block operations creation, deletion, and replication according to the instructions of NameNode.

Blocks.

Blocks are nothing but the smallest continuous location on your hard drive where data is stored.

In general, in any of the file system, you store the data as a collection of blocks. Similarly, HDFS stores each file as blocks which are scattered throughout the Apache Hadoop cluster.