## Experiment No: 07

**Title**: Installing Hadoop and implement program using MapReduce.

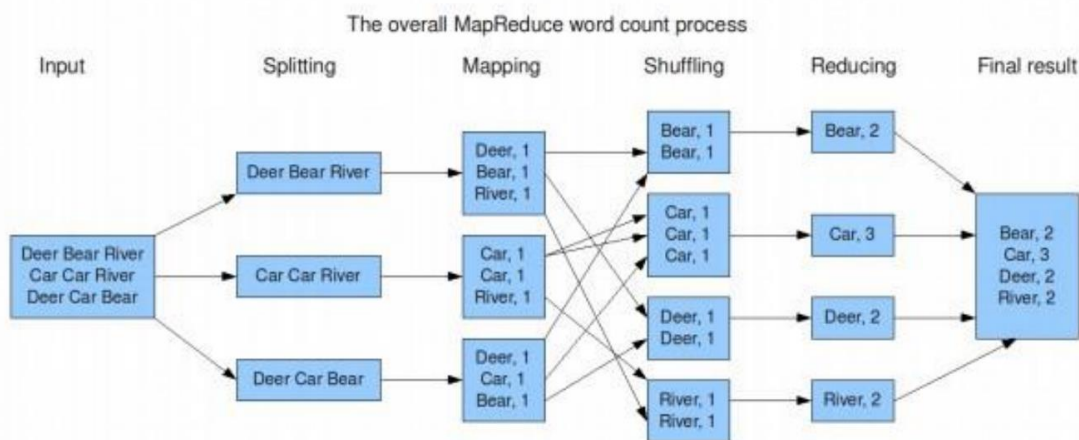**Aim:** To install Hadoop and implement program using MapReduce.

Theory:

Map-Reduce is a programming model that is mainly divided into two phases:
Map Phase and Reduce Phase.
It is designed for processing the data in parallel which is divided on various
machines(nodes). The Hadoop Java programs consist of Mapper class and
Reducer class along with the driver class. Hadoop Mapper is a function or task
which is used to process all input records from a file and generate the output
which works as input for Reducer. It produces the output by returning new keyvalue pairs.
The input data has to be converted to key-value pairs as Mapper
can not process the raw input records or tuples(key-value pairs). The mapper
also generates some small blocks of data while processing the input records as
a key-value pair. we will discuss the various process that occurs in Mapper,
There key features and how the key-value pairs are generated in the Mapper.
In MapReduce word count example, we find out the frequency of each word. Here, the role
of Mapper is to map the keys to the existing values and the role of Reducer is to aggregate
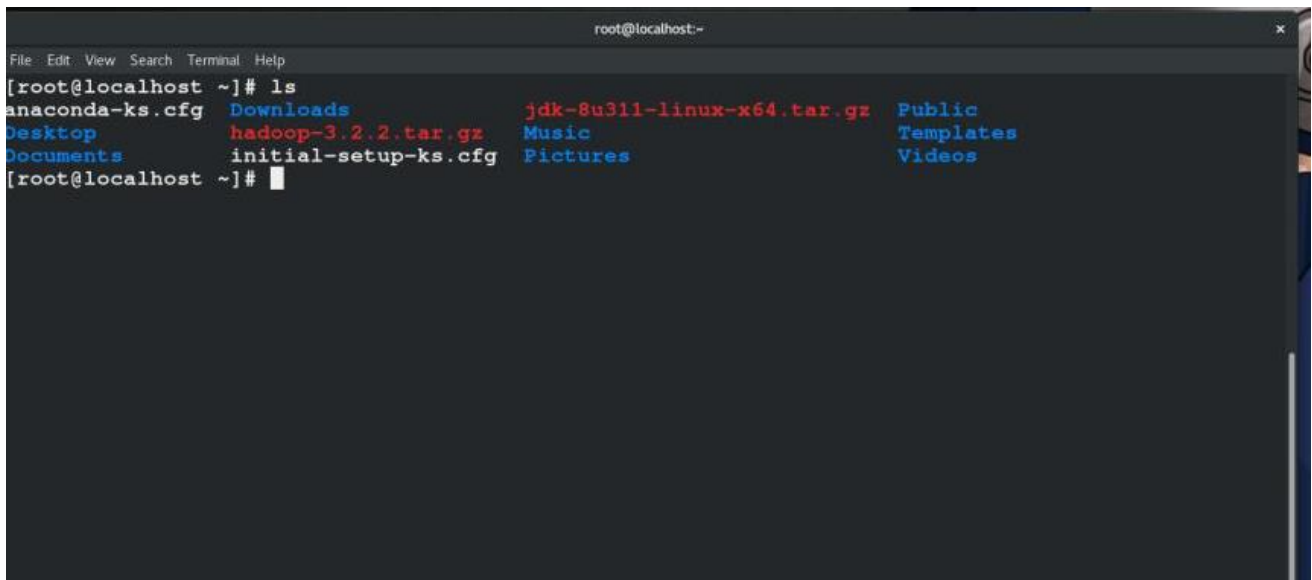the keys of common values.
Example :-



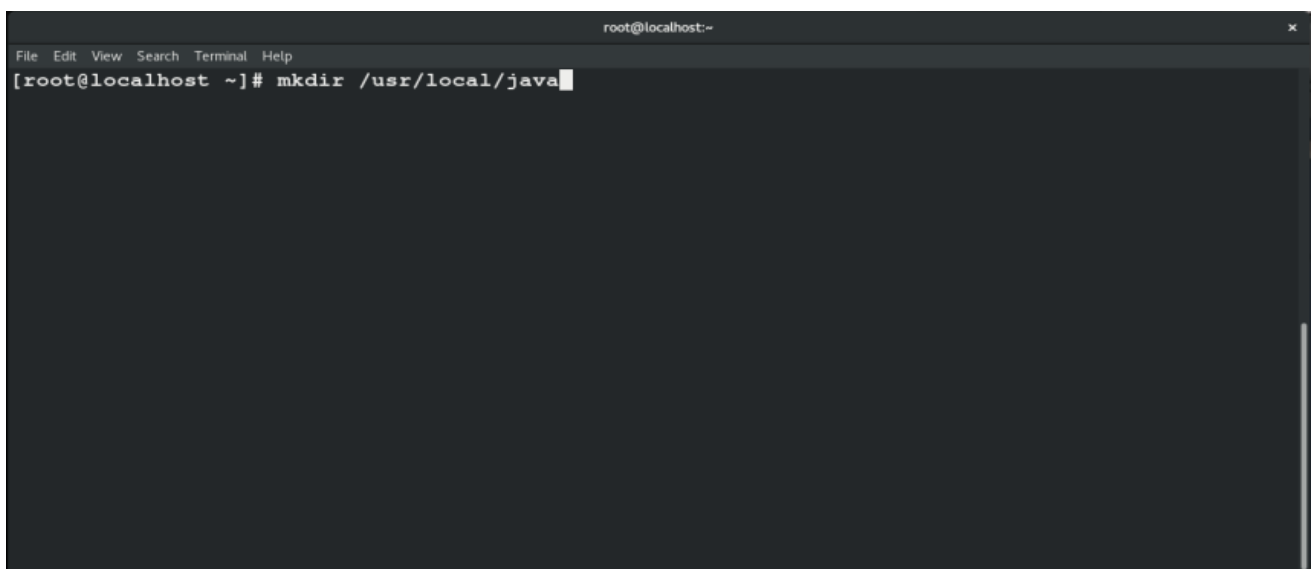The overall MapReduce word count process

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.

At the moment, Apache Hadoop 3.x fully supports Java 8. The OpenJDK 8 package in RedHat 8 contains both the runtime environment and development kit. Move jdk from desktop to /usr/local/java

**Practical:**

```
root@localhost:~
File Edit View Search Terminal Help
[root@localhost ~]# cp jdk-8u311-linux-x64.tar.gz /usr/local/java
```

Extract the jdk-8u311-linux-x64.tar.gz

```
root@localhost:/usr/local/java                                          x
File Edit View Search Terminal Help
[root@localhost java]# tar xvzf jdk-8u311-linux-x64.tar.gz --force
```

Rename the jdk-8u311-linux-x64.tar.gz into jdk

```
File Edit View Search Terminal Help
[root@localhost java]# ln -s jdk1.8.0_311 jdk
```

In the /etc/profiles we will also set up some of the required system variables and further inform our system regarding those updates. We also need to set oracle java as the default java



Now we need to update and install the alternatives variable for java and javac

Now we need to set the java and javac



```
root@localhost:/usr/local/java
File  Edit  View  Search  Terminal  Help
[root@localhost java]# update-alternatives --set java /usr/local/java/jdk/bin/java
```



```
root@localhost:/usr/local/java
File  Edit  View  Search  Terminal  Help
[root@localhost java]# update-alternatives --set javac /usr/local/java/jdk/bin/javac
```

Once the necessary things are done we will be restarting the /etc/profiles so that the updates will be implemented



```
root@localhost:/usr/local/java
File  Edit  View  Search  Terminal  Help
[root@localhost java]# . /etc/profile
```
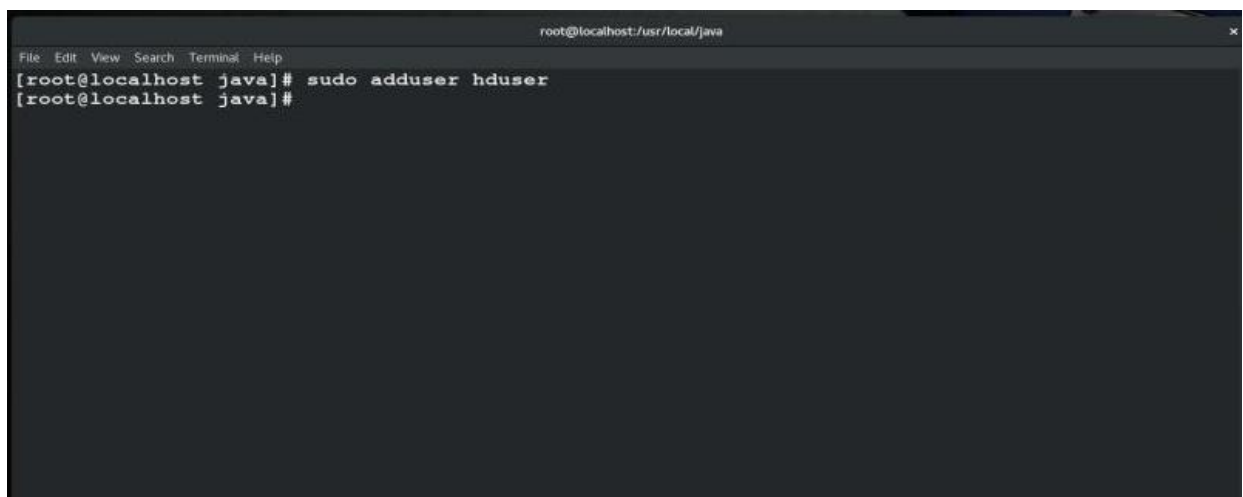
Java install successfully in our system

**Set Up a Non-Root User for Hadoop Environment**

It is advisable to create a non-root user, specifically for the Hadoop environment. A distinct user improves security and helps you manage your cluster more efficiently. To ensure the smooth functioning of Hadoop services, the user should have the ability to establish a passwordless SSH connection with the localhost.

**Create RedHat User**

Utilize the **adduser** command to create a new Hadoop user:

There are multiple situations where hduser might need the root power so for this we need to do the necessary updates in the /etc/sudoer file

```
##      user      MACHINE=COMMANDS
##
## The COMMANDS section may have other options added to it.
##
## Allow root to run any commands anywhere
root    ALL=(ALL)      ALL
hduser  ALL=(ALL)      ALL

## Allows members of the 'sys' group to run networking, software,
## service management apps and more.
# %sys ALL = NETWORKING, SOFTWARE, SERVICES, STORAGE, DELEGATING, PROCESSES, LOCATE, DRIVERS

## Allows people in group wheel to run all commands
%wheel  ALL=(ALL)      ALL

## Same thing without a password
# %wheel        ALL=(ALL)        NOPASSWD: ALL

## Allows members of the users group to mount and unmount the
## cdrom as root
# %users  ALL=/sbin/mount /mnt/cdrom, /sbin/umount /mnt/cdrom

## Allows members of the users group to shutdown this system
-- INSERT --                                              101,23-28      95%
```

**Install OpenSSH on Redhat**

Install the OpenSSH server and client using the following command:

**sudo yum install openssh-server openssh-client -y**

```
[root@localhost java]# yum  install openssh-server
Updating Subscription Management repositories.
Unable to read consumer identity
This system is not registered to Red Hat Subscription Management. You can use subscription-manag
er to register.
Repository 'AppStream' is missing name in configuration, using id.
Repository 'BaseOS' is missing name in configuration, using id.
Last metadata expiration check: 0:49:46 ago on Tuesday 01 February 2022 07:22:19 PM IST.
Package openssh-server-7.8p1-4.el8.x86_64 is already installed.
Dependencies resolved.
Nothing to do.
Complete!
[root@localhost java]#
```

**Enable Passwordless SSH for Hadoop User**

Generate an SSH key pair and define the location is is to be stored in:

**cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys**
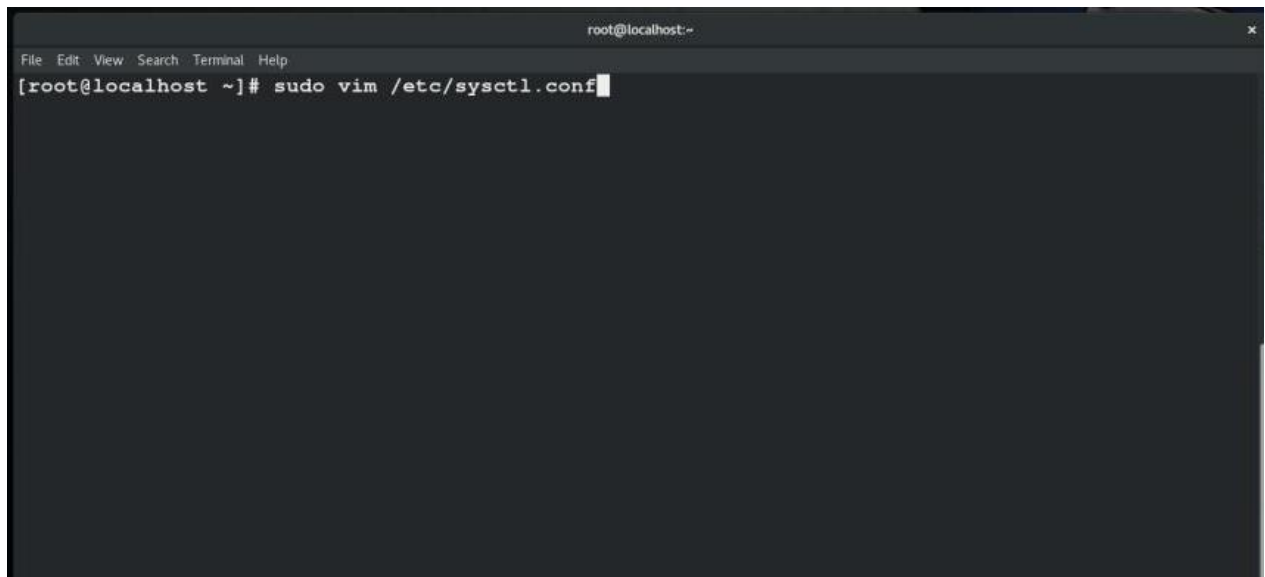
The system proceeds to generate and save the SSH key pair.

Use the cat command to store the public key as authorized_keys in the *ssh* directory:

We will also disable the ipv6 and only use the ipv4 in the machine

**to crosscheck we can use cat over /proc/sys/net/ipv6/conf/all/disable_ipv6 file**



Move hadoop tar file from ~ file to /usr/local

```
[root@localhost local]# tar -xzvf hadoop-3.2.2.tar.gz --force
```
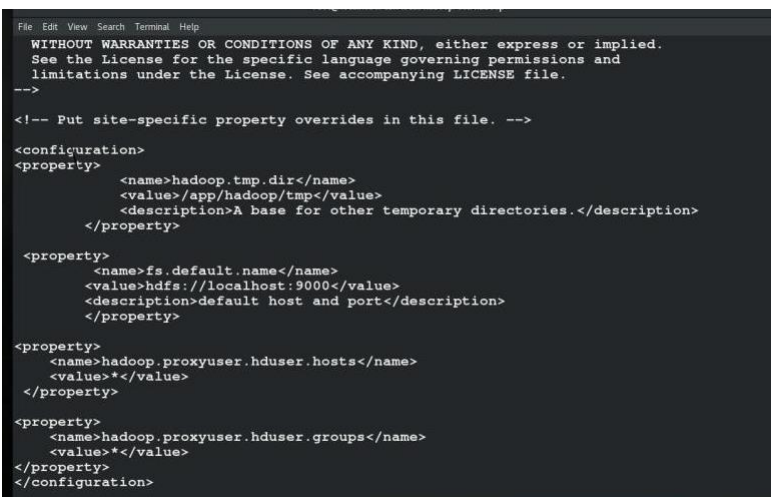
Now its time for the hadoop files
Configuration Changes in yarn-site.xml file
Edit **yarn-site.xml** with the following entries.



```
 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.nodemanager.aux-services</name>
         <value>mapreduce_shuffle</value>
    </property>
         <property>
    <name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
    <value>org.apache.hadoop.mapred.ShuffleHandler</value>
         </property>
 <property>
    <name>yarn.nodemanager.vmem-check-enabled</name>
    <value>false</value>
    <description>Whether virtual memory limits will be enforced for containers</description>
 </property>

 <property>
    <name>yarn.nodemanager.vmem-pmem-ratio</name>
    <value>4</value>
    <description>Ratio between virtual memory to physical memory when setting memory limits fo
r containers</description>
 </property>
</configuration>
~
-- INSERT --                                                          35,14          Bot
```
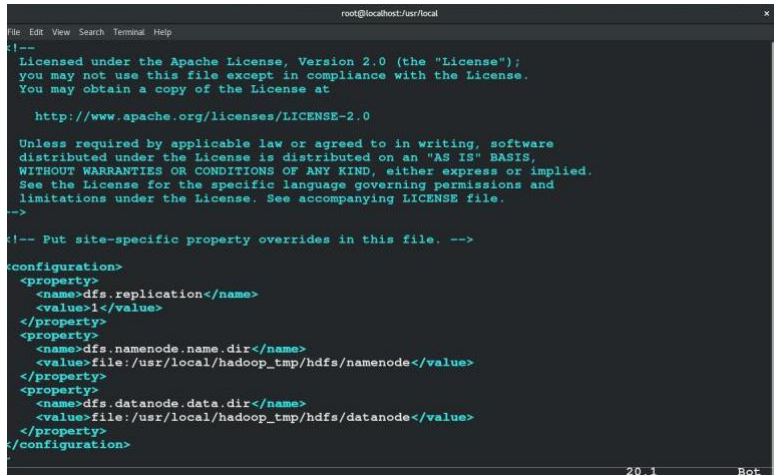
Configuration Changes in core-site.xml file

Edit the **core-site.xml** with vim or you can use any of the editors. The file is under

**/etc/hadoop** inside **hadoop** home directory and add following entries.



```
 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>hadoop.tmp.dir</name>
         <value>/app/hadoop/tmp</value>
         <description>A base for other temporary directories.</description>
         </property>

 <property>
         <name>fs.default.name</name>
         <value>hdfs://localhost:9000</value>
         <description>default host and port</description>
         </property>

<property>
    <name>hadoop.proxyuser.hduser.hosts</name>
    <value>*</value>
 </property>

<property>
    <name>hadoop.proxyuser.hduser.groups</name>
    <value>*</value>
</property>
</configuration>
```

**Configuration Changes in mapred-site.xml file**

Copy the mapred-site.xml from mapred-site.xml.template using cp command and then

edit the mapred-site.xml placed in /etc/hadoop under hadoop installation directory with
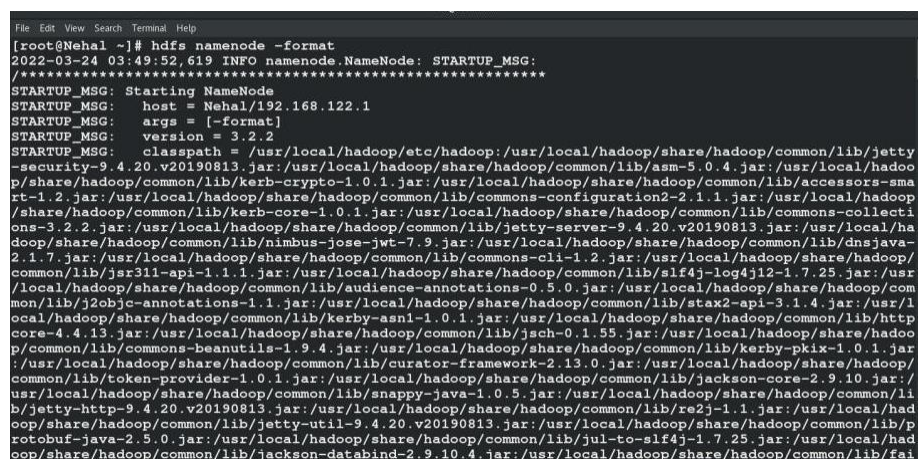
the following changes.



Now create a namenode and datanode folder and provide the all the necessary

permission to it

# Sudo mkdir -p /usr/local/hadoop_tmp/hdfs/namenode

# Sudo mkdir -p /usr/local/hadoop_tmp/hdfs/datanode

Starting the Hadoop Cluster

Format the namenode before using it for the first time. As hadoop users run the below

command to format the Namenode.



Once the Namenode has been formatted then start the HDFS using the
$ start-all.sh
All Services started successfully and all the node are

Working







Conclusion: Thus we have installed Hadoop and implemented program using MapReduce.