

PEOPLE'S DEMOCRATIC REPUBLIC OF ALGERIA
MINISTRY OF HIGHER EDUCATION AND SCIENTIFIC RESEARCH
SAAD DAHLEB BLIDA 01 UNIVERSITY
DEPARTMENT OF COMPUTER SCIENCE



MASTER'S INTELLIGENT SYSTEMS ENGINEERING

DATA WAREHOUSE AND BIG DATA

REPORT

**Implementation of a Hadoop Multi-Node Cluster
and Spark Framework for Image Compression**

Implemented and written by

Abdelatif Mekri

Halima Nfidsa

Imad Eddine Boukader

Nahla Yasmine Mihoubi

Academic year : 2024-2025

Table of content

Table of content.....	1
Introduction.....	3
Project Objectives.....	3
Technologies Used.....	3
Fundamentals.....	4
Image Compression.....	4
Image Compression Techniques.....	4
Lossless Compression.....	4
Lossy Compression.....	5
Comparison of Lossless and Lossy Compression.....	5
How Hadoop Work for Image Compression.....	6
Storing Images in HDFS.....	6
Image Compression Using MapReduce.....	6
How Spark and Hadoop Work Together for Image Compression.....	7
Role of Hadoop (HDFS).....	8
1. Distributed Storage:.....	8
2. Data Availability:.....	8
3. Input/Output for Compression:.....	9
Role of Spark.....	9
1. Reading Images:.....	9
2. In-Memory Processing:.....	9
3. Parallel Compression:.....	9
4. Saving Results:.....	9
Benefits of Combining Hadoop and Spark.....	9
Project Impact.....	10
Installation and Configuration of the Environment.....	10
A- Virtualisation Environment.....	10
B- Setting up the Machines.....	11
1- Operating System Selection and Installation.....	11
2- Hadoop Installation and Configuration.....	11
2-1 Installing JAVA.....	11
2-2 Installing Hadoop.....	12
2-3 Making 'slave' machines.....	14
2-4 Configuring SSH for Passwordless Login.....	14
2-5 Getting IP addresses and renaming hosts.....	15
2-6 setting up the hosts files.....	16
2-7 configuring hadoop on master and streaming the configuration to the slaves...	

17	
2-8 launching hadoop.....	20
3- Spark installation and configuration.....	22
3-1 Installing JAVA.....	22
3-2 Download Apache Spark.....	22
3-3 Installing Spark.....	22
3-4 Configure Environment Variables.....	23
3-5 Update Workers file.....	23
3-6 update workers with the new config.....	24
3-7 launching Spark.....	25
Code and running.....	26
I. First approach : Hadoop & spark.....	26
Code.....	26
A- Installing and importing the necessary libraries.....	26
B- Implementing the solution.....	26
Starting hadoop and spark.....	28
Check the HDFS.....	28
Upload the images to the designated folder in HDFS.....	29
Job submission.....	29
Running and waiting for job compilation.....	29
Retrieving results.....	30
II. Second approach : Hadoop (Map-Reduce).....	31
Code.....	31
Compiling the java files.....	32
Starting hadoop.....	32
Check the HDFS.....	32
Upload the images to the designated folder in HDFS.....	32
Job submission.....	33
Running and waiting for job compilation.....	33
Downloading the output files.....	33
Fix the output file to see the results.....	34
Results.....	34
Conclusion and Discussion.....	35
Notes.....	36

Keywords :

Hadoop ,spark , image compression , mapreduce ,distributed storage, HDFS

Introduction

Image compression is an essential technique in modern data processing, especially for large-scale multimedia datasets. It allows for the efficient reduction of file sizes while preserving the visual quality of images. In this project, we explore the application of **Big Data technologies—Hadoop and Spark**—to implement a parallelized image compression pipeline using the **JPEG** format. These formats are widely recognized for their ability to compress images efficiently, maintaining high image quality at lower file sizes. The primary goal of this project is to demonstrate the effectiveness of Hadoop and Spark in optimizing the compression process for large image datasets, focusing on distributed and in-memory processing capabilities, respectively. By leveraging these two technologies, we aim to highlight the scalability, speed, and efficiency improvements in handling and compressing vast amounts of image data.

Project Objectives

Setup and Configuration of Hadoop and Spark Environments: Install and configure both Hadoop and Spark frameworks, ensuring they are ready for processing large image datasets in a distributed and in-memory fashion.

1. **Development and Implementation of Image Compression Algorithms:** Implement efficient image compression algorithms using the JPEG and JPEG2000 formats, ensuring minimal loss of visual quality while achieving significant file size reduction.
2. **Integration of Hadoop and Spark for Enhanced Performance:** Utilize Hadoop's distributed storage system and MapReduce framework alongside Spark's in-memory processing engine to parallelize the image compression workflow, optimizing resource usage and reducing processing times.
3. **Comparative Analysis of Hadoop vs. Spark:** Analyze the performance of Hadoop and Spark in the context of image compression, comparing aspects such as processing time, scalability, resource utilization, and overall efficiency when handling large datasets.

Technologies Used

- **Hadoop**



Hadoop  is an open-source Big Data framework designed to efficiently store and process large datasets across distributed clusters. In this project, Hadoop's components were used to facilitate the storage and parallel processing of images:

- **HDFS (Hadoop Distributed File System):** HDFS provides a robust and fault-tolerant distributed storage system, enabling the storage of massive

image files across multiple nodes. This makes it well-suited for large-scale image data processing, ensuring that even when some nodes fail, the data remains accessible and can be recovered.

- **MapReduce:** Hadoop's MapReduce framework was leveraged to divide the image compression task into smaller sub-tasks, which can be processed independently and in parallel. This parallelization speeds up the overall compression process and ensures that the system can handle larger datasets more efficiently.

- **Spark**  :

Apache Spark is a high-performance, in-memory data processing engine known for its speed and scalability. It processes data much faster than Hadoop by storing intermediate results in memory rather than writing them to disk. In this project, Spark was used to enhance the image compression process:

- **RDDs (Resilient Distributed Datasets):** RDDs are the core data abstraction in Spark. They allow for the distribution of image data across multiple nodes and enable parallel processing of image compression tasks. RDDs provide fault tolerance and the ability to perform transformations and actions on data in parallel.
- **Parallel Image Compression:** By utilizing Spark's powerful transformation capabilities (such as `map()`, `filter()`, and `reduce()`), the compression process was parallelized across all nodes in the cluster. This approach significantly accelerates the overall compression time, making it possible to handle large datasets efficiently and reduce computational overhead.

Fundamentals

Image Compression

Image compression is the process of reducing the size of an image file while maintaining its quality as much as possible. It is essential for efficient storage, transmission, and processing of images, especially in applications like cloud storage, web applications, and big data analytics.

Image Compression Techniques

Lossless Compression

Lossless compression methods retain the original image quality by encoding the data efficiently without losing any information. Common techniques include:

- **Run-Length Encoding (RLE)**
- **Huffman Coding**
- **Lempel-Ziv-Welch (LZW)**
- **PNG and GIF formats**

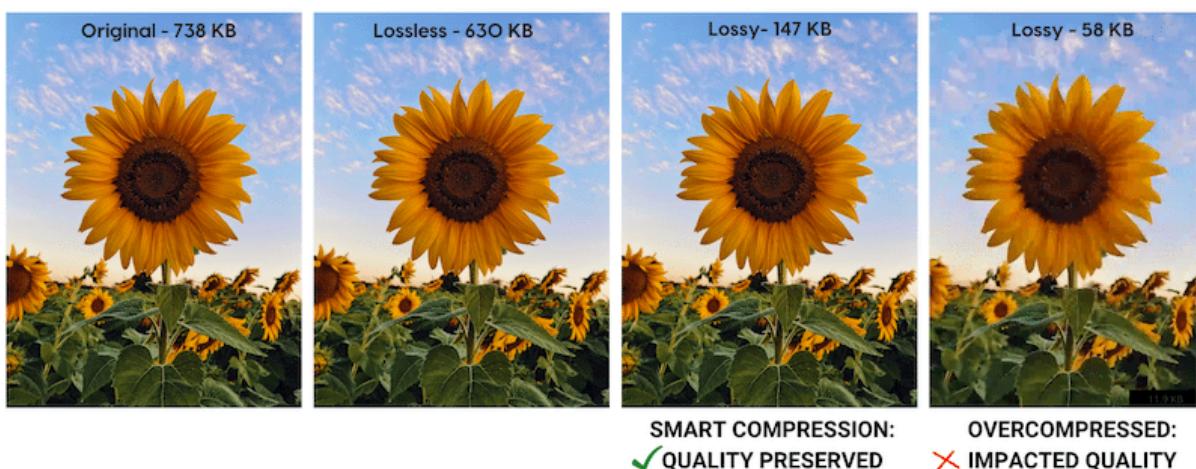
Lossy Compression

Lossy compression techniques reduce file size significantly by discarding some image data, leading to a slight loss in quality. Common methods include:

- **Discrete Cosine Transform (DCT)** (used in JPEG)
- **Wavelet Transform** (used in JPEG2000)
- **Quantization and Downsampling**

Comparison of Lossless and Lossy Compression

Type of Compression	Lossless Compression	Lossy Compression
Principle	Algorithms used to precisely reconstruct original data from compressed data.	Uses approximate estimations to represent content.
Usage	Text, programs, images, and sound.	Images, audio, and video.
Applications	RAW, BMP, PNG, WAV, FLAC, ALAC, et	JPEG, GIF, MP3, MP4, OGG, H.264, MKV, etc.
Storage Accuracy	High (original content preserved).	Reliable but with reduced quality (data discarded).



How Hadoop Work for Image Compression



Storing Images in HDFS

Images (JPEG files) are stored in Hadoop Distributed File System (HDFS), which breaks them into blocks and distributes them across the cluster.



A **traditional storage system**, such as a local file system (e.g., NTFS, Ext4), operates on a single machine and is limited by the system's hardware capacity. It provides fast access to files but lacks built-in fault tolerance and scalability. In contrast, the **Hadoop Distributed File System (HDFS)** is designed for distributed environments, storing data across multiple nodes to handle large-scale datasets efficiently. HDFS follows a **block-based architecture**, where files are split into large blocks (typically 128MB) and replicated across different nodes to ensure fault tolerance and high availability. Unlike traditional storage, HDFS is optimized for **sequential reads and parallel processing**, making it ideal for big data applications.

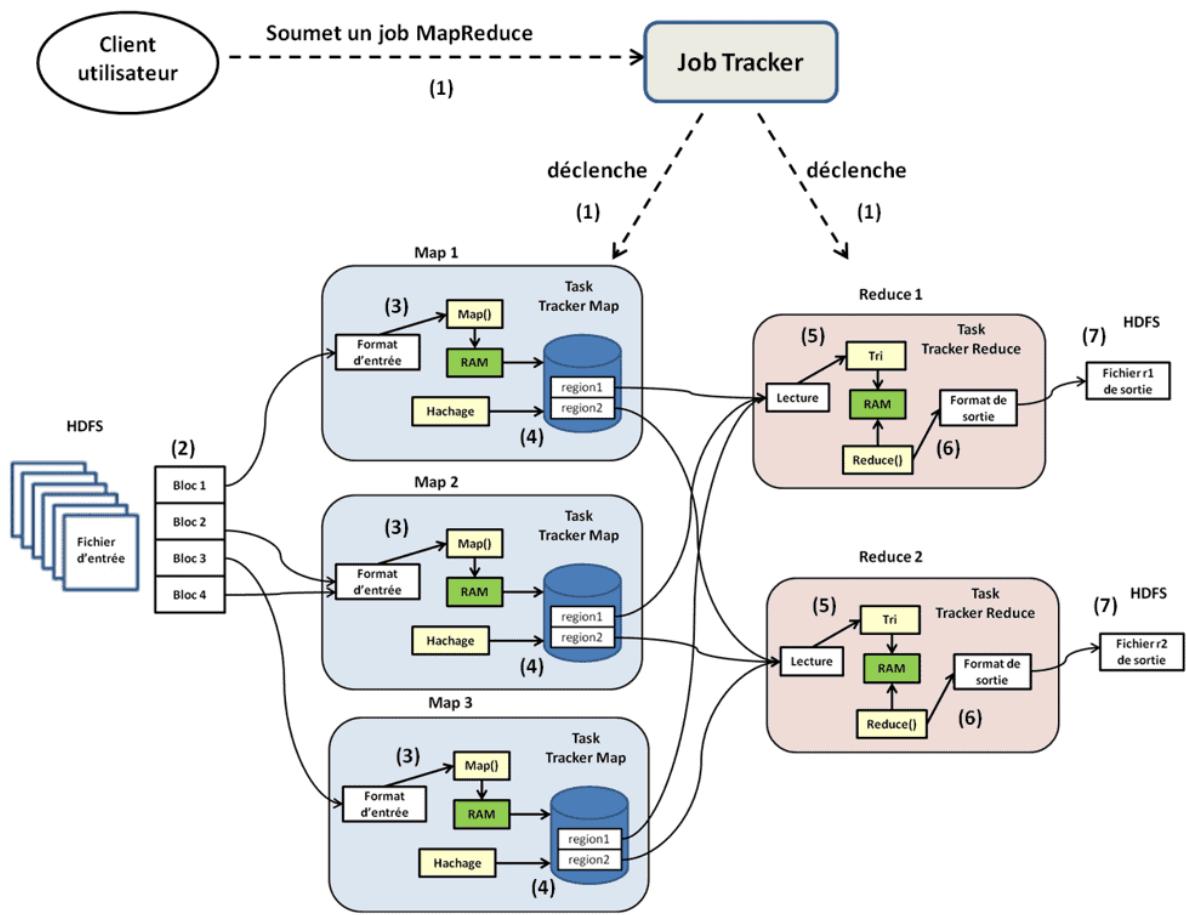
Additionally, while normal file systems allow frequent updates and modifications, HDFS follows a **write-once, read-many** approach, which enhances consistency and throughput in large-scale data processing frameworks like **MapReduce and Spark**. Its distributed nature allows for horizontal scalability, ensuring that storage and processing capabilities grow seamlessly with the addition of new nodes.

Image Compression Using MapReduce

Hadoop's MapReduce framework can be used to process images in parallel across the cluster. The process involves:

Map Phase: Each node processes a subset of images (or parts of images) to apply compression algorithms.

Reduce Phase: The compressed images are collected and stored back in HDFS.



to summarize we have 4 Steps for JPEG Compression in Hadoop

Step 1: Load images into HDFS.

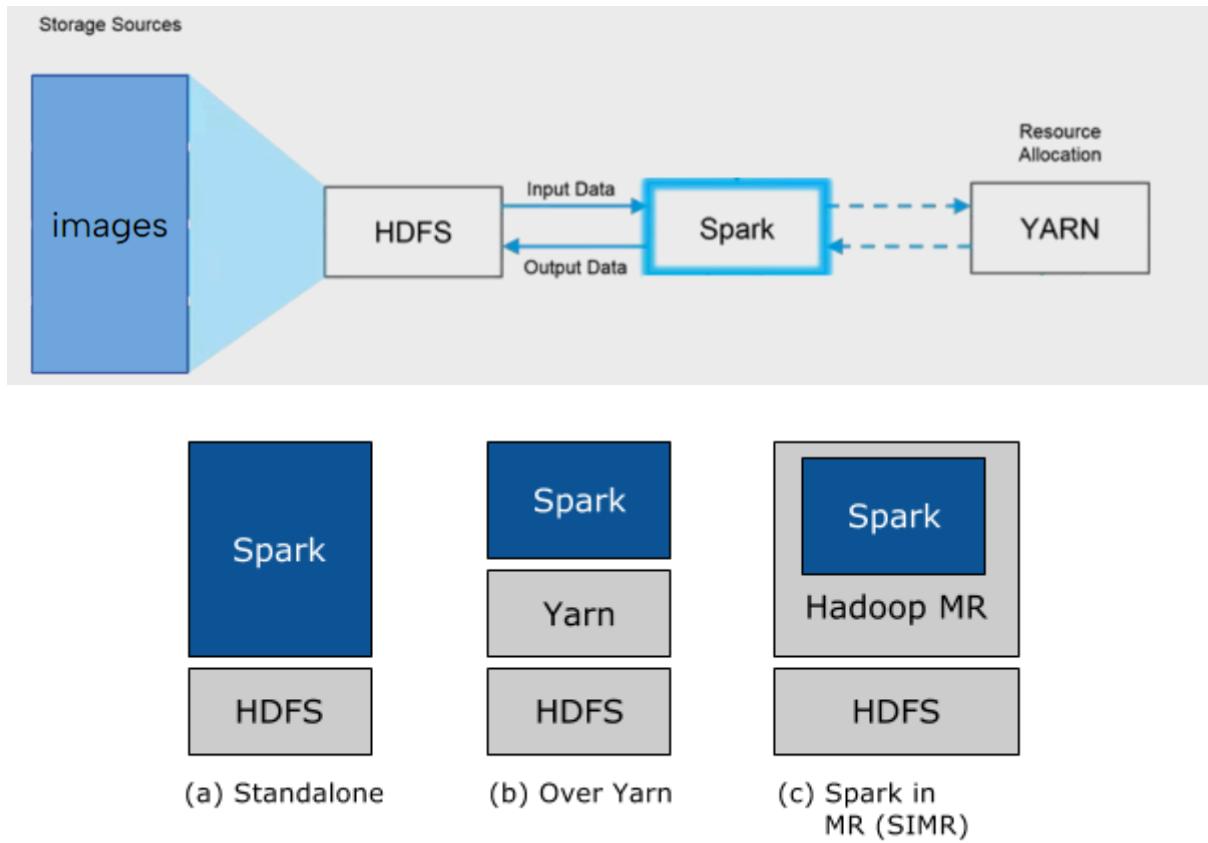
Step 2: Write a MapReduce job where:

The Mapper reads the images, applies compression (e.g., reducing quality or resizing), and outputs the compressed images.

The Reducer collects the compressed images and writes them back to HDFS.

How Spark and Hadoop Work Together for Image Compression

When combining Spark and Hadoop, their respective strengths Hadoop's distributed storage (HDFS) and Spark's in-memory processing are leveraged to build an efficient and scalable pipeline for image compression.



There are three ways to deploy Spark with Hadoop: **Standalone deployment**, where Spark runs alongside Hadoop MapReduce with manually allocated resources; **Hadoop YARN deployment**, which allows Spark to run on YARN without additional installation, integrating seamlessly into the Hadoop ecosystem; and **Spark In MapReduce (SIMR)**, which enables users without YARN to launch Spark jobs within MapReduce for quick experimentation. We chose the **Standalone deployment** because of its simplicity and ease of use, making it the preferred choice for many Hadoop users.

Role of Hadoop (HDFS)

Hadoop Distributed File System (HDFS) serves as the backbone for data storage. Here's how it contributes:

1. Distributed Storage:
 - Images are divided into blocks and distributed across a cluster of machines for redundancy and parallel access.
 - This ensures fault tolerance and scalability for large datasets.
2. Data Availability:
 - Spark reads the images directly from HDFS, leveraging its parallel access capabilities to process data faster.

3. Input/Output for Compression:
 - HDFS acts as both the input source for uncompressed images and the output destination for compressed images.

Role of Spark

Apache Spark handles the actual image compression with its powerful distributed computation framework:

1. Reading Images:
 - Spark reads the image files from HDFS into its distributed data structures, such as RDDs (Resilient Distributed Datasets) or DataFrames.
2. In-Memory Processing:
 - Unlike Hadoop's MapReduce, which writes intermediate results to disk, Spark performs all operations in memory, making it much faster.
 - For image compression, Spark processes multiple images simultaneously by partitioning the data across the cluster.
3. Parallel Compression:
 - Spark executes a map function across the RDD partitions, where each partition processes a subset of the images.
 - Compression is applied using libraries like Pillow or OpenCV, converting images to the desired format (e.g., JPEG).
4. Saving Results:
 - The compressed images are written back to HDFS, preserving the benefits of distributed storage.

Benefits of Combining Hadoop and Spark

- 1. Scalability:**
 - HDFS can handle petabytes of data, while Spark processes it efficiently in parallel.
- 2. Speed:**
 - Spark's in-memory processing makes it faster than Hadoop's disk-based MapReduce.
- 3. Fault Tolerance:**
 - Hadoop ensures data redundancy in HDFS, while Spark automatically handles task failures.
- 4. Ease of Integration:**
 - Spark can directly access data from HDFS without requiring additional setup.

Project Impact

By combining the strengths of **Hadoop** and **Spark**, this project demonstrates how distributed and in-memory processing can revolutionize the way large-scale image datasets are handled. Hadoop's ability to store vast amounts of data across multiple nodes and Spark's capacity to process data quickly in memory enables a highly efficient image compression pipeline. The comparison between Hadoop and Spark will provide insights into which technology is more suitable for different types of image processing tasks, ultimately helping to guide future decisions for large-scale data processing applications.

This work not only showcases the capabilities of Hadoop and Spark but also highlights the broader impact of Big Data technologies in the field of multimedia data processing, with potential applications in cloud storage, streaming services, and content delivery networks.

Installation and Configuration of the Environment

A- Virtualisation Environment

For this project, the environment was initially set up using Oracle Virtual Machine (VirtualBox). However,  VMware Workstation was chosen as the preferred virtualization solution due to its advanced features and better performance for multi-cluster setups. A total of five virtual machines (VMs) were configured, with one serving as the *master* node and the remaining four acting as *slave* nodes. This setup allowed for the distributed processing and testing of the Hadoop and Spark frameworks in a controlled and scalable environment.

Why VMware Workstation Was Chosen Over Oracle Virtual Machine

VMware Workstation offered several advantages over Oracle Virtual Machine, making it a better fit for this project:

1. Enhanced Performance: VMware provides better resource allocation and utilization, ensuring smoother operation of multiple VMs simultaneously.
2. Superior Networking Features: VMware's advanced networking capabilities simplified the configuration of communication between the master and slave nodes.
3. Snapshot and Cloning Functionality: VMware's snapshot feature allowed for easy backup and restoration of VM states, while cloning enabled quick replication of slave nodes.
4. Stability and Compatibility: VMware Workstation demonstrated greater stability and compatibility with the Hadoop and Spark frameworks, reducing setup and runtime issues.

5. User-Friendly Interface: VMware's intuitive interface streamlined the management of multiple VMs, saving time during configuration and deployment.

B- Setting up the Machines

1- Operating System Selection and Installation

ubuntu 20.04 LTS was chosen as the operating system for all virtual machines due to its stability, long-term support (LTS), and compatibility with the project requirements. Ubuntu 20.04 provided seamless support for Java 8 installation, which was essential for running Hadoop and Spark. Additionally, its robust package management system and extensive community support made it an ideal choice for setting up and maintaining the multi-cluster environment.

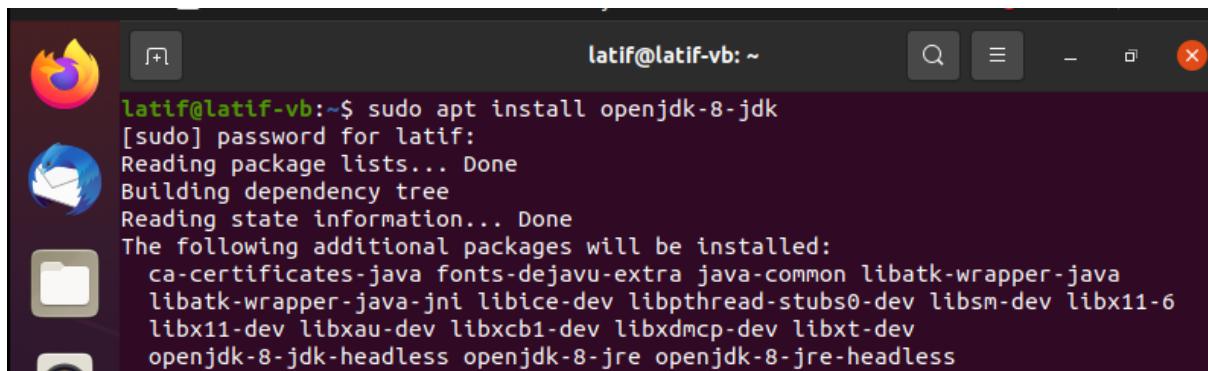
2- Hadoop Installation and Configuration

To set up the Hadoop framework, the following steps were performed on the master and slave nodes. The installation process was carried out on Ubuntu 20.04 LTS, ensuring compatibility with Java 8, which is a prerequisite for Hadoop.

2-1 Installing JAVA

Since Hadoop requires Java to run, OpenJDK 8 was installed on the master node using the following command:

```
sudo apt install openjdk-8-jdk
```

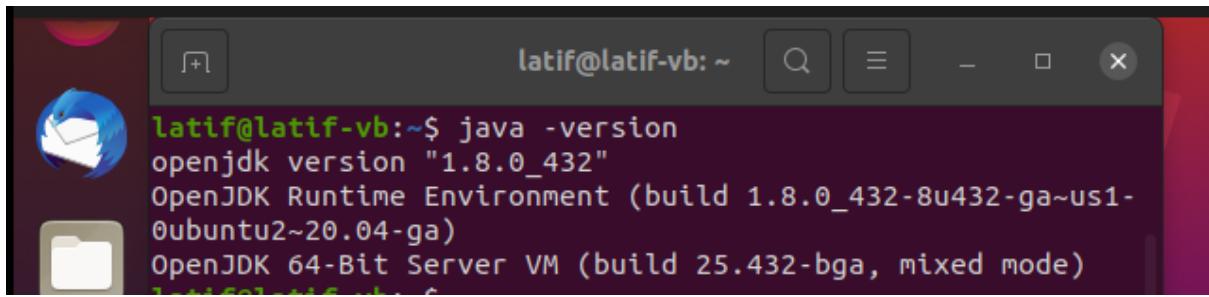


The screenshot shows a terminal window on an Ubuntu desktop. The title bar says 'latif@latif-vb: ~'. The terminal output is as follows:

```
latif@latif-vb:~$ sudo apt install openjdk-8-jdk
[sudo] password for latif:
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following additional packages will be installed:
  ca-certificates-java fonts-dejavu-extra java-common libatk-wrapper-java
  libatk-wrapper-java-jni libice-dev libpthread-stubs0-dev libsm-dev libx11-6
  libx11-dev libxau-dev libxcb1-dev libxdmcp-dev libxt-dev
  openjdk-8-jdk-headless openjdk-8-jre openjdk-8-jre-headless
```

Check the installed version of java with :

```
Java -version
```



```
latif@latif-vb:~$ java -version
openjdk version "1.8.0_432"
OpenJDK Runtime Environment (build 1.8.0_432-8u432-ga~us1-0ubuntu2~20.04~ga)
OpenJDK 64-Bit Server VM (build 25.432-bga, mixed mode)
```

2-2 Installing Hadoop

Open a browser on your system and navigate to the official Apache Hadoop website:

- <https://hadoop.apache.org/releases.html>

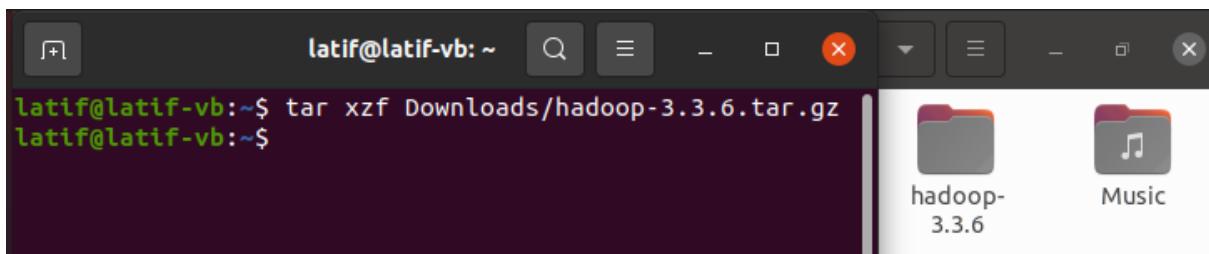
Select the latest stable release (in our case : Hadoop 3.3.6).

Download the **binary release** (.tar.gz file) by clicking the appropriate link for your system.

- **Extract the File**

After downloading, extract the Hadoop .tar.gz , using the commande

```
tar xzf Downloads/hadoop-3.2.1.tar.gz
```



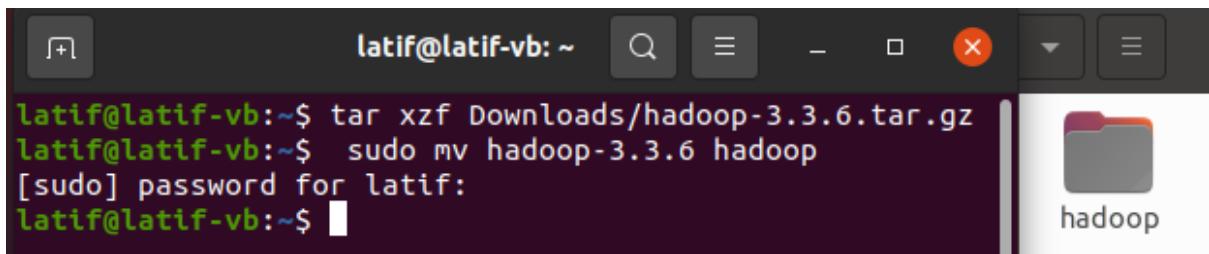
```
latif@latif-vb:~$ tar xzf Downloads/hadoop-3.3.6.tar.gz
latif@latif-vb:~$
```

The file explorer shows a folder named "hadoop-3.3.6" and a "Music" folder.

- **Rename the extracted folder**

Renaming it to something simpler or more convenient (e.g., **hadoop**), using the commande :

```
mv hadoop-3.3.6 hadoop
```



```
latif@latif-vb:~$ tar xzf Downloads/hadoop-3.3.6.tar.gz
latif@latif-vb:~$ sudo mv hadoop-3.3.6 hadoop
[sudo] password for latif:
latif@latif-vb:~$
```

The file explorer shows a folder named "hadoop".

- **Start to configure the Java path:**

on Hadoop's virtual environment , access it using the commande :

```
sudo nano ~/hadoop/etc/hadoop/hadoop-env.sh
```

Add the following line to the file for it to access the java path

```
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64/
```



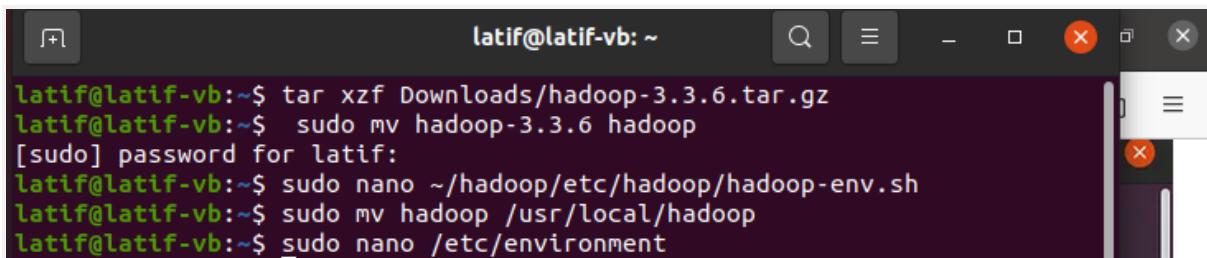
The screenshot shows a terminal window titled "latif@latif-vb: ~". It displays a file named "hadoop-env.sh" being edited with "GNU nano 4.8". The file contains several comments and one line of code:

```
## {YARN_xyz|HDFS_xyz} > HADOOP_xyz > hard-coded defaults
##
# Many of the options here are built from the perspective that users
# may want to provide OVERWRITING values on the command line.
# For example:
#
# export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64/
#
```

- **Move the Folder to an Appropriate Location**

Move the renamed folder to `/usr/local` for an easy access later on , using :

```
sudo mv hadoop /usr/local/hadoop
```



The screenshot shows a terminal window titled "latif@latif-vb: ~". It displays a series of commands run at the prompt:

```
latif@latif-vb:~$ tar xzf Downloads/hadoop-3.3.6.tar.gz
latif@latif-vb:~$ sudo mv hadoop-3.3.6 hadoop
[sudo] password for latif:
latif@latif-vb:~$ sudo nano ~hadoop/etc/hadoop/hadoop-env.sh
latif@latif-vb:~$ sudo mv hadoop /usr/local/hadoop
latif@latif-vb:~$ sudo nano /etc/environment
```

- **set up hadoop path**

on machine's environment, open the environment file with:

```
sudo nano /etc/environment
```

Add the following lines

```
'PATH="/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games:/usr/local/games:/usr/local/hadoop/bin:/usr/local/hadoop/sbin"
JAVA_HOME="/usr/lib/jvm/java-8-openjdk-amd64/jre" '
```

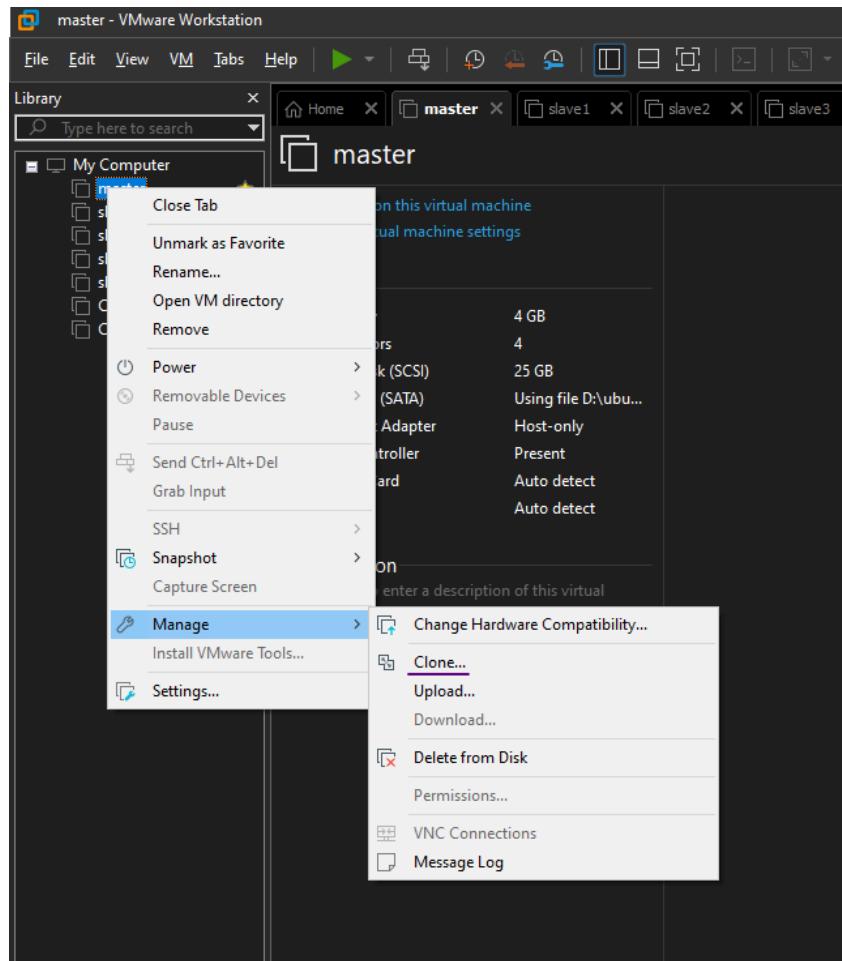


The screenshot shows a terminal window titled "latif@latif-vb: ~". It displays the contents of the "/etc/environment" file being edited with "GNU nano 4.8". The file contains the following lines:

```
PATH="/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games:/usr/local/games:/usr/local/hadoop/bin:/usr/local/hadoop/sbin"
JAVA_HOME="/usr/lib/jvm/java-8-openjdk-amd64/jre"
```

2-3 Making 'slave' machines

In the context of setting up a Hadoop cluster, configuring the "slave" machines (or more appropriately termed "worker" nodes) is an essential part of creating the distributed system. These nodes handle the actual data storage and processing in Hadoop.



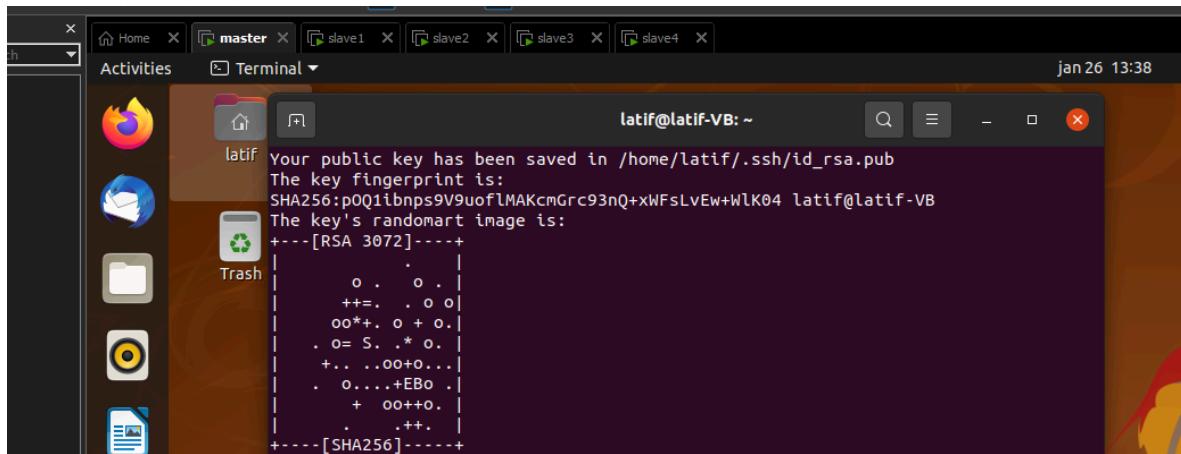
2-4 Configuring SSH for Passwordless Login

To allow the nodes in a Hadoop cluster to communicate seamlessly, we need to configure SSH for passwordless login. This ensures that the master node can execute commands on worker nodes without requiring a password

- **Generate an SSH Key Pair on the Master Node**

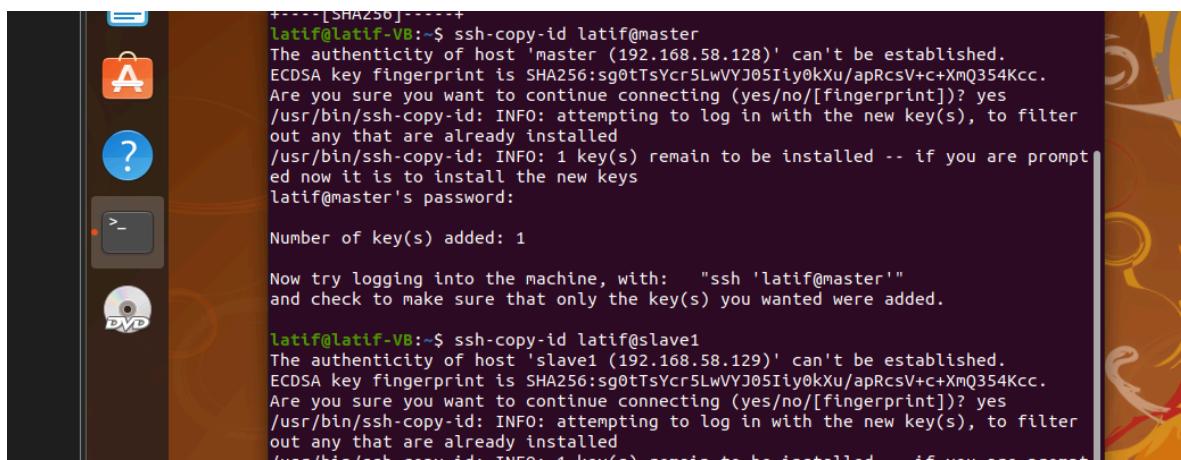
On the master node, generate an **SSH** key pair :

```
ssh-keygen -t rsa -P ""
```



- **Copy the Public Key to All Nodes**

Use the `ssh-copy-id` command to copy the master node's public key to itself and all worker nodes.



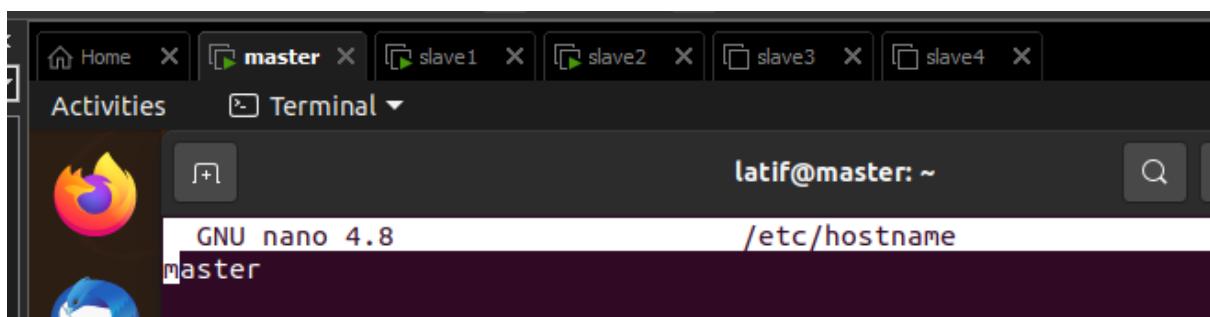
2-5 Getting IP addresses and renaming hosts

To properly set up the Hadoop cluster, obtaining the IP addresses of all nodes and renaming their hostnames is essential. This step ensures that nodes in the cluster can communicate effectively using hostnames instead of raw IP addresses. Here's how to do it:

- **Edit the Hostname File** Open the `/etc/hostname` file:

`sudo nano /etc/hostname`

Replace the existing hostname with the desired new hostname (e.g., `master`).



- **Get IP addresses** : Use the `ip a` (or `ip ad` / `ip add`)

```
latif@slave1:~$ ip a
1: lo: <LOOPBACK,UP,LOWER_UP> mtu 65536 qdisc noqueue state UNKNOWN group default qlen 1000
    link/loopback 00:00:00:00:00:00 brd 00:00:00:00:00:00
    inet 127.0.0.1/8 scope host lo
        valid_lft forever preferred_lft forever
    inet6 ::1/128 scope host
        valid_lft forever preferred_lft forever
2: ens33: <BROADCAST,MULTICAST,UP,LOWER_UP> mtu 1500 qdisc fq_codel state UNKNOWN group default qlen 1000
    link/ether 00:0c:29:6a:23:e4 brd ff:ff:ff:ff:ff:ff
    altname enp2s1
    inet 192.168.58.129/24 brd 192.168.58.255 scope global dynamic noprefixroute ens33
```

2-6 setting up the hosts files

To ensure that each node in the cluster is familiar with the other nodes, we edited the `/etc/hosts` file on all nodes (master and workers). This step is essential for name resolution within the Hadoop cluster. The following actions were performed:

- **Locate and Edit the Hosts File**

Open the `/etc/hosts` file on each node

`sudo nano /etc/hosts`

- **Add Hostnames and IP Addresses**

We added entries for all the nodes in the cluster (master and workers), mapping their IP addresses to their hostnames.

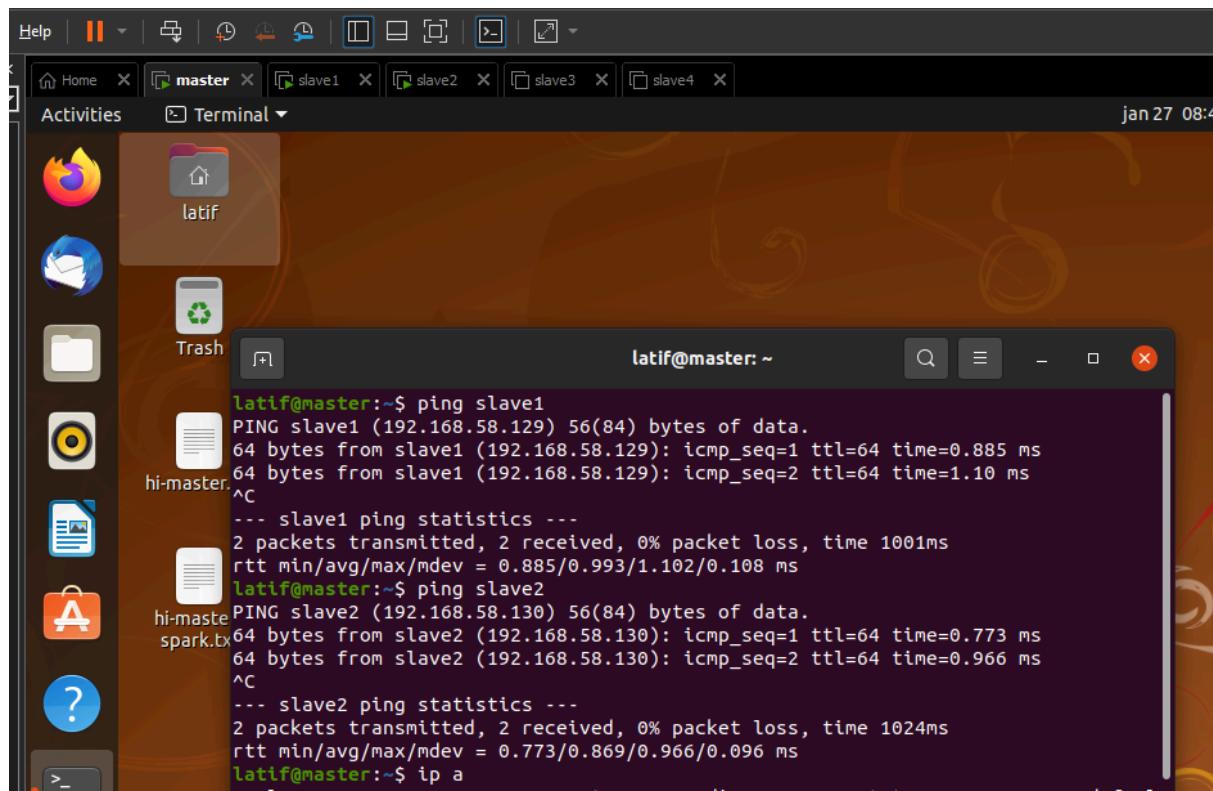
```
GNU nano 4.8          /etc/hosts          Modified
127.0.0.1      localhost
127.0.1.1      latif-VB

192.168.58.128     master
192.168.58.129     slave1
192.168.58.130     slave2
192.168.58.131     slave3
192.168.58.132     slave4

# The following lines are desirable for IPv6 capable hosts
::1      ip6-localhost ip6-loopback
fe00::0 ip6-localnet
ff00::0 ip6-mcastprefix
ff02::1 ip6-allnodes
ff02::2 ip6-allrouters
```

- **Verify the Configuration**

We tested the setup by pinging each node from the others using their hostnames

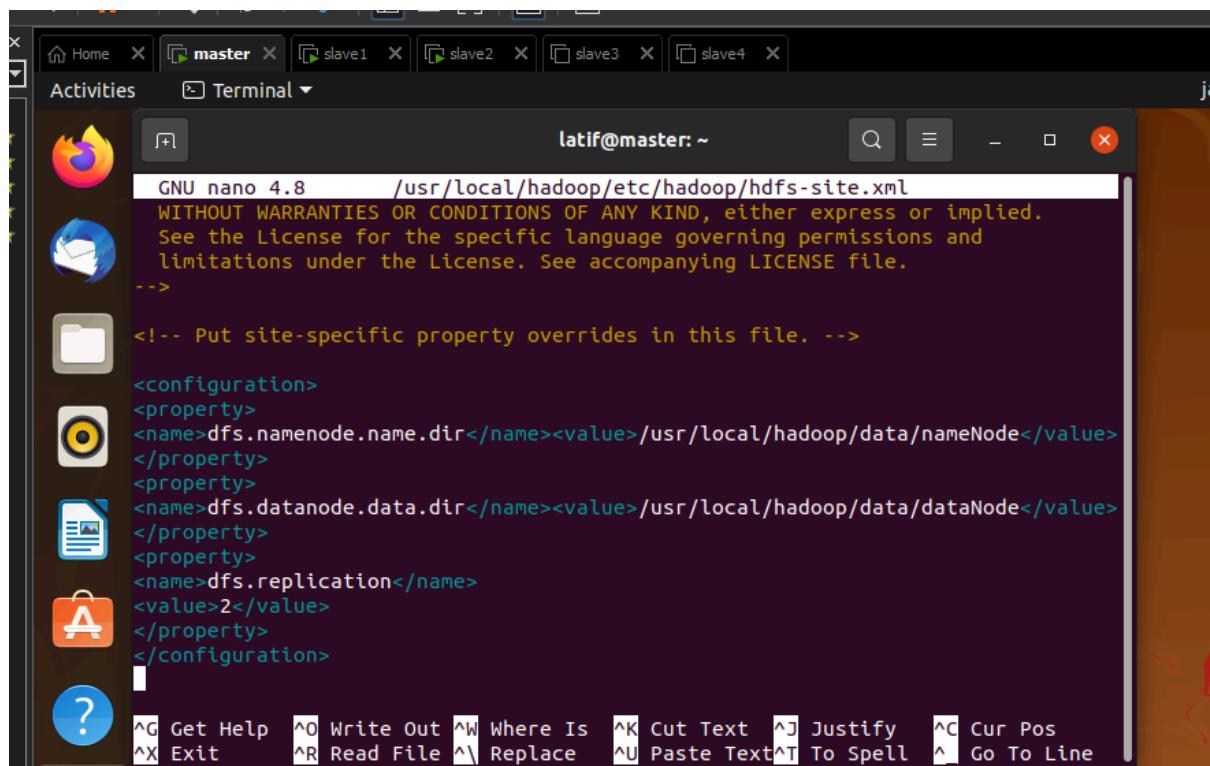


2-7 configuring hadoop on master and streaming the configuration to the slaves

- **Setting up the configuration files and identifying workers**

Inorder to set up a functioning hadoop we need to add the configuration to the
`hdfs-site.xml` , `core-site.xml` , `yarn-site.xml` , `workers`

```
sudo nano /usr/local/hadoop/etc/hadoop/hdfs-site.xml
```



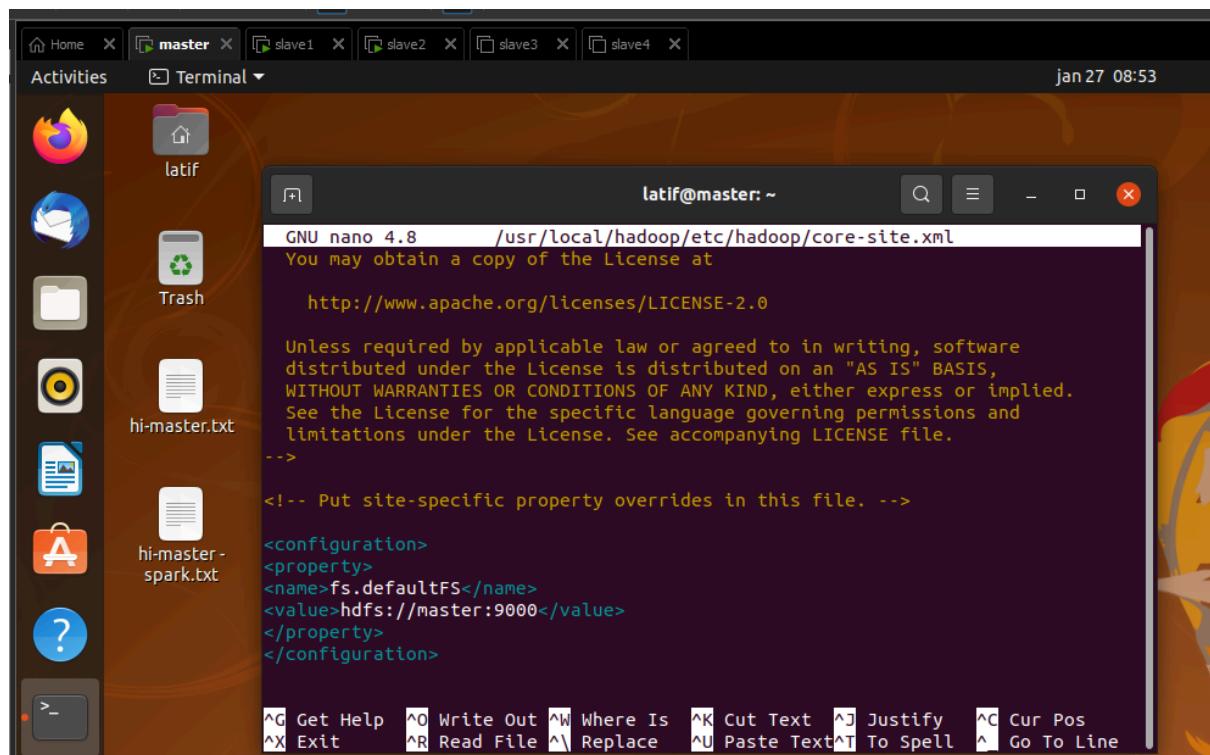
A screenshot of a Linux desktop environment, likely Ubuntu, showing a terminal window titled "latif@master: ~". The terminal is running the nano text editor on the file "/usr/local/hadoop/etc/hadoop/hdfs-site.xml". The screen displays the XML configuration for HDFS, including properties for namenode and datanode data directories, and replication settings. The terminal interface includes a menu bar, a dock with icons for Home, master, slave1, slave2, slave3, and slave4, and a bottom row of keyboard shortcuts.

```
GNU nano 4.8      /usr/local/hadoop/etc/hadoop/hdfs-site.xml
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
<name>dfs.namenode.name.dir</name><value>/usr/local/hadoop/data/nameNode</value>
</property>
<property>
<name>dfs.datanode.data.dir</name><value>/usr/local/hadoop/data/dataNode</value>
</property>
<property>
<name>dfs.replication</name>
<value>2</value>
</property>
</configuration>
```

`sudo nano /usr/local/hadoop/etc/hadoop/core-site.xml`



A screenshot of a Linux desktop environment, likely Ubuntu, showing a terminal window titled "latif@master: ~". The terminal is running the nano text editor on the file "/usr/local/hadoop/etc/hadoop/core-site.xml". The screen displays the XML configuration for Hadoop's core services, including the default file system (fs.defaultFS). The terminal interface includes a menu bar, a dock with icons for Home, master, slave1, slave2, slave3, and slave4, and a bottom row of keyboard shortcuts. The desktop background shows a cartoon character.

```
GNU nano 4.8      /usr/local/hadoop/etc/hadoop/core-site.xml
You may obtain a copy of the License at
http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
<name>fs.defaultFS</name>
<value>hdfs://master:9000</value>
</property>
</configuration>
```

`sudo nano /usr/local/hadoop/etc/hadoop/yarn-site.xml`

```
latif@latif-VB: ~
GNU nano 4.8      /usr/local/hadoop/etc/hadoop/yarn-site.xml  Modified
<?xml version="1.0"?>
<!--
Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at

http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
</configuration>
<!-- Site specific YARN configuration properties -->
<property>
<name>yarn.resourcemanager.hostname</name>
<value>master</value>
</property>
</configuration>

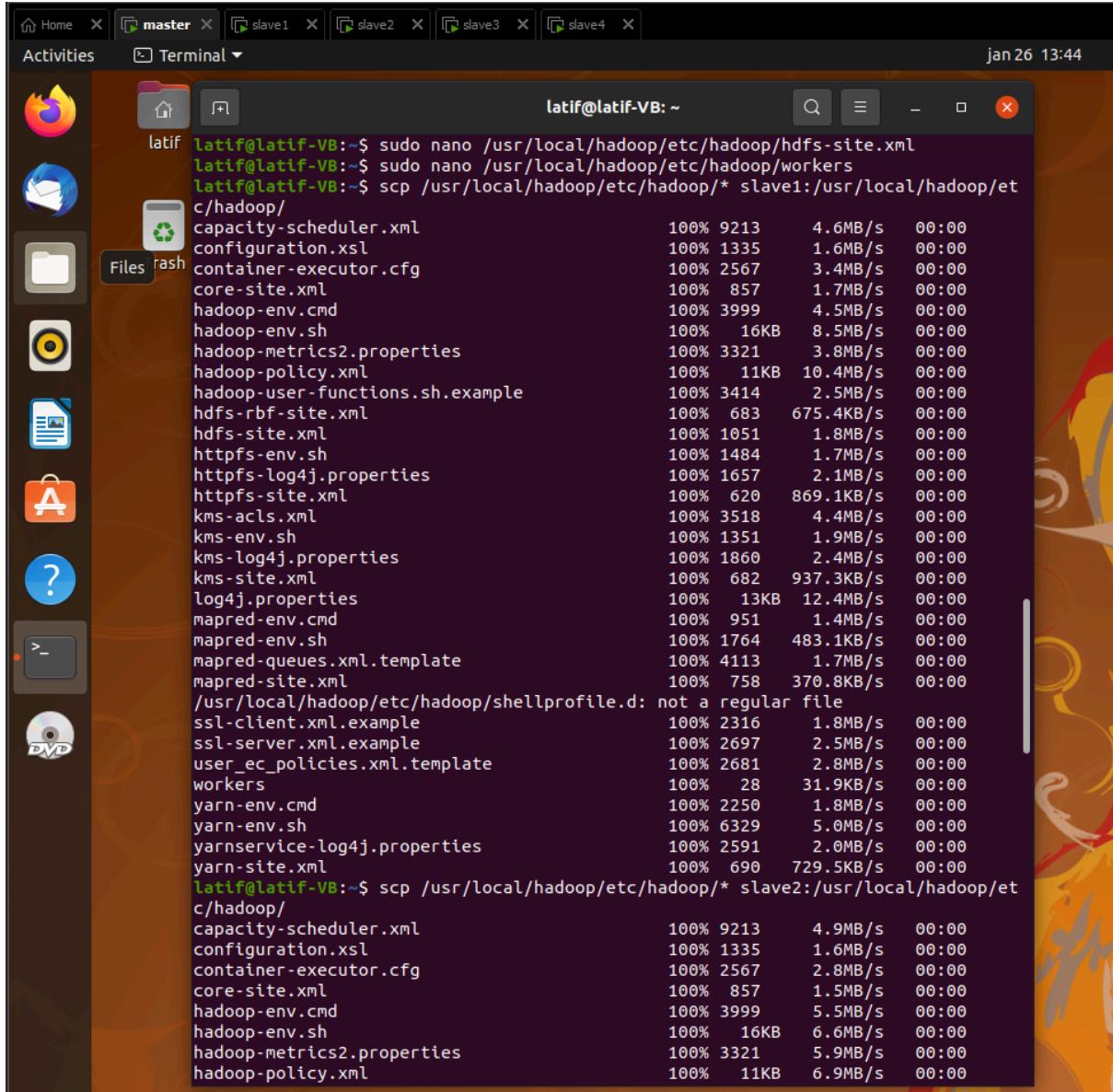
File Name to Write: /usr/local/hadoop/etc/hadoop/yarn-site.xml
^G Get Help      M-D DOS Format    M-A Append      M-B Backup File
^C Cancel       M-M Mac Format    M-P Prepend    ^T To Files
```

```
sudo nano /usr/local/hadoop/etc/hadoop/workers
```

```
latif@master: ~
GNU nano 4.8      /usr/local/hadoop/etc/hadoop/workers
slave1
slave2
slave3
slave4
```

- Streaming the updated configuration :
Using the command to update the new configuration seamlessly to all the other nodes (aka. Slaves), to ensure a correct configuration

```
scp/usr/local/hadoop/etc/hadoop/*h-secondary1:/usr/local/hadoop/etc/hadoop/
```

A screenshot of a Linux desktop environment. On the left is a dock with various icons: Home, master, slave1, slave2, slave3, slave4, Activities, Terminal, Files, Dash, Home, Applications, Help, and a DVD icon. The main area shows a terminal window titled 'latif@latif-VB: ~'. The user has run the command 'scp /usr/local/hadoop/etc/hadoop/* slave1:/usr/local/hadoop/etc/hadoop/' and 'scp /usr/local/hadoop/etc/hadoop/* slave2:/usr/local/hadoop/etc/hadoop/'. The terminal displays the progress of the file transfers, showing 100% completion for each file. The files transferred include configuration files like 'capacity-scheduler.xml', 'core-site.xml', 'hadoop-env.cmd', 'hadoop-env.sh', 'hadoop-metrics2.properties', 'hadoop-policy.xml', etc., and shell scripts like 'mapred-env.cmd', 'mapred-env.sh', 'mapred-queues.xml.template', 'mapred-site.xml', etc.

```
latif@latif-VB:~$ sudo nano /usr/local/hadoop/etc/hadoop/hdfs-site.xml
latif@latif-VB:~$ sudo nano /usr/local/hadoop/etc/hadoop/workers
latif@latif-VB:~$ scp /usr/local/hadoop/etc/hadoop/* slave1:/usr/local/hadoop/etc/hadoop/
capacity-scheduler.xml          100% 9213    4.6MB/s  00:00
configuration.xsl               100% 1335    1.6MB/s  00:00
container-executor.cfg          100% 2567    3.4MB/s  00:00
core-site.xml                   100% 857     1.7MB/s  00:00
hadoop-env.cmd                 100% 3999    4.5MB/s  00:00
hadoop-env.sh                  100% 16KB    8.5MB/s  00:00
hadoop-metrics2.properties     100% 3321    3.8MB/s  00:00
hadoop-policy.xml              100% 11KB    10.4MB/s 00:00
hadoop-user-functions.sh.example 100% 3414    2.5MB/s  00:00
hdfs-rbf-site.xml              100% 683     675.4KB/s 00:00
hdfs-site.xml                  100% 1051    1.8MB/s  00:00
httpfs-env.sh                  100% 1484    1.7MB/s  00:00
httpfs-log4j.properties        100% 1657    2.1MB/s  00:00
httpfs-site.xml                100% 620     869.1KB/s 00:00
kms-acls.xml                  100% 3518    4.4MB/s  00:00
kms-env.sh                     100% 1351    1.9MB/s  00:00
kms-log4j.properties           100% 1860    2.4MB/s  00:00
kms-site.xml                   100% 682     937.3KB/s 00:00
log4j.properties               100% 13KB    12.4MB/s 00:00
mapred-env.cmd                 100% 951     1.4MB/s  00:00
mapred-env.sh                  100% 1764    483.1KB/s 00:00
mapred-queues.xml.template    100% 4113    1.7MB/s  00:00
mapred-site.xml                100% 758     370.8KB/s 00:00
/usr/local/hadoop/etc/hadoop/shellprofile.d: not a regular file
ssl-client.xml.example         100% 2316    1.8MB/s  00:00
ssl-server.xml.example         100% 2697    2.5MB/s  00:00
user_ec_policies.xml.template 100% 2681    2.8MB/s  00:00
workers                         100% 28      31.9KB/s 00:00
yarn-env.cmd                   100% 2250    1.8MB/s  00:00
yarn-env.sh                     100% 6329    5.0MB/s  00:00
yarnservice-log4j.properties   100% 2591    2.0MB/s  00:00
yarn-site.xml                   100% 690     729.5KB/s 00:00
latif@latif-VB:~$ scp /usr/local/hadoop/etc/hadoop/* slave2:/usr/local/hadoop/etc/hadoop/
capacity-scheduler.xml          100% 9213    4.9MB/s  00:00
configuration.xsl               100% 1335    1.6MB/s  00:00
container-executor.cfg          100% 2567    2.8MB/s  00:00
core-site.xml                   100% 857     1.5MB/s  00:00
hadoop-env.cmd                 100% 3999    5.5MB/s  00:00
hadoop-env.sh                  100% 16KB    6.6MB/s  00:00
hadoop-metrics2.properties     100% 3321    5.9MB/s  00:00
hadoop-policy.xml              100% 11KB    6.9MB/s  00:00
```

2-8 launching hadoop

Start Hadoop Services

Ensure the Hadoop services are running by executing the following commands on the master node:

- Start the Hadoop Distributed File System (HDFS):
`start-dfs.sh`
- Start the YARN resource manager:
`start-yarn.sh`
- Confirm that the services are running by checking the logs or using:
`jps`

```

latif@master:~$ start-dfs.sh
Starting namenodes on [master]
Starting datanodes
Starting secondary namenodes [master]
latif@master:~$ start-yarn.sh
Starting resourcemanager
Starting nodemanagers
latif@master:~$ jps
2803 Jps
2506 ResourceManager
2075 NameNode
2319 SecondaryNameNode

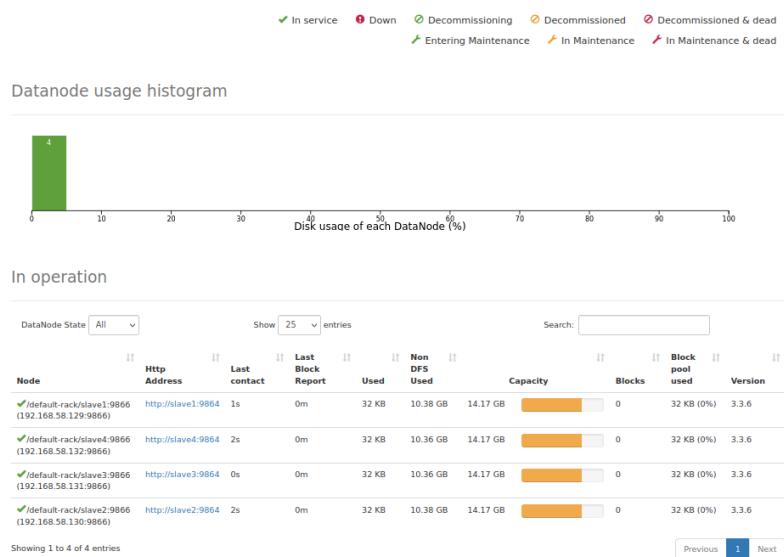
```

Namenode information Problem loading page +

localhost:9870/dfshealth.html#tab-datanode

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

Datanode Information



Entering Maintenance

Nodes of the cluster Namenode information +

192.168.58.128:8088/cluster/nodes

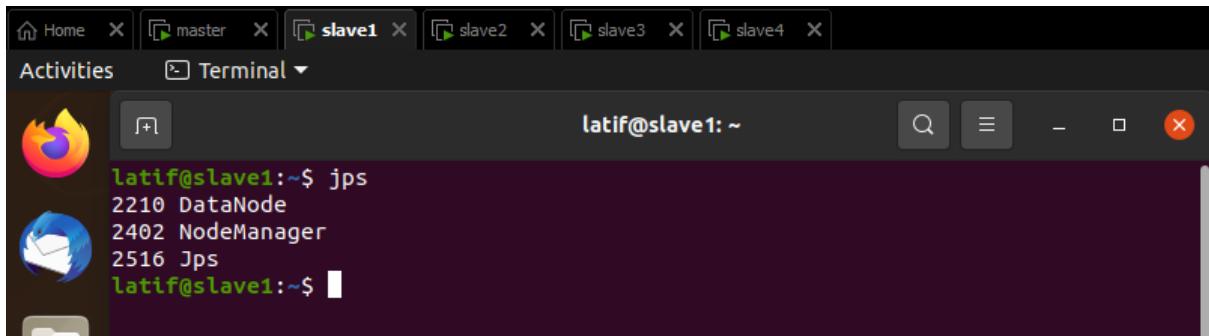
Nodes of the cluster

Node Labels	Rack	Node State	Node Address	Node HTTP Address	Last health-update	Health-report	Containers	Allocation Tags	Mem Used	Mem Avail	Phys Mem Used %	Vcores Used %	Vcores Avail	Phys Vcores Used %	Version
/default-rack	RUNNING	slave1:38277	slave1:8042		Mon Jan 27 12:44:41 +0100 2025	0	0 B	8 GB	44	0	8	0	3.3.6		
/default-rack	RUNNING	slave2:44657	slave2:8042		Mon Jan 27 12:44:41 +0100 2025	0	0 B	8 GB	85	0	8	0	3.3.6		
/default-rack	RUNNING	slave3:39093	slave3:8042		Mon Jan 27 12:44:42 +0100 2025	0	0 B	8 GB	86	0	8	0	3.3.6		
/default-rack	RUNNING	slave4:43945	slave4:8042		Mon Jan 27 12:44:41 +0100 2025	0	0 B	8 GB	36	0	8	0	3.3.6		

Showing 1 to 4 of 4 entries

- Checking hadoop and yarn on slaves

Using jps



```
latif@slave1:~$ jps
2210 DataNode
2402 NodeManager
2516 Jps
latif@slave1:~$
```

3- Spark installation and configuration

3-1 Installing JAVA

Spark requires Java to operate, and we already have the JDK installed.

3-2 Download Apache Spark

Spark requires Java to operate, and we already have the JDK installed.

1. Head over to the [Apache Spark Download page](#).
2. Choose a Spark release: 3.4.4 (Oct 27 2024)
3. Choose a package type: Pre-built for Apache Hadoop 3.3 and later / Pre-built for Apache Hadoop 3.3 and later (Scala 2.13)
4. Download Spark: [spark-3.4.4-bin-hadoop3.tgz](#)

3-3 Installing Spark

- **Extract the Downloaded File**

After downloading the `.tgz` file, the next step is to extract it using the `tar` command

```
tar -xvzf Downloads/spark-3.4.4-bin-hadoop3.tgz
```

- **Rename the spark folder**

To make it easier for further usage

```
sudo mv spark-3.4.4-bin-hadoop3 spark
```

- **Move the Extracted Folder**

To make Spark accessible system-wide, move the extracted folder to the `directory /usr/local`.

```
sudo mv spark /usr/local/spark
```

```

spark-3.4.4-bin-hadoop3/R/lib/SparkR/tests/testthat/
Tspark-3.4.4-bin-hadoop3/R/lib/SparkR/tests/testthat/test_basic.R
spark-3.4.4-bin-hadoop3/R/lib/SparkR/worker/
spark-3.4.4-bin-hadoop3/R/lib/SparkR/worker/daemon.R
spark-3.4.4-bin-hadoop3/R/lib/SparkR/worker/worker.R
spark-3.4.4-bin-hadoop3/R/lib/SparkR/help/
spark-3.4.4-bin-hadoop3/R/lib/SparkR/help/paths.rds
spark-3.4.4-bin-hadoop3/R/lib/SparkR/help/SparkR.rdb
spark-3.4.4-bin-hadoop3/R/lib/SparkR/help/SparkR.rdx
spark-3.4.4-bin-hadoop3/R/lib/SparkR/help/AnIndex
spark-3.4.4-bin-hadoop3/R/lib/SparkR/help/aliases.rds
spark-3.4.4-bin-hadoop3/R/lib/SparkR/html/
spark-3.4.4-bin-hadoop3/R/lib/SparkR/html/00Index.html
spark-3.4.4-bin-hadoop3/R/lib/SparkR/html/R.css
spark-3.4.4-bin-hadoop3/R/lib/SparkR/INDEX
spark-3.4.4-bin-hadoop3/R/lib/sparkr.zip
latif@master:~$ mv spark-3.4.4-bin-hadoop3 /usr/local/spark
mv: cannot move 'spark-3.4.4-bin-hadoop3' to '/usr/local/spark': Permission denied
latif@master:~$ mv spark /usr/local/spark
mv: cannot move 'spark' to '/usr/local/spark': Permission denied
latif@master:~$ sudo mv spark /usr/local/spark
[sudo] password for latif:
latif@master:~$ 

```

3-4 Configure Environment Variables

To ensure that Spark commands are recognized system-wide, update the `.bashrc` file to include Spark's `bin` directory in the system's PATH.

1. Open the `.bashrc` file for editing:

```
sudo nano ~/.bashrc
```

2. Add the following lines

```
export SPARK_HOME=/usr/local/spark
```

```
export PATH=$SPARK_HOME/bin:$PATH
```

```

export JAVA_HOME /usr/lib/jvm/java-8-openjdk-amd64
export PATH=$PATH:/usr/lib/jvm/java-8-openjdk-amd64/bin
export HADOOP_HOME=/usr/local/hadoop
export PATH=$PATH:$HADOOP_HOME/bin
export PATH=$PATH:$HADOOP_HOME/sbin
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export YARN_HOME=$HADOOP_HOME
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native
export HADOOP_OPTS="-Djava.library.path=$HADOOP_HOME/lib/native"
export HADOOP_STREAMING=$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-3.2.1.jar
export HADOOP_LOG_DIR=$HADOOP_HOME/logs
export PSSH_RCMD_TYPE=ssh

export SPARK_HOME=/usr/local/spark
export PATH=$SPARK_HOME/bin:$SPARK_HOME/sbin

```

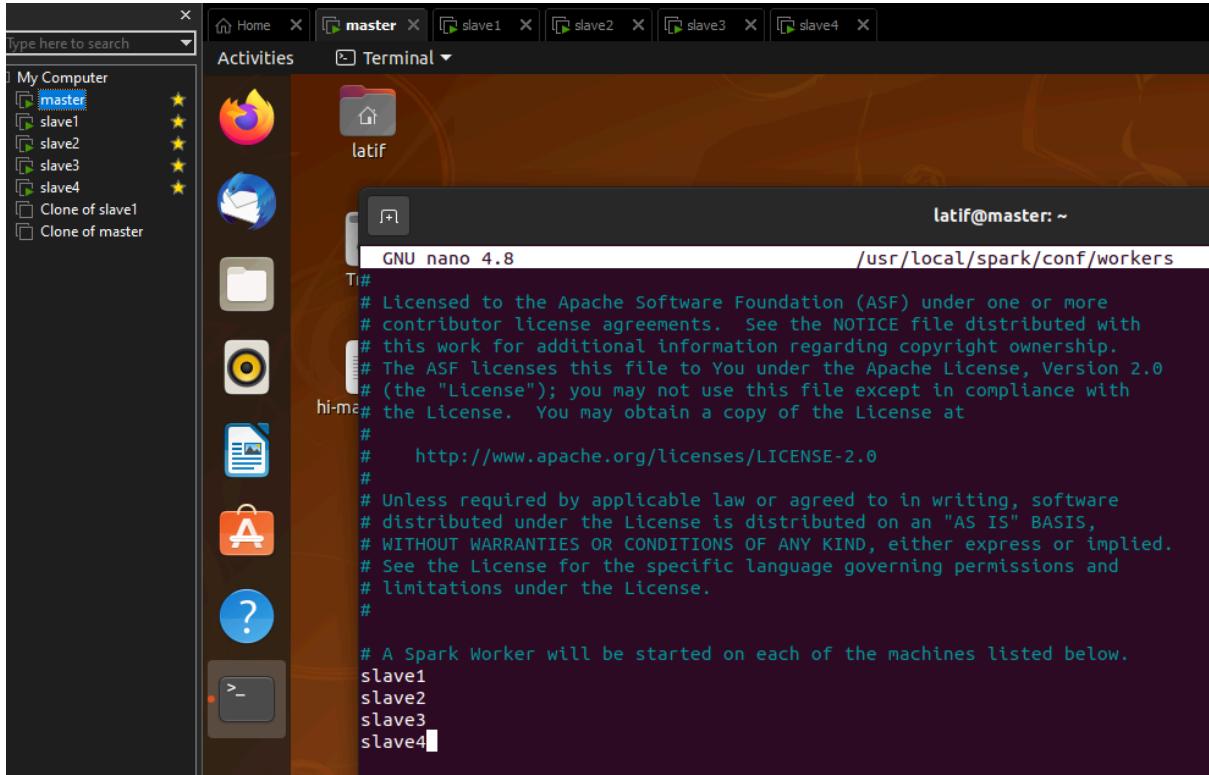
3-5 Update Workers file

When setting up a Spark cluster, we need to specify the slave nodes (or worker nodes) in the `slaves` file.

- Open the `workers` file

```
nano $SPARK_HOME/conf/workers
```

- Add the hostnames



3-6 update workers with the new config

- **spark-env.sh File:** Open the file for editing:

```
sudo nano spark-env.sh
```

- **Set the Required Environment Variables:** Add the following lines to configure Spark::

Example configuration:

```

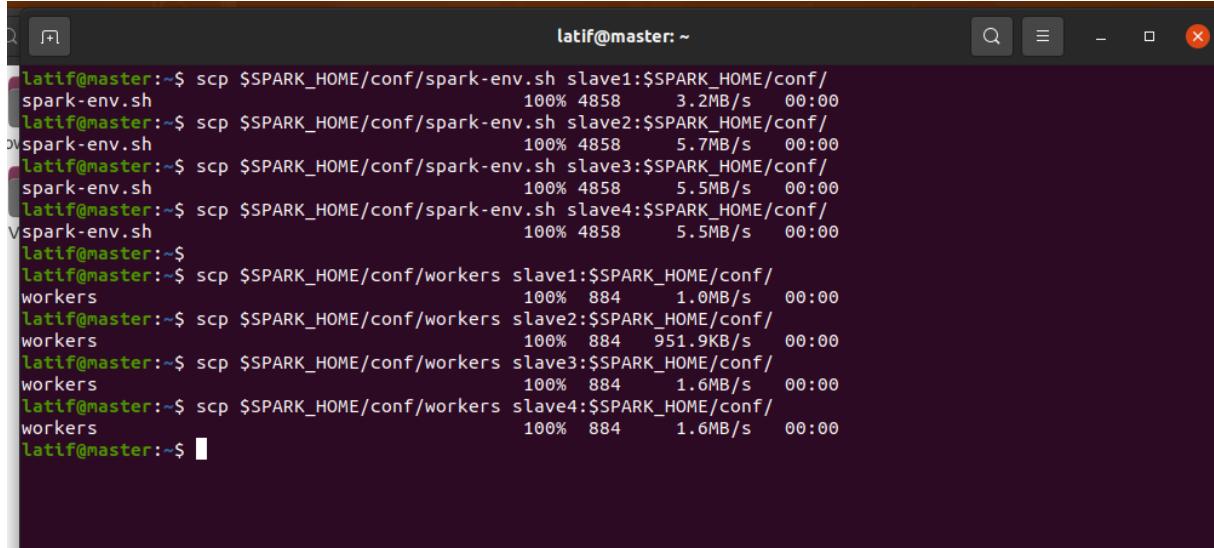
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export SPARK_MASTER_HOST=192.168.58.128
export SPARK_WORKER_CORES=4
export SPARK_WORKER_MEMORY=8G
export SPARK_WORKER_INSTANCES=1
export SPARK_LOCAL_DIRS=/usr/local/spark/tmp

```

Distribute the spark-env.sh File to All Nodes: Copy the updated `spark-env.sh` file to the `conf` directory on all nodes (workers) in the cluster.

example:

```
scp $SPARK_HOME/conf/spark-env.sh slave2:/usr/local/spark/conf/
```



```
latif@master:~$ scp $SPARK_HOME/conf/spark-env.sh slave1:$SPARK_HOME/conf/
spark-env.sh                                100% 4858      3.2MB/s  00:00
latif@master:~$ scp $SPARK_HOME/conf/spark-env.sh slave2:$SPARK_HOME/conf/
spark-env.sh                                100% 4858      5.7MB/s  00:00
latif@master:~$ scp $SPARK_HOME/conf/spark-env.sh slave3:$SPARK_HOME/conf/
spark-env.sh                                100% 4858      5.5MB/s  00:00
latif@master:~$ scp $SPARK_HOME/conf/spark-env.sh slave4:$SPARK_HOME/conf/
spark-env.sh                                100% 4858      5.5MB/s  00:00
latif@master:~$ 
latif@master:~$ scp $SPARK_HOME/conf/workers slave1:$SPARK_HOME/conf/
workers                                     100%  884      1.0MB/s  00:00
latif@master:~$ scp $SPARK_HOME/conf/workers slave2:$SPARK_HOME/conf/
workers                                     100%  884     951.9KB/s  00:00
latif@master:~$ scp $SPARK_HOME/conf/workers slave3:$SPARK_HOME/conf/
workers                                     100%  884      1.6MB/s  00:00
latif@master:~$ scp $SPARK_HOME/conf/workers slave4:$SPARK_HOME/conf/
workers                                     100%  884      1.6MB/s  00:00
latif@master:~$
```

3-7 launching Spark

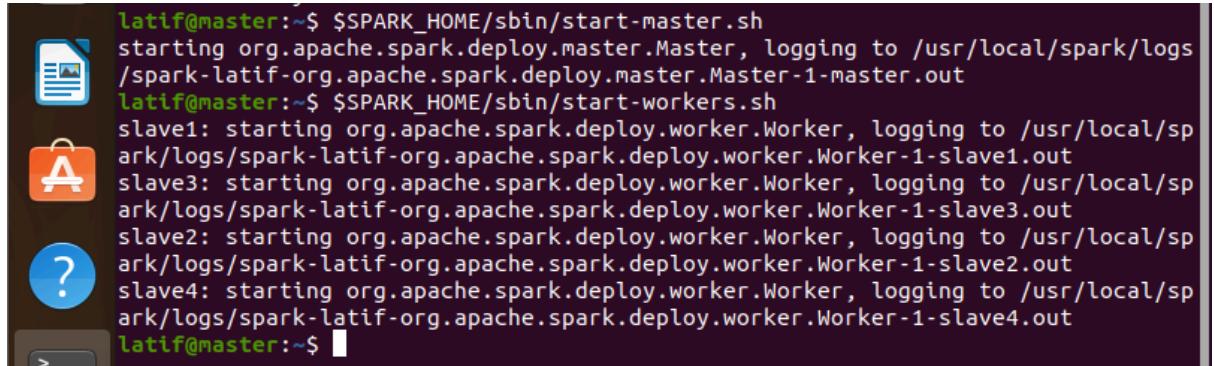
From the master node we run both commands

- **To start the master node**

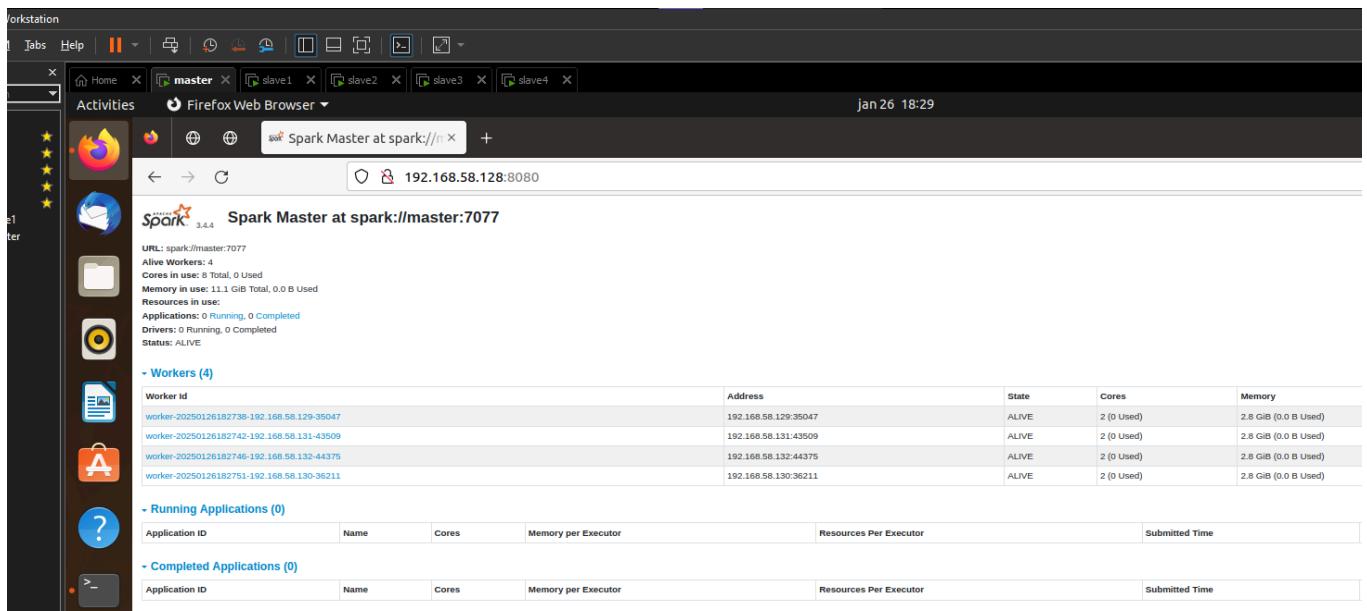
```
$SPARK_HOME/sbin/start-master.sh
```

- **To start slaves (workers) nodes**

```
$SPARK_HOME/sbin/start-workers.sh
```



```
latif@master:~$ $SPARK_HOME/sbin/start-master.sh
starting org.apache.spark.deploy.master.Master, logging to /usr/local/spark/logs
/spark-latif-org.apache.spark.deploy.master.Master-1-master.out
latif@master:~$ $SPARK_HOME/sbin/start-workers.sh
slave1: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/sp
ark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave1.out
slave3: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/sp
ark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave3.out
slave2: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/sp
ark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave2.out
slave4: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/sp
ark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave4.out
latif@master:~$
```



Code and running

I. First approach : Hadoop & spark

Code

A- Installing and importing the necessary libraries

- Install Python Dependencies

```
sudo apt install -y python3 python3-pip
```

- Install PySpark and dependencies:

```
pip3 install pyspark pillow
```

- Enable JPEG2000 Support in PIL (Pillow)

By default, Pillow **does not support JPEG2000**. You need to install **OpenJPEG**:

```
sudo apt install -y libopenjp2-7-dev
```

B- Implementing the solution

1. Initialize SparkSession for connecting to Hadoop cluster

- Set application name ("JPEG2000 Compression")
- Set master to "yarn" for cluster management
- Create Spark context (sc)

2. Define base path on HDFS for storing compressed images

```
hdfs_output_base = "hdfs://master:9000/user/hadoop/compressed_images"
```

3. Get Hadoop FileSystem handle from Spark for interacting with HDFS

```
conf = sc._jsc.hadoopConfiguration()
fs = sc._jvm.org.apache.hadoop.fs.FileSystem.get(conf)
Path = spark._jvm.org.apache.hadoop.fs.Path # Java Path class
```

4. Define a function `save_to_hdfs(filename, data)`:

- Create full path for the image on HDFS
- Open an output stream to the target HDFS location
- Write the binary data (compressed image) to HDFS
- Close the stream

5. Define a function `compress_image(image_bytes)`:

- Try to open image from binary data
img = Image.open(io.BytesIO(image_bytes))
- Compress the image using JPEG2000 format
output_buffer = io.BytesIO()
img.save(output_buffer, format="JPEG2000")
- Return compressed image as bytes
return output_buffer.getvalue()
- Handle errors if any

6. Define the input path for images stored in HDFS

```
input_path = "hdfs://master:9000/user/hadoop/input_images"
```

7. Load images from HDFS as binary files using Spark's `binaryFiles` function

```
images = sc.binaryFiles(input_path)
```

8. Map over the loaded images:

- Compress each image using `compress_image`
compressed_images = images.mapValues(lambda data: compress_image(data))
- Store compressed data paired with the image filename

9. Collect the results back to the driver (assuming dataset is small enough)

```
results = compressed_images.collect()
```

10. Loop through the collected results:

- For each image, extract the filename
- Save compressed data to HDFS using `save_to_hdfs`

11. Stop the SparkSession to release resources

```
spark.stop()
```

Starting hadoop and spark

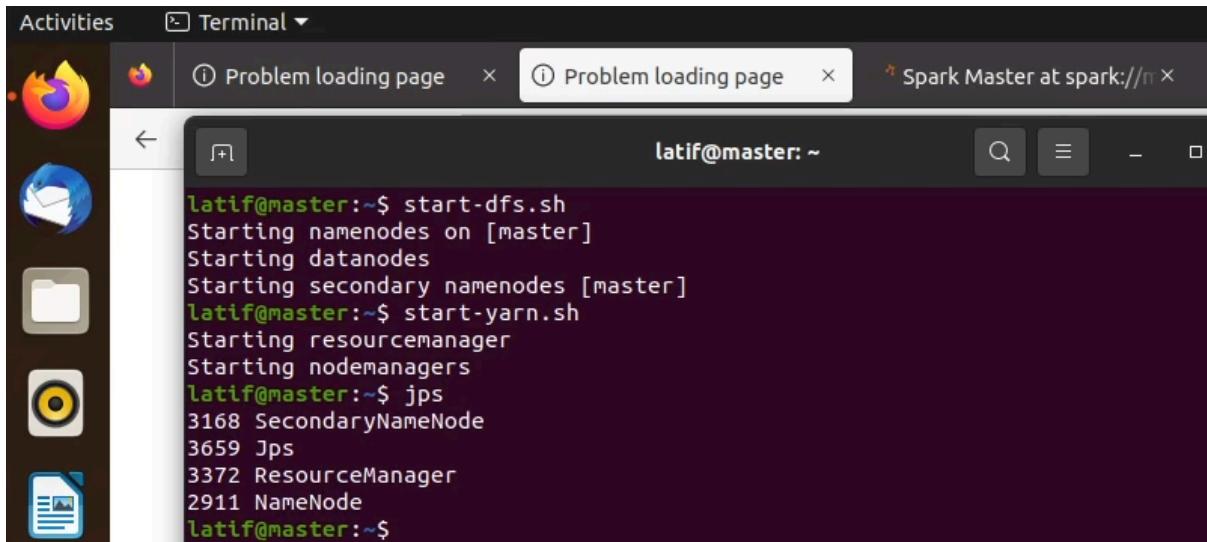
From the master node we start the HDFS and the YARN;

- **Start the Hadoop Distributed File System (HDFS):**

```
start-dfs.sh
```

- **Start the YARN resource manager:**

```
start-yarn.sh
```



```
Activities Terminal ▾
latif@master:~$ start-dfs.sh
Starting namenodes on [master]
Starting datanodes
Starting secondary namenodes [master]
latif@master:~$ start-yarn.sh
Starting resourcemanager
Starting nodemanagers
latif@master:~$ jps
3168 SecondaryNameNode
3659 Jps
3372 ResourceManager
2911 NameNode
latif@master:~$
```

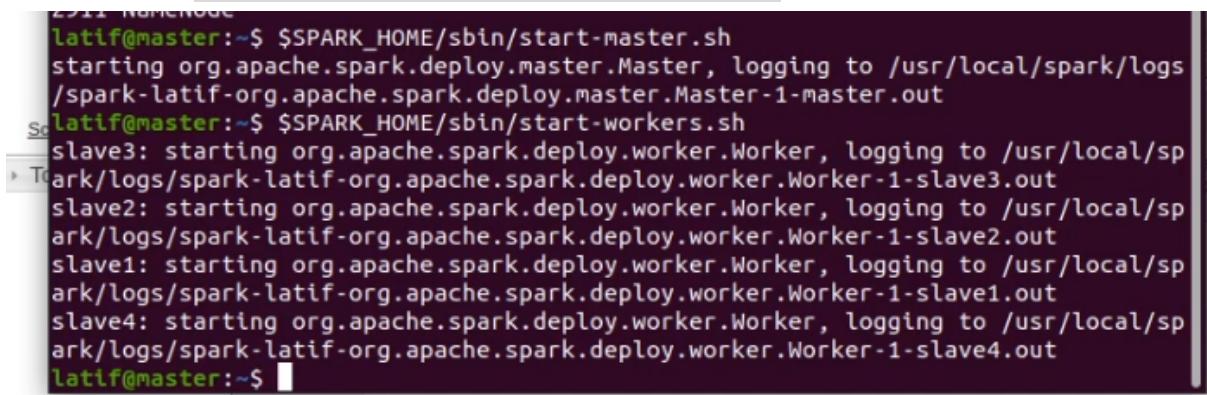
From the master node we run both commands

- **To start the master node**

```
$SPARK_HOME/sbin/start-master.sh
```

- **To start slaves (workers) nodes**

```
$SPARK_HOME/sbin/start-workers.sh
```



```
latif@master:~$ $SPARK_HOME/sbin/start-master.sh
starting org.apache.spark.deploy.master.Master, logging to /usr/local/spark/logs/spark-latif-org.apache.spark.deploy.master.Master-1-master.out
latif@master:~$ $SPARK_HOME/sbin/start-workers.sh
slave3: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave3.out
slave2: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave2.out
slave1: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave1.out
slave4: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-latif-org.apache.spark.deploy.worker.Worker-1-slave4.out
latif@master:~$
```

Check the HDFS

Checking if the *input_images* and *compresed_images* folders exist in the HDFS , if not to be created using

```
hdfs dfs -mkdir -p /user/hadoop/compressed_images
```

```
hdfs dfs -mkdir -p /user/hadoop/input_images
```

```
latif@master:~/user/hadoop/input_images$ hdfs dfs -ls /user/hadoop/
Found 3 items
drwxr-xr-x 2 latif supergroup 0 2025-02-01 21:11 /user/hadoop/compre
sed_images
drwxr-xr-x 2 latif supergroup 0 2025-02-01 20:41 /user/hadoop/input_i
mages
drwxr-xr-x 2 latif supergroup 0 2025-02-01 17:34 /user/hadoop/jars
Status:latif@master:~$
```

Upload the images to the designated folder in HDFS

Using the command below allows the upload of the files from the local machine to the HDFS system

```
hdfs dfs -put ~/Desktop/Hadoop-job/images/*
/user/hadoop/input_images
```

```
latif@master:~/Desktop/mapreduce-job$ hdfs dfs -put ~/Desktop/Hadoop-job/images/* /user/hadoop/input_images
latif@master:~/Desktop/mapreduce-job$ hdfs dfs -ls /user/hadoop/input_images
Found 3 items
-rw-r--r-- 2 latif supergroup 2108235 2025-02-01 23:18 /user/hadoop/input_images/20241231_111322.jpg
-rw-r--r-- 2 latif supergroup 2464886 2025-02-01 23:18 /user/hadoop/input_images/20250107_093226.jpg
-rw-r--r-- 2 latif supergroup 240045 2025-02-01 23:18 /user/hadoop/input_images/IMG_20241231_202142_665.jpg
```

Job submission

`spark-submit` command, it communicates with the cluster manager (YARN) to allocate resources on available nodes.

```
spark-submit --master yarn --deploy-mode cluster
Desktop/Hadoop-job/jpeg2000_compression.py
```

```
latif@master:~$ spark-submit --master yarn --deploy-mode cluster Desktop/Hadoop-
job/jpeg2000_compression.py
```

Running and waiting for job compilation

When a Spark job is submitted, it is compiled and prepared for execution across the cluster's nodes. Once the compilation is complete, the job is scheduled by the cluster manager, and tasks are distributed to available worker nodes for parallel processing.

The screenshot shows the Hadoop Cluster Application Overview page. The application details are as follows:

- User:** latif
- Name:** jpeg2000_compression.py
- Application Type:** SPARK
- Application Priority:** 0 (Higher Integer value indicates higher priority)
- YarnApplicationState:** FINISHED
- Queue:** default
- FinalStatus Reported by AM:** SUCCEEDED
- Started:** Sat Feb 01 21:24:34 +0100 2025
- Length:** Sat Feb 01 21:24:36 +0100 2025
- Finished:** Sat Feb 01 21:25:27 +0100 2025
- Elapsed:** 52secs
- Tracking URL:** History
- Log Aggregation Status:** DISABLED
- Application Timeout (Remaining Time):** Unlimited
- Diagnostics:**
- Unmanaged Application:** false
- Application Node Label expression:** <Not set>
- AM container Node Label expression:** <DEFAULT_PARTITION>

Application Metrics section shows:

- Total Resource Preempted: <memory:0, vCores:0>
- Total Number of Non-AM Containers Preempted: 0
- Total Number of AM Containers Preempted: 0
- Resource Preempted from Current Attempt: <memory:0, vCores:0>
- Number of Non-AM Containers Preempted from Current Attempt: 0
- Aggregate Resource Allocation: 246484 MB-seconds, 119 vcore-seconds
- Aggregate Preempted Resource Allocation: 0 MB-seconds, 0 vcore-seconds

Table showing Application Metrics:

Attempt ID	Started	Node	Logs	Nodes blacklisted by the app	Nodes blacklisted by the system
appattempt_1738441000373_0001_000001	Sat Feb 1 21:24:34 +0100 2025	http://slave2:8042	Logs	0	0

Showing 1 to 1 of 1 entries

Retrieving results

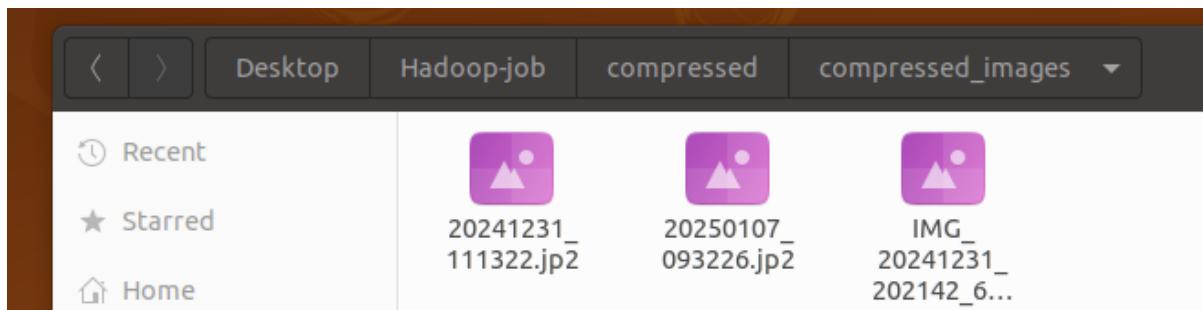
We can check if the results of the job (in our case the compression job) gave a result or no by checking the designated folder for saving the results on the HDFS ,by using the command

```
hdfs dfs -ls /user/hadoop/compressed_images
```

```
latif@master:~$ hdfs dfs -ls /user/hadoop/compressed_images
Found 3 items
-rw-r--r-- 2 latif supergroup 10040834 2025-02-01 21:11 /user/hadoop/compressed_images/20241231_111322.jp2
-rw-r--r-- 2 latif supergroup 11141779 2025-02-01 21:11 /user/hadoop/compressed_images/20250107_093226.jp2
-rw-r--r-- 2 latif supergroup 2188856 2025-02-01 21:11 /user/hadoop/compressed_images/IMG_20241231_202142_665.jp2
```

The results can be saved on one of the local machines just by running

```
hdfs dfs -get /user/hadoop/compressed_images
~/Desktop/Hadoop-job/compressed
```



II. Second approach : Hadoop (Map-Reduce)

Code

Three parts of code were implemented supporting **JPEG** compression to assure the functionality of the map reduce job



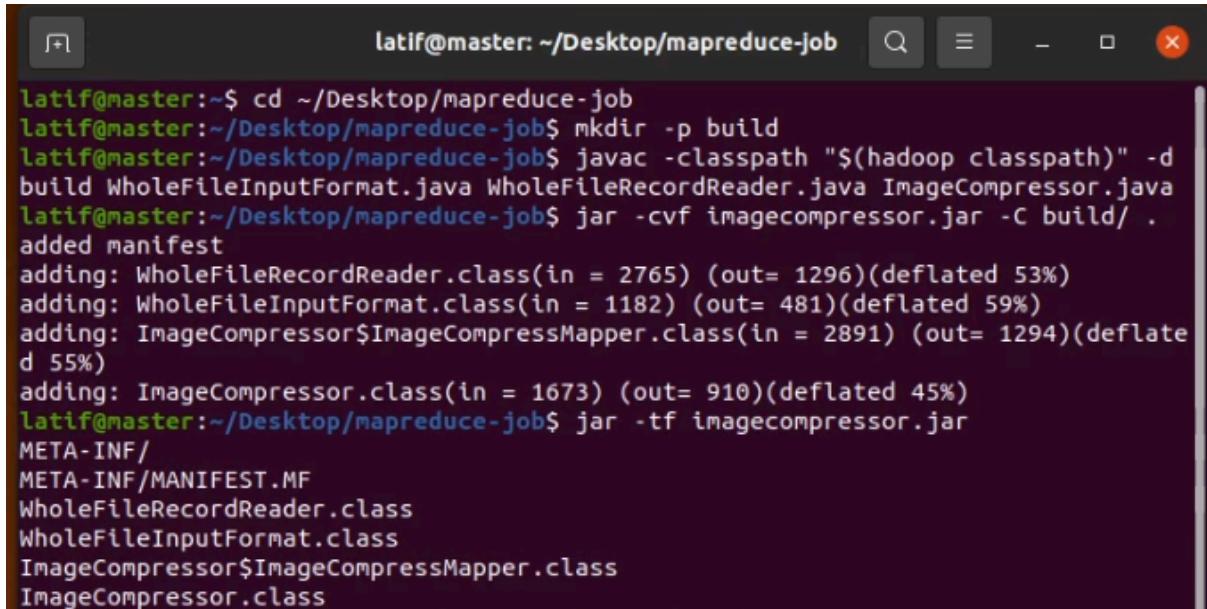
code	WholeFileInputFormat.java	WholeFileRecordReader.java	ImageCompressor.java
job	This custom input format ensures that each file is treated as a single input record instead of being split into smaller chunks	This custom record reader reads an entire file into memory and stores its contents as a BytesWritable object.	<p>This file contains:</p> <ol style="list-style-type: none"> 1. <code>ImageCompressMapper</code> (Mapper class) → Responsible for compressing images. 2. Main method → Configures and runs the MapReduce job.
Key points	<ul style="list-style-type: none"> • Extends <code>FileInputFormat<Text, BytesWritable></code> to specify the input format. • Overrides <code>isSplitable()</code> → returns false, meaning that each file is processed as a whole. • Overrides <code>createRecordReader()</code> → provides a <code>WholeFileRecordReader</code> to read the file as a whole. 	<ul style="list-style-type: none"> • Reads the entire file into a byte array. • Stores the file path as the key (<code>Text</code>) and file content as the value (<code>BytesWritable</code>). • Ensures that each file is processed only once (processed flag). 	<p><code>ImageCompressMapper</code></p> <ul style="list-style-type: none"> • Reads the image from <code>BytesWritable</code>. • Uses <code>ImageIO</code> to convert it to a <code>BufferedImage</code>. • Compresses the image to JPEG format with 50% quality. • Writes the compressed image data as output. <p><code>Main Method</code></p> <ul style="list-style-type: none"> • Sets up a Hadoop MapReduce job. • Uses <code>WholeFileInputFormat</code> to read full image files. • Runs only the Mapper function (map-only job) with <code>setNumReduceTasks(0)</code>. • Takes input and output paths as arguments
Importance	<ul style="list-style-type: none"> • Ensures that image files are processed completely rather than split, which is crucial for image processing. 	<ul style="list-style-type: none"> • Essential for handling binary files like images, ensuring the entire file is available for processing. 	<ul style="list-style-type: none"> • Image compression reduces storage size while maintaining reasonable quality. • Efficient processing with a map-only job, ideal for large-scale distributed image compression

Compiling the java files

To prepare the previous files to be used in the job , the task of compiling them is assured by javac

```
javac -classpath $(hadoop classpath) -d build/ WholeFileInputFormat.java  
WholeFileRecordReader.java ImageCompressor.java
```

- After compilation , we need to build the jar that's going to be introduced later on



A terminal window titled "latif@master: ~/Desktop/mapreduce-job". The window shows the command line history and output of the following commands:

```
latif@master:~$ cd ~/Desktop/mapreduce-job  
latif@master:~/Desktop/mapreduce-job$ mkdir -p build  
latif@master:~/Desktop/mapreduce-job$ javac -classpath "$(hadoop classpath)" -d build WholeFileInputFormat.java WholeFileRecordReader.java ImageCompressor.java  
latif@master:~/Desktop/mapreduce-job$ jar -cvf imagecompressor.jar -C build/ .  
added manifest  
adding: WholeFileRecordReader.class(in = 2765) (out= 1296)(deflated 53%)  
adding: WholeFileInputFormat.class(in = 1182) (out= 481)(deflated 59%)  
adding: ImageCompressor$ImageCompressMapper.class(in = 2891) (out= 1294)(deflate d 55%)  
adding: ImageCompressor.class(in = 1673) (out= 910)(deflated 45%)  
latif@master:~/Desktop/mapreduce-job$ jar -tf imagecompressor.jar  
META-INF/  
META-INF/MANIFEST.MF  
WholeFileRecordReader.class  
WholeFileInputFormat.class  
ImageCompressor$ImageCompressMapper.class  
ImageCompressor.class
```

Starting hadoop

From the master node we start the HDFS and the YARN;

- **Start the Hadoop Distributed File System (HDFS):**
`start-dfs.sh`
- **Start the YARN resource manager:**
`start-yarn.sh`

Check the HDFS

Checking if the *input_images* and *compresed_images* folders exist in the HDFS , if not to be created using

```
hdfs dfs -mkdir -p /user/hadoop/compressed_images  
hdfs dfs -mkdir -p /user/hadoop/input_images
```

Upload the images to the designated folder in HDFS

Using the command below allows the upload of the files from the local machine to the HDFS system

```
hdfs dfs -put ~/Desktop/Hadoop-job/images/*  
/user/hadoop/input_images
```

```

latif@master:~/Desktop/mapreduce-job$ hdfs dfs -put ~/Desktop/Hadoop-job/images/* /user/hadoop/input_images
latif@master:~/Desktop/mapreduce-job$ hdfs dfs -ls /user/hadoop/input_images
Found 3 items
-rw-r--r-- 2 latif supergroup 2108235 2025-02-01 23:18 /user/hadoop/input_images/20241231_111322.jpg
-rw-r--r-- 2 latif supergroup 2464886 2025-02-01 23:18 /user/hadoop/input_images/20250107_093226.jpg
-rw-r--r-- 2 latif supergroup 240045 2025-02-01 23:18 /user/hadoop/input_images/IMG_20241231_202142_665.jpg

```

Job submission

Hadoop jar command communicates with the cluster manager (YARN) to allocate resources on available nodes.

```

hadoop jar imagecompressor.jar ImageCompressor
/usr/hadoop/input_images /user/hadoop/compressed_images

```

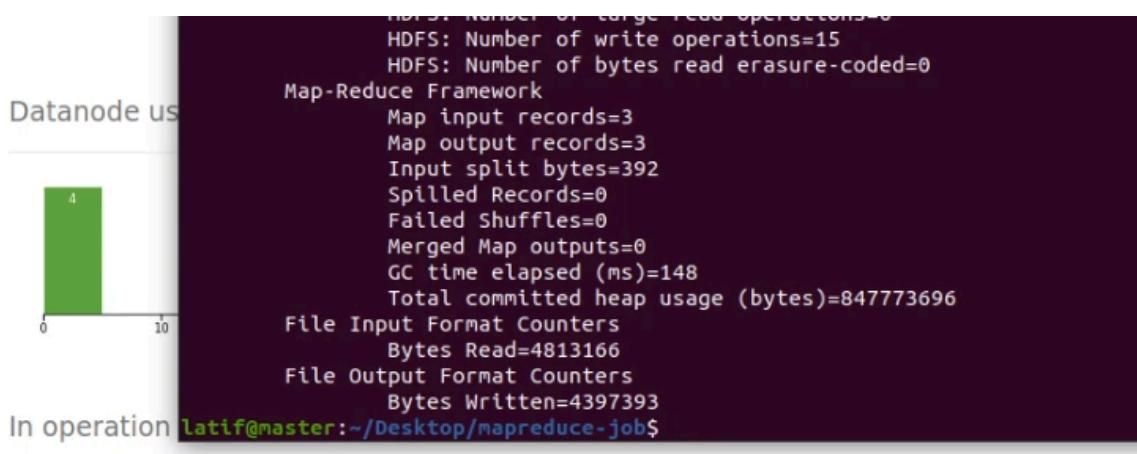
```

latif@master:~/Desktop/mapreduce-job$ hadoop jar /home/latif/Desktop/mapreduce-job/imagecompressor.jar ImageCompressor
/usr/hadoop/input_images /user/hadoop/compressed_images
2025-02-01 23:19:28,371 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2025-02-01 23:19:28,508 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2025-02-01 23:19:28,508 INFO impl.MetricsSystemImpl: JobTracker metrics system started

```

Running and waiting for job compilation

When a mapreduce job is submitted, it is compiled and prepared for execution across the cluster's nodes. Once the compilation is complete, the job is scheduled by the cluster manager, and tasks are distributed to available worker nodes for parallel processing.



Downloading the output files

Using the command

```

hdfs dfs -get /user/hadoop/compressed_images
~/Desktop/mapreduce-job/compressed
Bytes Written=4397393
latif@master:~/Desktop/mapreduce-job$ hdfs dfs -get /user/hadoop/compressed_images ~/Desktop/mapreduce-job/compressed

```

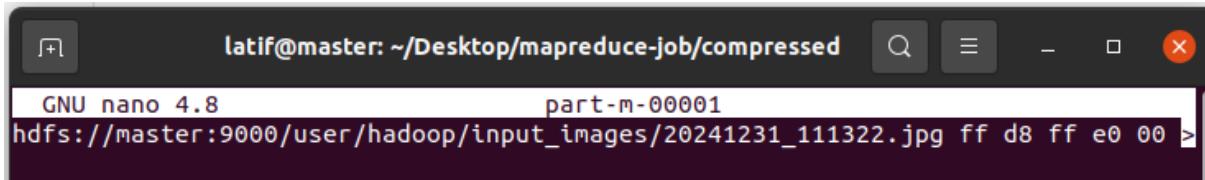
Fix the output file to see the results

Hadoop outputs images data in an incompatible format (e.g., hex or with extra metadata). This process ensures we get a clean binary JPEG file.

To visualise the file , we need to go through the following steps

-open image using sudo nano

-remove the first line which can look something like



```
latif@master: ~/Desktop/mapreduce-job/compressed
GNU nano 4.8
part-m-00001
hdfs://master:9000/user/hadoop/input_images/20241231_111322.jpg ff d8 ff e0 00 >
```

-use this after to reproduce

```
xxd -r -p part-m-00001 output.jpeg
```

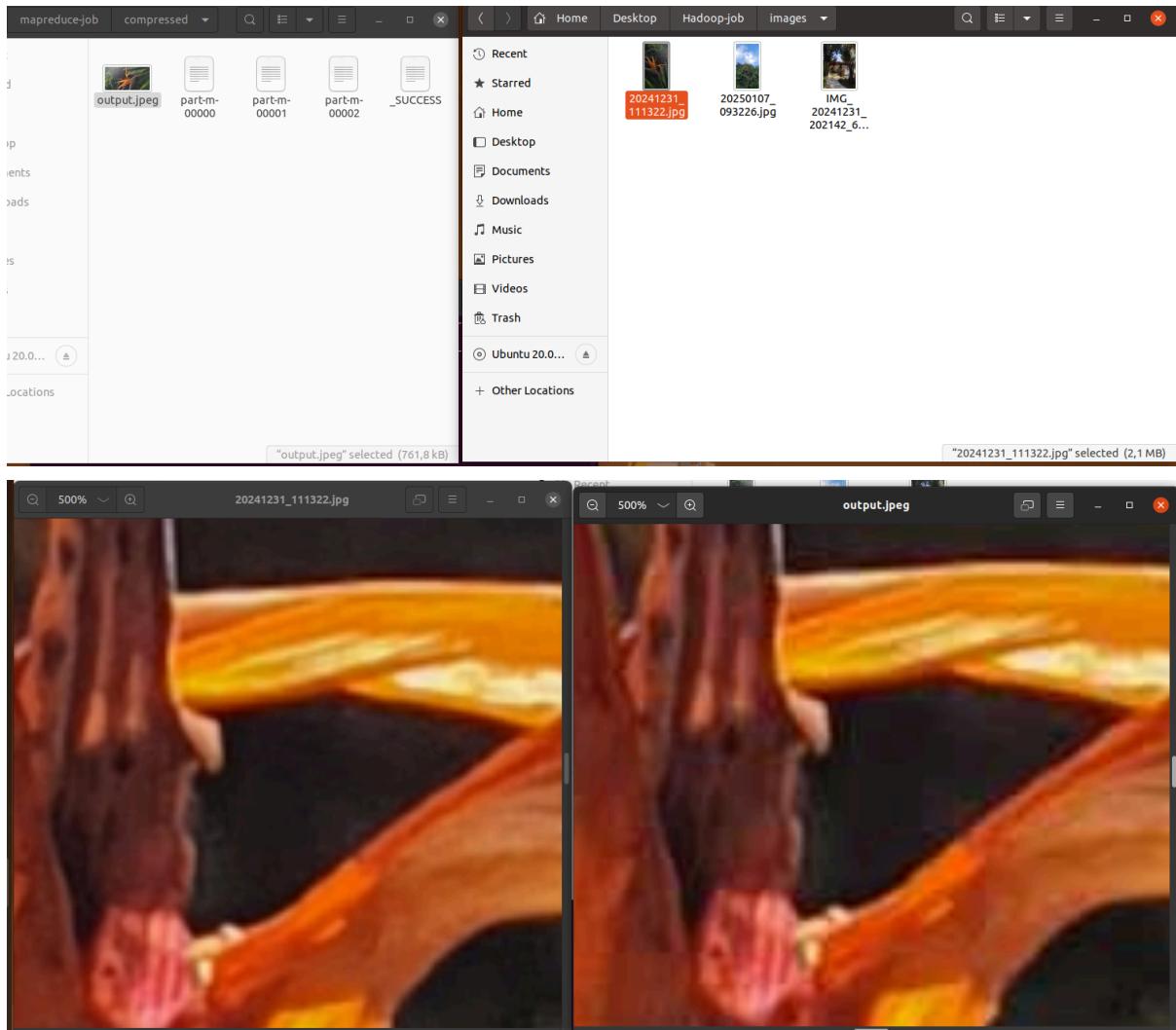
–And get something like this after



Results

The compression effectively reduced image file sizes while maintaining high visual quality.

For example, a **flower image originally sized at 2.1MB was compressed to 761.8KB**, demonstrating a significant reduction in storage space. At normal viewing distances, the quality appeared nearly identical to the original. However, upon zooming beyond **400% scale**, subtle differences became noticeable, such as minor loss of fine details and texture. This result confirms that **JPEG achieves high compression efficiency while preserving essential image quality**, making it a viable solution for reducing storage needs without a significant visual downgrade.



Conclusion and Discussion

In this project, we compared two approaches for processing and compressing images using: **JPEG MapReduce alone** and **JPEG2000 on Hadoop with Spark**. Both methods successfully performed the task, demonstrating the effectiveness of distributed computing for handling large-scale data.

When running MapReduce, the execution was relatively smooth on our setup, and the system remained responsive throughout the process. On the other hand, running Hadoop with Spark introduced significant delays before the actual job execution. The initialization phase of Spark, including job scheduling and resource allocation, caused noticeable slowdowns, even leading to system freezes at times.

However, despite the performance issues observed in our specific setup—running five virtual machines on a single host—both approaches yielded promising results. It is important to consider that Spark is designed for large-scale distributed processing and typically requires **more powerful hardware** to fully leverage its advantages, such as in-memory

computation and optimized parallel execution. If executed on a more capable cluster, **Spark's performance could have been significantly better**, potentially surpassing MapReduce in efficiency and speed.

Ultimately, while our findings suggest that MapReduce was more stable on our limited hardware, we **cannot definitively say one approach is superior** based solely on this experiment. Both technologies are built to handle **big data challenges**, which typically require enterprise-grade infrastructure rather than a modest personal setup. In a more powerful environment, Spark could demonstrate its full potential, making it a strong alternative to traditional MapReduce for large-scale data processing

Notes

To provide further insight into the project, the following resources are available:

- **Demo Videos:** [Demo videos](#)
 - This includes recorded demonstrations showcasing the model in action, its predictions, and key results. It provides a visual overview of how the approach performs in real-world scenarios.
- **GitHub Repository:** [Github repository](#)
 - The full implementation, including data preprocessing, model training, evaluation scripts, and additional resources, can be found in this repository. Feel free to explore the code, run experiments, or contribute improvements.