**HADOOP**

Hadoop is an open-source framework that allows to store and process big data in a distributed environment across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. It is provided by Apache to process and analyze very huge volume of data. It is written in Java.

Hadoop Ecosystem has been hailed for its reliability and scalability. With the massive increase in information, it becomes increasingly difficult for the database systems to accommodate growing information. Hadoop provides a scalable and a fault-tolerant architecture that allows massive information to be stored without any loss. Hadoop fosters two types of scalability:

Vertical Scalability - In vertical scaling, we add more resources (like CPUs) to the single node. In this way, we increase the hardware capacity of our Hadoop system. We can further add more RAM and CPU to it in order to enhance its power and make it more robust.

Horizontal Scalability - In Horizontal Scaling, we add more nodes or systems to the distributed software system. Unlike vertical scalability's method of increasing capacity, we can add more machines without halting the system. This eliminates the issue of downtime and gives maximum efficiency while scaling out. This also renders multiple machines that are working in parallel

**COMPONENTS**

**Hadoop Distributed File system**

HDFS is the storage unit of Hadoop. Hadoop is a collection of master-slave networks. In HDFS there are two daemons - name node and data node that run on the master and slave nodes respectively. In HDFS data is distributed over several machines and replicated to ensure their durability to failure and high availability to parallel application. It is cost effective as it uses commodity hardware. It involves the concept of blocks, data nodes and node name

Blocks - A Block is the minimum amount of data that it can read or write. HDFS blocks are 128 MB by default and this is configurable. Files n HDFS are broken into block-sized chunks which are stored as independent units.

Name Node - HDFS works in master-worker pattern where the name node acts as master. Name Node is controller and manager of HDFS as it knows the status and the metadata of all the files in HDFS; the metadata information being file permission, names and location of each block. The metadata are small, so it is stored in the memory of name node, allowing faster access to data. More over the HDFS cluster is accessed by multiple clients concurrently, so all this information is handled by a single machine. The file system operations like opening, closing, renaming etc. are executed by it.

Data Node - They store and retrieve blocks when they are told to; by client or name node. They report back to name node periodically, with list of blocks that they are storing. The data node being a commodity hardware also does the work of block creation, deletion and replication as stated by the name node.

**YARN**

It is used for resource management and job scheduling. In a multi-node cluster, it is difficult to manage, allocate and release the resources. Hadoop Yarn allows to manage and control these resources very efficiently.

Yet Another Resource Manager takes programming to the next level beyond Java, and makes it interactive to let another application Hbase, Spark etc. to work on it. Different Yarn applications can co-exist on the same cluster so MapReduce, Hbase, Spark all can run at the same time bringing great benefits for manageability and cluster utilization.

**Map-Reduce**

This part of Hadoop is responsible for high-level data processing. It facilitates processing of a large amount of data over the cluster of nodes. It is the core component of Hadoop, which divides the big data into small chunks and process them parallelly.

MapReduce is a framework using which we can write functions to process massive quantities of data in parallel on giant clusters of commodity hardware in a dependable manner. It is a data processing tool which is used to process the data parallelly in a distributed form.

The MapReduce is a paradigm which has two phases, the mapper phase, and the reducer phase. In the Mapper, the input is given in the form of a key-value pair. The output of the Mapper is fed to the reducer as input. The reducer runs only after the Mapper is over. The reducer too takes input in key-value format, and the output of reducer is the final output.\

Flow

- Map function

The map function process the upcoming key-value pairs and generated the corresponding output key-value pairs. The map input and output type may be different from each other.

- Partition function

The partition function assigns the output of each Map function to the appropriate reducer. The available key and value provide this function. It returns the index of reducers.

- Shuffling and Sorting

The data are shuffled between/within nodes so that it moves out from the map and get ready to process for reduce function. Sometimes, the shuffling of data can take much computation time.

The sorting operation is performed on input data for Reduce function. Here, the data is compared using comparison function and arranged in a sorted form.

- Reduce function

The Reduce function is assigned to each unique key. These keys are already arranged in sorted order. The values associated with the keys can iterate the Reduce and generates the corresponding output.

- Output writer

Once the data flow from all the above phases, Output writer executes. The role of Output writer is to write the Reduce output to the stable storage.


**Analytics Modules**

Hive

Apache Hive is a data ware house system for Hadoop that runs SQL like queries called HQL (Hive query language) which gets internally converted to map reduce jobs. Hive was developed by Facebook. It supports Data Definition Language, Data Manipulation Language and user defined functions. It converts the query into map-reduce functions.


Pig

Pig is a high-level data flow platform for executing Map Reduce programs of Hadoop. It was developed by Yahoo. It converts the query into map-reduce functions.


Spark

Spark is a unified analytics engine for large-scale data processing including built-in modules for SQL, streaming, machine learning and graph processing. Unlike MapReduce, it can deal with real-time processing. It is an open-source framework used for faster data processing. It has Spark SQL as its very own query language. Spark is a fast and general processing engine compatible with Hadoop data. It can run in Hadoop clusters through YARN or Spark's standalone

mode, and it can process data in HDFS, HBase, Cassandra, Hive. Many organizations run Spark on clusters of thousands of nodes.

Sqoop

Sqoop is an open source framework provided by Apache. It is a command-line interface application for transferring data between relational databases and Hadoop.  It is used to import data from relational databases such as MySQL, Oracle to Hadoop HDFS, and export from Hadoop file system to relational databases.

An example use case of Sqoop is an enterprise that runs a nightly Sqoop import to load the day's data from a production transactional RDBMS into a Hive data warehouse for further analysis.