# Big Data Hadoop & Spark

Session 1 Assignment

## Introduction

**Big Data**

According to Gartner, Big Data is a high-volume, high-velocity, and high-variety

information asset that demands cost-effective, innovative forms of information

processing for enhanced insight and decision making

**3 Vs of Big Data**

| Volume | Variety | Velocity |
|---|---|---|
| Data size | Data Sources | Speed of Change |
| • Terabyte | Structured | Batch |
| • Records | Unstructured | Near-Time |
| • Transactions | Semi Structured | Real time |
| • Tables/Files | All of the above | Streams |

### Exploding of the Problem

Big Data constitutes large data sets in petabytes & zettabyte which cannot be processed by a single machine within an expected timeframe.

## Big Data Challenges

- Big Data Storage
- Big Data Processsing
- Manual Distributed Computing

**Possible Solutions**

- Scale up
- Scale out

**Scale up**

- Increase the configuration of a single system, like disk capacity, RAM, data transfer speed, etc.
- Complex, costly, and a time consuming process

**Scale Out**
- Use multiple commodity (economical) machines and distribute the load of storage/processing among them.
- Economical and quick to implement as it focuses on distribution of load
- Instead of having a single system with 10 TB of storage and 80 GB of RAM, use 40 machines with 256 GB of storage and 2 GB of RAM.

## Solution For Big Data Explosion-Hadoop

### Apache Hadoop
Apache Hadoop is an open source framework which provides an automated distributed computing environment that supports storage of big data sets. It does that storage using a cluster of commodity machines. It then analyses this stored big data using a very simple programming model.

The storage mechanism is known as **HDFS** (Hadoop Distributed File System). It is based on google GFS(Google File System) white paper.

The analytical mechanism is known as **Map Reduce** and is based on google map reduce white paper.

### Apache Hadoop Philosophies

There are 4 basic philosophies on which hadoop works.

a) All the basic software that helps start a hadoop cluster is a software daemons.
b) All the above daemons are based on master and slave architecture.
c) The entire hadoop framework is divided into 2 broad parts - storage (HDFS) and processing (Map Reduce).

- Single Node Hadoop Installation (dev and testing)
- Mulinode Installation (Production)

### Hadoop 2.x

- **HDFS (Storage Part)**
-       **Master Daemon** - **Namenode** (High End Admin Machine) (1 in number)
-       Back up Master Daemon - Secondary Namenode (High End Admin Machine) (1 in number)
-       **Slave Daemons** - **Datanode** (Commodity Machines) (**Many in number)**
- **Map Reduce (Processing Part)** - YARN (Yet Another Resource Negotiater)
-       **Master Daemon** - **ResourceManager** (High End Admin Machine) (1 in number)
-       **Slave Daemon** - **NodeManager** (Commodity Machines) (Many in number)

d) Hadoop is a batch oriented system which can never be plugged behind and online transaction processing system. Moreover, it a write once, ready many times data storage mechanism. This means, you can never update the data. If you really want to update, you need to delete the previous version and upload a new copy.
Data Node is a slave daemon software for the storage part of hadoop.
Likewise, Node manager is a slave daemon software for the processing part of hadoop.
When you refer to hardware in hadoop, you always refer to it as either a commodity machines or a slave node.
Important point to remember, there is a difference between "Slave daemon" and "slave node" (which btw is also called as commodity hardware.)
Both the datanode slave daemon and node manager slave daemon run on a commodity machine or a slave machine (which is a hardware)
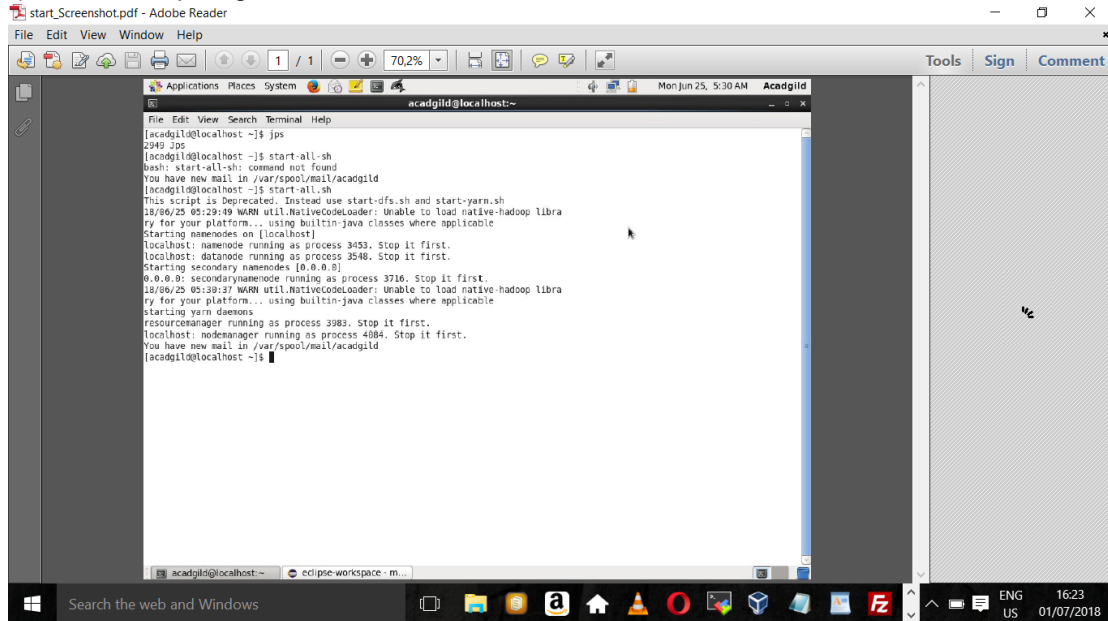
- Structured Data (Big and Small)
- Unstructured Data (Big and Small)
- Semi - Structure Data (Email, XML)

**10 GB - 1 Machine -1 hr to process the data**
**10 machines - 1 GB of Data in each machine (60/10=6 minutes)**

**Cluster - a combination of machines which collabrates to give you a combined processing power**
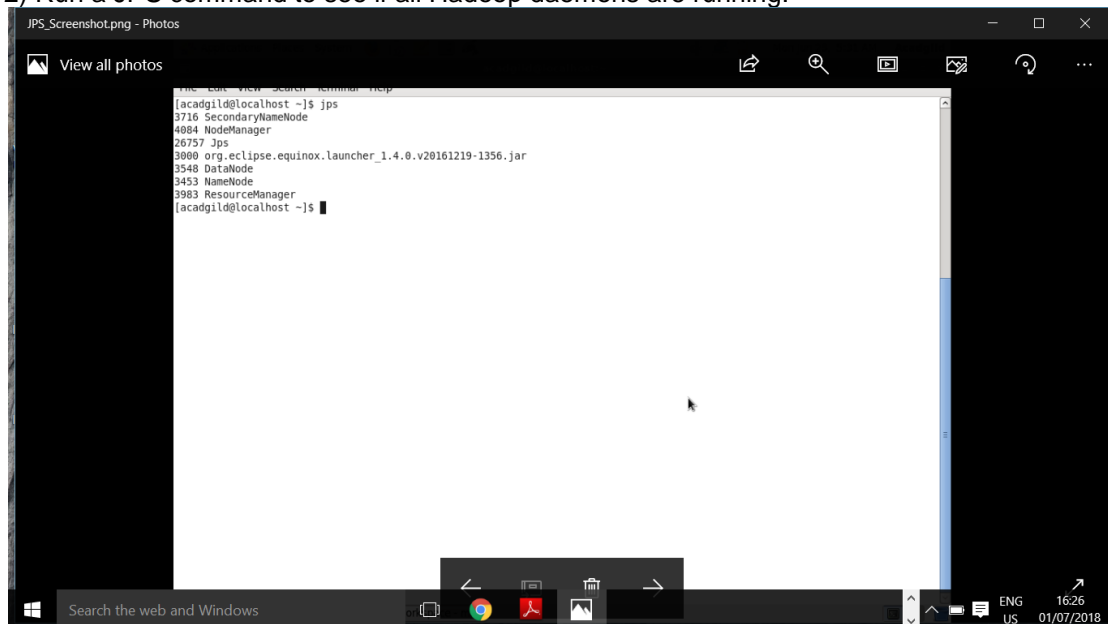
## TASK 2

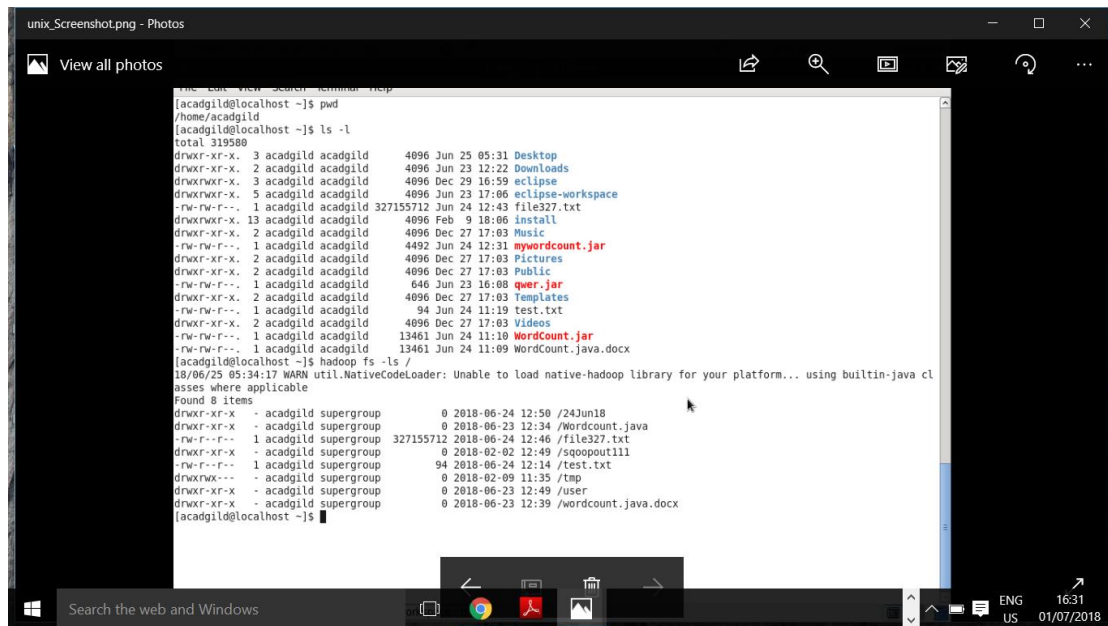1) Start Hadoop single node on AcadGild VM. The command is start-all.sh.



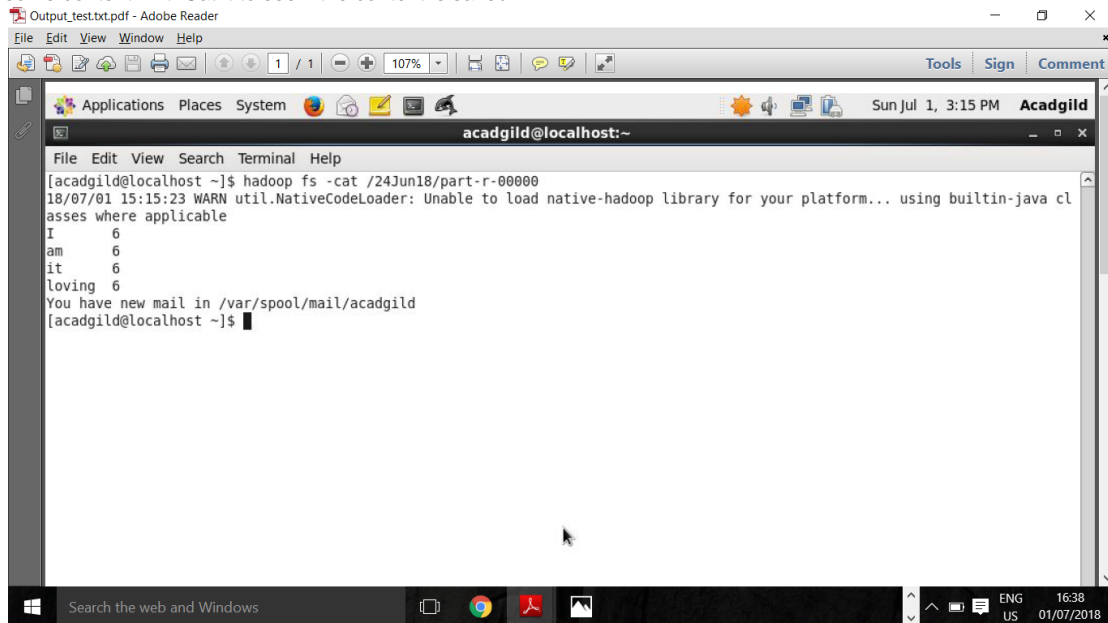2) Run a JPS command to see if all Hadoop daemons are running.



3) Run few Unix commands like pwd, ls -ls, etc.

## Explanation of Commands

- hadoop fs -ls / - Another useful command to display the list of files and directory in HDFS File path
- hadoop fs – mkdir /direcory_name – creates the directory in HDFS
- hadoop fs -touchz /directory-name/file_name – creates the file HDFS with a size of 0 bytes
- hadoop fs –cat /directory-name/file_name – displays file content to stdout
- hadoop fs –count /directory-name/ – displays the count of directories, files, bytes for path provided
- hadoop fs –put XYZ.txt /directory-name/XYZ.txt – Copies the local file to HDFS on the specified path
- hadoop fs –get /directory-name/XYZ.txt ./XYZ.txt – reverse of put
- hadoop fs -rm -r to remove the files and hadoop fs –rm to remove directories

4) Create a file from the terminal using nano editor (example: nano test.txt), and add some content in it. Cat it to see if the content is saved.



5)Open the hdfs web page by typing localhost:50070 in the browser. Check all the details of the HDFS.