# <u>HADOOP</u>

## <u>SET UP A SINGLE HADOOP CLUSTER AND SHOW THE PROCESS USING WEB UI</u>

**AIM:**

To set-up one node Hadoop cluster.

**PROCEDURE:**

1.  System Update
2.  Install Java
3.  Add a dedicated Hadoop user
4.  Install SSH and setup SSH certificates
5.  Check if SSH works
6.  Install Hadoop
7.  Modify Hadoop config files
8.  Format Hadoop filesystem
9.  Start Hadoop
10. Check Hadoop through web UI
11. Stop Hadoop

**THEORY**

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 a single server to thousands of machines, each offering local computation and storage.

**HADOOP ARCHITECTURE**
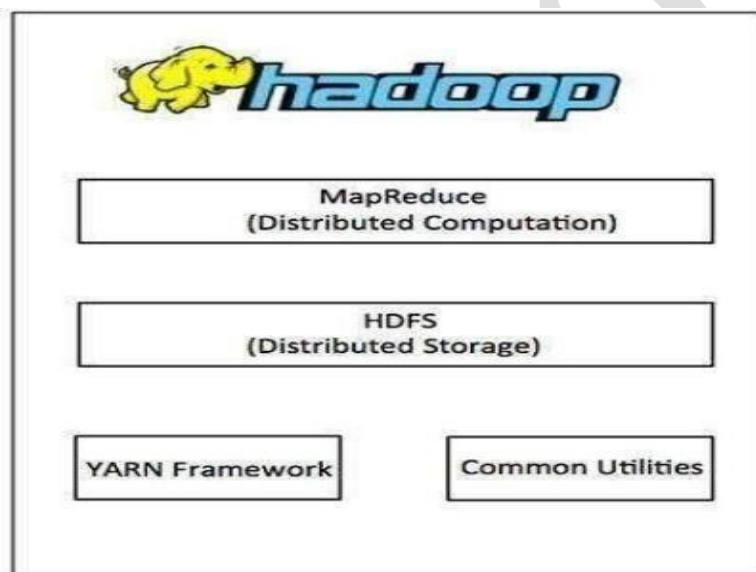
Hadoop framework includes following four modules:

Hadoop Common: These are Java libraries and utilities required by other Hadoop modules. These libraries provide filesystem and OS level abstractions and contain the necessary Java files and scripts required to start Hadoop.

Hadoop YARN: This is a framework for job scheduling and cluster resource management.

Hadoop Distributed File System (HDFS): A distributed file system that provides high-throughput access to application data.
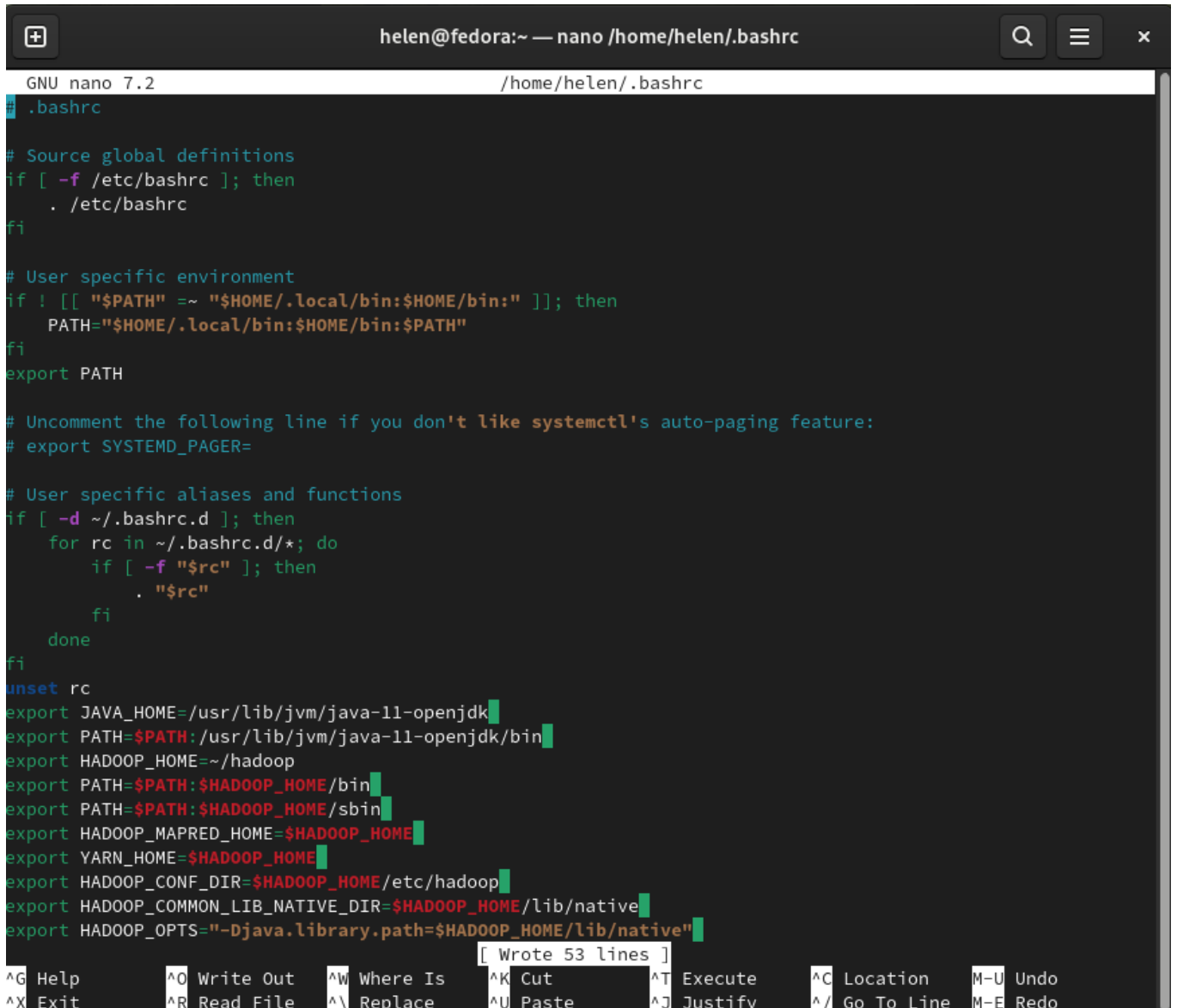
Hadoop MapReduce: This is a YARN-based system for parallel processing of large data sets.

We can use following diagram to depict these four components available in Hadoop framework.

**PROCEDURE**

**$ nano ~/.bashrc**

```
helen@fedora:~ — nano /home/helen/.bashrc

GNU nano 7.2                            /home/helen/.bashrc
# .bashrc

# Source global definitions
if [ -f /etc/bashrc ]; then
    . /etc/bashrc
fi

# User specific environment
if ! [[ "$PATH" =~ "$HOME/.local/bin:$HOME/bin:" ]]; then
    PATH="$HOME/.local/bin:$HOME/bin:$PATH"
fi
export PATH

# Uncomment the following line if you don't like systemctl's auto-paging feature:
# export SYSTEMD_PAGER=

# User specific aliases and functions
if [ -d ~/.bashrc.d ]; then
    for rc in ~/.bashrc.d/*; do
        if [ -f "$rc" ]; then
            . "$rc"
        fi
    done
fi
unset rc
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk
export PATH=$PATH:/usr/lib/jvm/java-11-openjdk/bin
export HADOOP_HOME=~/hadoop
export PATH=$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"
                              [ Wrote 53 lines ]
^G Help      ^O Write Out   ^W Where Is   ^K Cut      ^T Execute   ^C Location    M-U Undo
^X Exit      ^R Read File   ^\ Replace    ^U Paste    ^J Justify   ^/ Go To Line  M-E Redo
```

## $ nano $HADOOP_HOME/etc/hadoop/hadoop-env.sh

```
  GNU nano 7.2                    /home/helen/hadoop/etc/hadoop/hadoop-env.sh
#
# Licensed to the Apache Software Foundation (ASF) under one
# or more contributor license agreements.  See the NOTICE file
# distributed with this work for additional information
# regarding copyright ownership.  The ASF licenses this file
# to you 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.

# Set Hadoop-specific environment variables here.

##
## THIS FILE ACTS AS THE MASTER FILE FOR ALL HADOOP PROJECTS.
## SETTINGS HERE WILL BE READ BY ALL HADOOP COMMANDS.  THEREFORE,
## ONE CAN USE THIS FILE TO SET YARN, HDFS, AND MAPREDUCE
## CONFIGURATION OPTIONS INSTEAD OF xxx-env.sh.
##
## Precedence rules:
##
## {yarn-env.sh|hdfs-env.sh} > hadoop-env.sh > hard-coded defaults
##
## {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:

^G Help        ^O Write Out   ^W Where Is    ^K Cut         ^T Execute     ^C Location    M-U Undo
^X Exit        ^R Read File   ^\ Replace     ^U Paste       ^J Justify     ^/ Go To Line  M-E Redo
```

## $nano $HADOOP_HOME/etc/hadoop/core-site.xml
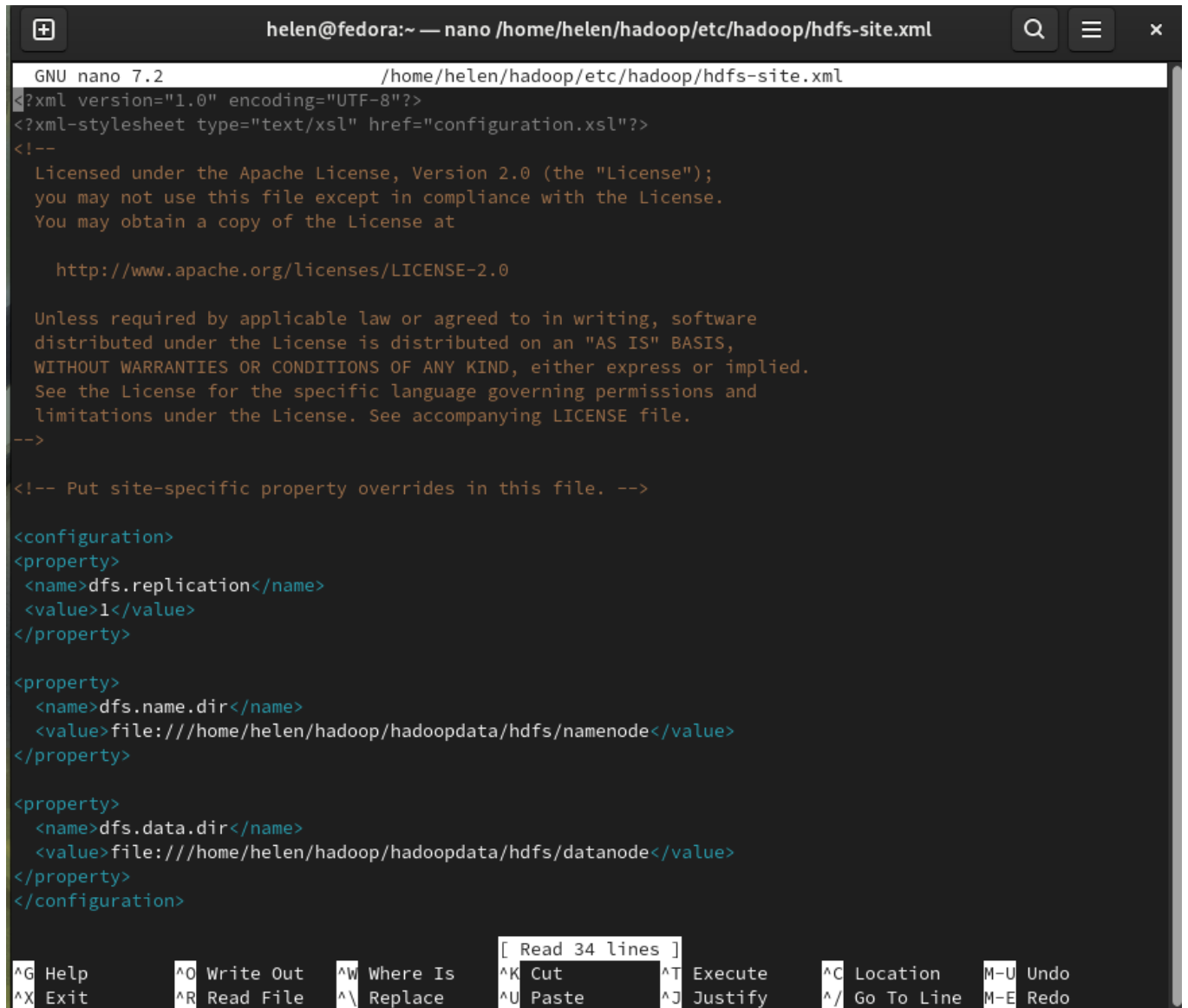
```
  GNU nano 7.2                    /home/helen/hadoop/etc/hadoop/core-site.xml
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
  Licensed under the Apache License, Version 2.0 (the "License");
  you may not use this file except in compliance 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.
-->

<!-- Put site-specific property overrides in this file. -->
<configuration>
<property>
  <name>fs.default.name</name>
  <value>hdfs://localhost:9000</value>
</property>
</configuration>
```

**$nano $HADOOP_HOME/etc/hadoop/hdfs-site.xml**

```
helen@fedora:~ — nano /home/helen/hadoop/etc/hadoop/hdfs-site.xml

  GNU nano 7.2                    /home/helen/hadoop/etc/hadoop/hdfs-site.xml
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
  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.
-->

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

<configuration>
<property>
 <name>dfs.replication</name>
 <value>1</value>
</property>

<property>
  <name>dfs.name.dir</name>
  <value>file:///home/helen/hadoop/hadoopdata/hdfs/namenode</value>
</property>

<property>
  <name>dfs.data.dir</name>
  <value>file:///home/helen/hadoop/hadoopdata/hdfs/datanode</value>
</property>
</configuration>

                              [ Read 34 lines ]
^G Help        ^O Write Out   ^W Where Is    ^K Cut       ^T Execute    ^C Location    M-U Undo
^X Exit        ^R Read File   ^\ Replace     ^U Paste     ^J Justify    ^/ Go To Line  M-E Redo
```
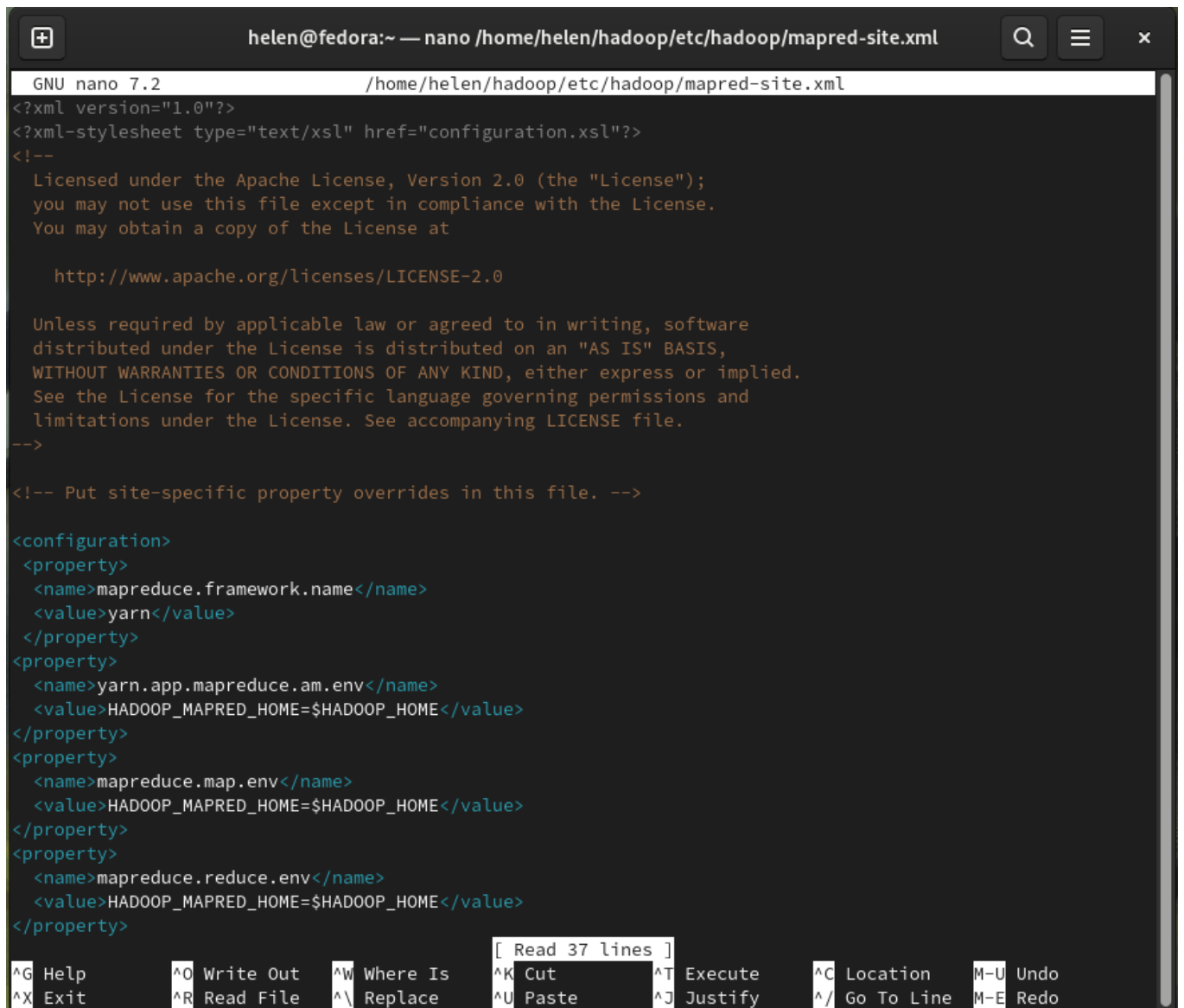
**$nano $HADOOP_HOME/etc/hadoop/mapred-site.xml**

```
GNU nano 7.2                         /home/helen/hadoop/etc/hadoop/mapred-site.xml
<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
  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.
-->

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

<configuration>
 <property>
  <name>mapreduce.framework.name</name>
  <value>yarn</value>
 </property>
<property>
  <name>yarn.app.mapreduce.am.env</name>
  <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
</property>
<property>
  <name>mapreduce.map.env</name>
  <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
</property>
<property>
  <name>mapreduce.reduce.env</name>
  <value>HADOOP_MAPRED_HOME=$HADOOP_HOME</value>
</property>
                                     [ Read 37 lines ]
^G Help        ^O Write Out   ^W Where Is    ^K Cut         ^T Execute     ^C Location    M-U Undo
^X Exit        ^R Read File   ^\ Replace     ^U Paste       ^J Justify     ^/ Go To Line  M-E Redo
```

**$nano $HADOOP_HOME/etc/hadoop/yarn-site.xml**

```
helen@fedora:~ — nano /home/helen/hadoop/etc/hadoop/yarn-site.xml

  GNU nano 7.2                    /home/helen/hadoop/etc/hadoop/yarn-site.xml
<?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>
 <property>
  <name>yarn.nodemanager.aux-services</name>
  <value>mapreduce_shuffle</value>
 </property>
</configuration>
```

**$ start-all.sh**

```
helen@fedora:~ — bash /home/helen/hadoop/sbin/start-all.sh

helen@fedora:~$ nano ~/.bashrc
helen@fedora:~$ nano $HADOOP_HOME/etc/hadoop/hadoop-env.sh
helen@fedora:~$ nano $HADOOP_HOME/etc/hadoop/core-site.xml
helen@fedora:~$ nano $HADOOP_HOME/etc/hadoop/hdfs-site.xml
helen@fedora:~$ nano $HADOOP_HOME/etc/hadoop/mapred-site.xml
helen@fedora:~$ nano $HADOOP_HOME/etc/hadoop/yarn-site.xml
helen@fedora:~$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as helen in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [fedora]
Starting resourcemanager
Starting nodemanagers
helen@fedora:~$
```

**$ jps**

```
helen@fedora:~$ jps
3716 NameNode
3908 DataNode
4485 NodeManager
4345 ResourceManager
5081 Jps
4108 SecondaryNameNode
```

## localhost:9870

### Overview

| | |
|---|---|
| **Version** | 3.3.6 |
| **Compiled** | 2023-06-18T08:22Z by ubuntu from (HEAD detached at release-3.3.6-RC1) |
| **NameNode Address** | localhost:9000 |
| **Started** | Wed Aug 14 21:51:32 -0400 2024 |
| **Last Checkpoint** | Never |
| **Checkpoint Period** | 3600 seconds |
| **Checkpoint Transactions** | 1000000 |

Checkpoint Image URI

- file:///tmp/hadoop-kali/dfs/namesecondary

Checkpoint Editlog URI

- file:///tmp/hadoop-kali/dfs/namesecondary

Hadoop, 2023.

## localhost:8088



**RESULT:**

Thus, Hadoop has been successfully installed.