## Group A

Assignments based on Hadoop

## Assignment No. 01

| Performance | Understanding | Regularity | Total | Dated Sign of Subject Teacher |
|---|---|---|---|---|
| 03 | 01 | 01 | 05 | |
| | | | | |

**Date of Performance:** ...................................      **Date of Completion**: ..........................

**Title:** Single node/Multiple node Hadoop Installation.

**Objectives:**

1. To understand Big data primitives and fundamentals.
2. To understand the different Big data processing techniques

**Problem Statement:**

Hadoop Installation on:
  a. Single Node
  b. Multiple Node

**Outcomes:**

*Students will be able to,*

1. Apply the Big data primitives and fundamentals for application development.
2. Explore different Big data processing techniques with use cases

**Software and Hardware requirements:**

1. Software: Ubuntu OS / Windows, Hadoop 2.6.0, jdk7, Eclipse
2. Hardware: Processor, Ethernet Connection or WiFi, RAM 1GB, HDD, Sound Card, camera, microphone (depending upon website selection)
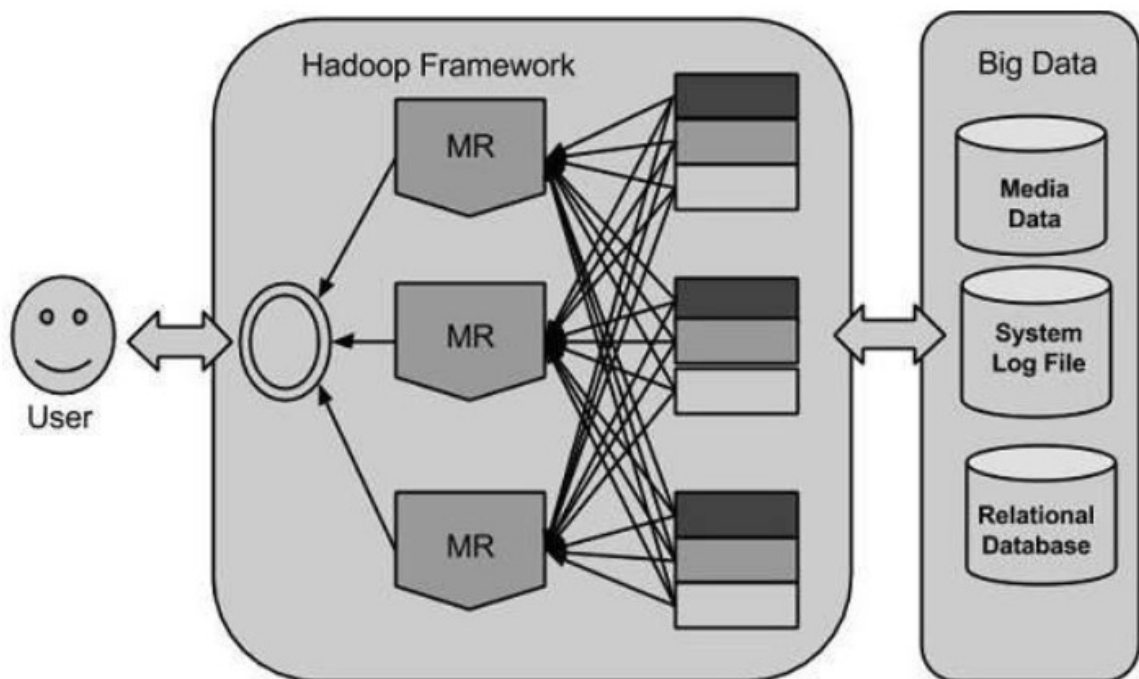
**Theory:**

Introduction 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. Big Data Big data means really a big data; it is a collection of large datasets that cannot be processed using traditional computing techniques. Big data is not merely a data; rather it has become a complete subject, which involves various tools, techniques and frameworks. 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.
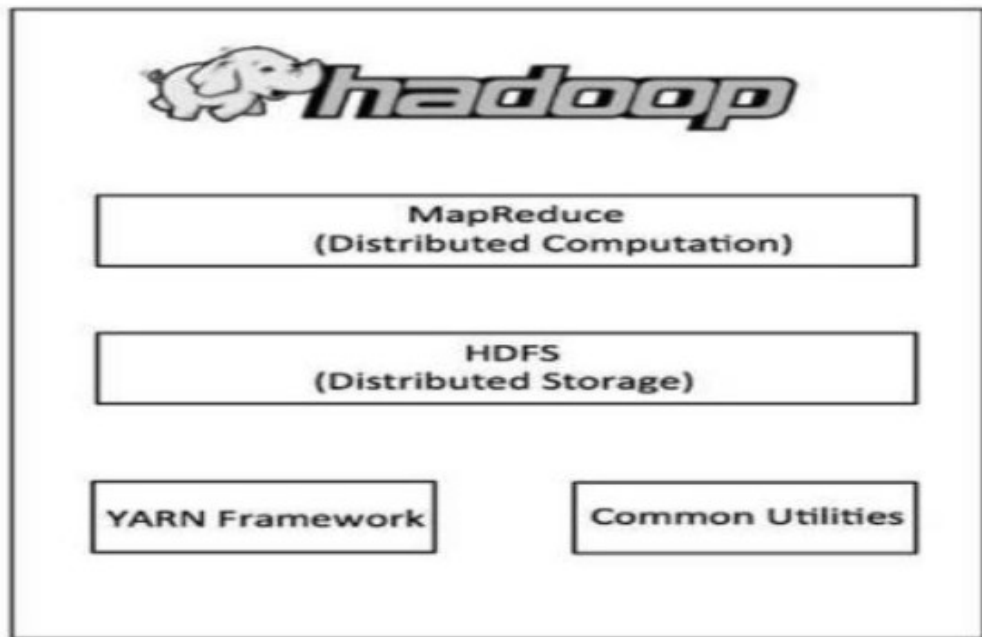
**Hadoop:**

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. A Hadoop frame-worked 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. Hadoop runs applications using the MapReduce algorithm, where the data is processed in parallel on different CPU nodes. In short, Hadoop framework is capable enough to develop applications capable of running on clusters of computers and they could perform complete statistical analysis for a huge amount of data.



**Hadoop Architecture:**

Hadoop framework includes following four modules:

- **Hadoop Common:** These are Java libraries and utilities required by other Hadoop modules. These libraries provide file system and OS level abstractions and contains the necessary Java files and scripts required to start Hadoop.

- **Hadoop YARN**: This is a framework for job scheduling and clusterresource management. Hadoop Distributed File System (HDFS): A distributed file system that provides high-throughput access to application data.

- **Hadoop MapReduce:** This is YARN-based system for parallel processing of large data sets.

## Steps of Installation and configuration of Hadoop:

### Step 1: Download and install Java

Hadoop is built on Java, so you must have Java installed on your PC. You can get the most recent version of Java from the official website. After downloading, follow the installation wizard to install Java on your system.

JDK: https://www.oracle.com/java/technologies/javase-downloads.html

### Step 2: Download Hadoop

Hadoop can be downloaded from the Apache Hadoop website. Make sure to have the latest stable release of Hadoop. Once downloaded, extract the contents to a convenient location.

Hadoop: https://hadoop.apache.org/releases.html

### Step 3: Set Environment Variables

You must configure environment variables after downloading and unpacking Hadoop. Launch the Start menu, type "Edit the system environment variables," and select the result. This will launch the System Properties dialogue box. Click on "Environment Variables" button to open.

Click "New" under System Variables to add a new variable. Enter the variable name "HADOOP_HOME" and the path to the Hadoop folder as the variable value. Then press "OK."

Then, under System Variables, locate the "Path" variable and click "Edit." Click "New" in the Edit Environment Variable window and enter "%HADOOP_HOME%bin" as the variable value. To close all the windows, use the "OK" button.

## Step 4: Setup Hadoop

You must configure Hadoop in this phase by modifying several configuration files. Navigate to the "etc/hadoop" folder in the Hadoop folder. You must make changes to three files:

- core-site.xml
- hdfs-site.xml
- mapred-site.xml

**Open each file in a text editor and edit the following properties:**

### In core-site.xml

```xml
<configuration>
  <property>
    <name>fs.default.name</name>
    <value>hdfs://localhost:9000</value>
  </property>
</configuration>
```

### In hdfs-site.xml

```xml
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file:/hadoop-3.3.1/data/namenode</value>
  </property>
  <property>
    <name>dfs.datanode.data.dir</name>
    <value>file:/hadoop-3.3.1/data/datanode</value>
  </property>
</configuration>
```

### In mapred-site.xml

```xml
<configuration>
  <property>
    <name>mapred.job.tracker</name>
    <value>localhost:54311</value>
  </property>
</configuration>
```

Save the changes in each file.

## Step 5: Format Hadoop NameNode

You must format the NameNode before you can start Hadoop. Navigate to the Hadoop bin folder using a command prompt. Execute this command:

```
hadoop namenode -format
```

## Step 6: Start Hadoop

To start Hadoop, open a command prompt and navigate to the Hadoop bin folder. Run the following command:

```
start-all.cmd
```

This command will start all the required Hadoop services, including the NameNode, DataNode, and JobTracker. Wait for a few minutes until all the services are started.

## Step 7: Verify Hadoop Installation

To ensure that Hadoop is properly installed, open a web browser and go to http://localhost:50070/. This will launch the web interface for the Hadoop NameNode. You should see a page with Hadoop cluster information.

**CONCLUSION:**

We have studied Hadoop installation and configuration