# Apache Hadoop 3.3.0 installation on Ubuntu Part 1

**Apache Hadoop 3.3.0 installation on Ubuntu Part 1**

With this tutorial, we will learn the complete process to install Hadoop 3.3.1 on Ubuntu 20.

**Supported Java Versions**

- Apache Hadoop 3.3 and upper supports Java 8 and Java 11 (runtime only)
- Please compile Hadoop with Java 8. Compiling Hadoop with Java 11 is not supported: HADOOP-16795 – Java 11 compile support OPEN
- Apache Hadoop from 3.0.x to 3.2.x now supports only Java 8
- Apache Hadoop from 2.7.x to 2.10.x support both Java 7 and 8

**Required software for Linux include:**

- Java must be installed. Recommended Java versions are described at HadoopJavaVersions.
- ssh must be installed and sshd must be running to use the Hadoop scripts that manage remote Hadoop daemons if the optional start and stop scripts are to be used.

**Steps for Installing JAVA 8 on Ubuntu**

Step 1 – Install Java 8 on Ubuntu

The OpenJDK 8 is available under default Apt repositories. You can simply install Java 8 on an Ubuntu system using the following commands.

```
1. $sudo apt update
2. $sudo apt install openjdk-8-jdk -y
```

Step 2 – Verify Java Installation

You have successfully installed Java 8 on your system. Let's verify the installed and current active version using the following command.

```
1. $java -version
2. openjdk version "1.8.0_252"
3. OpenJDK Runtime Environment (build 1.8.0_252-8u252-b09-1ubuntu1-b09)
4. OpenJDK 64-Bit Server VM (build 25.252-b09, mixed mode)
```
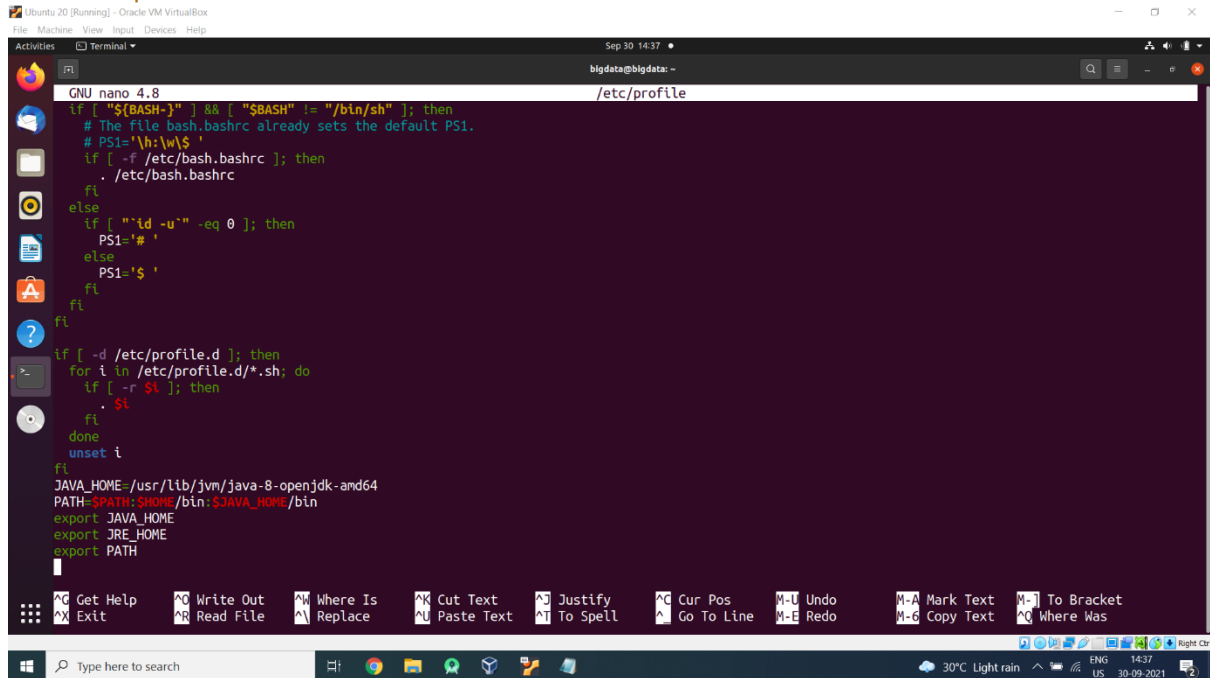
Step 3 – Setup JAVA_HOME and JRE_HOME Variable

As you have installed Java on your Linux system, You must have to set JAVA_HOME and JRE_HOME environment variables,

Edit the system Path file /etc/profile

```
sudo nano /etc/profile
```

Add the following lines at the end

1. JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
2. PATH=$PATH:$HOME/bin:$JAVA_HOME/bin
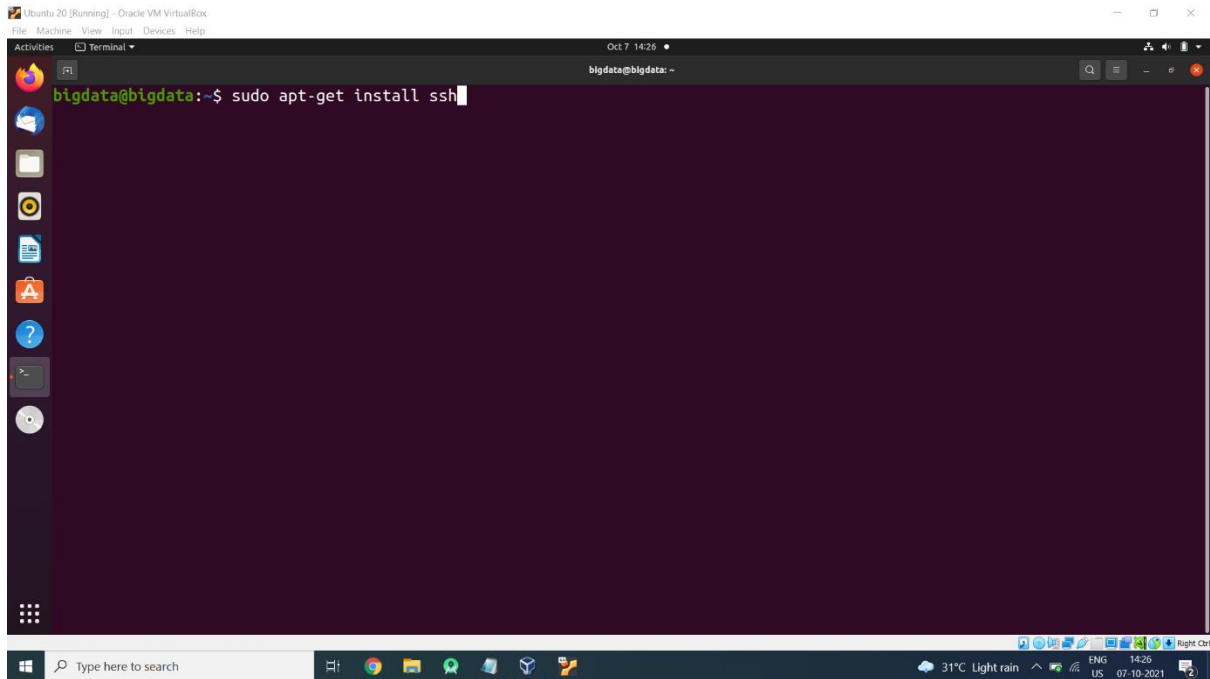3. export JAVA_HOME
4. export JRE_HOME
5. export PATH



## Steps for Installing ssh on Ubuntu

Secure Shell (SSH) is a cryptographic network protocol for operating network services securely over an unsecured network. Typical applications include remote command-line, login, and remote command execution, but any network service can be secured with SSH.

Install ssh on your system using the below command:

```
sudo apt-get install ssh
```

Type the password for the sudo user and then press Enter.

Install pdsh on your system using the below command:

```
sudo apt-get install pdsh
```



Type 'Y' and then press Enter to continue with the installation process.

Open the .bashrc file in the nano editor using the following command:

```
nano .bashrc
```

Now set the PDSH_RCMD_TYPE environment variable to ssh



## Steps for Installing Hadoop on Ubuntu

- Create a directory for example

$mkdir /home/bigdata/hadoop

- Move to hadoop directory

$cd /home/bigdata/hadoop

Download Hadoop (Link will change with respect to country so please get the download link from hadoop website ie https://hadoop.apache.org/releases.html



A new web page will get open and copy the link



In Ubuntu terminal type

```
$wget https://dlcdn.apache.org/hadoop/common/hadoop-3.3.1/hadoop-3.3.1.tar.gz
```

Then type

1. `$tar xvf hadoop-3.3.1.tar.gz`
2. `$cd hadoop-3.3.1`

```
bigdata@bigdata:~/hadoop$ pwd
/home/bigdata/hadoop
bigdata@bigdata:~/hadoop$ ls -ltr
total 591012
-rw-rw-r--  1 bigdata bigdata 605187279 Jun 15 15:25 hadoop-3.3.1.tar.gz
drwxr-xr-x 11 bigdata bigdata      4096 Oct  4 15:09 hadoop-3.3.1
bigdata@bigdata:~/hadoop$ cd hadoop-3.3.1/
bigdata@bigdata:~/hadoop/hadoop-3.3.1$ ls -ltr
total 116
-rw-rw-r-- 1 bigdata bigdata   175 May 21 21:41 README.txt
-rw-rw-r-- 1 bigdata bigdata  1541 May 21 21:41 NOTICE.txt
-rw-rw-r-- 1 bigdata bigdata 29473 Jun 15 10:32 NOTICE-binary
-rw-rw-r-- 1 bigdata bigdata 15217 Jun 15 10:32 LICENSE.txt
-rw-rw-r-- 1 bigdata bigdata 23450 Jun 15 10:32 LICENSE-binary
drwxr-xr-x 3 bigdata bigdata  4096 Jun 15 10:45 sbin
drwxr-xr-x 3 bigdata bigdata  4096 Jun 15 10:45 etc
drwxr-xr-x 2 bigdata bigdata  4096 Jun 15 11:22 licenses-binary
drwxr-xr-x 3 bigdata bigdata  4096 Jun 15 11:22 lib
drwxr-xr-x 2 bigdata bigdata  4096 Jun 15 11:22 include
drwxr-xr-x 2 bigdata bigdata  4096 Jun 15 11:22 bin
drwxr-xr-x 4 bigdata bigdata  4096 Jun 15 11:22 libexec
drwxr-xr-x 4 bigdata bigdata  4096 Jun 15 11:48 share
drwxrwxr-x 3 bigdata bigdata  4096 Oct  4 15:25 logs
bigdata@bigdata:~/hadoop/hadoop-3.3.1$
```

Edit the file etc/hadoop/hadoop-env.sh to define some parameters as follows:

```
$cd etc/
```

`$nano hadoop-env.sh`

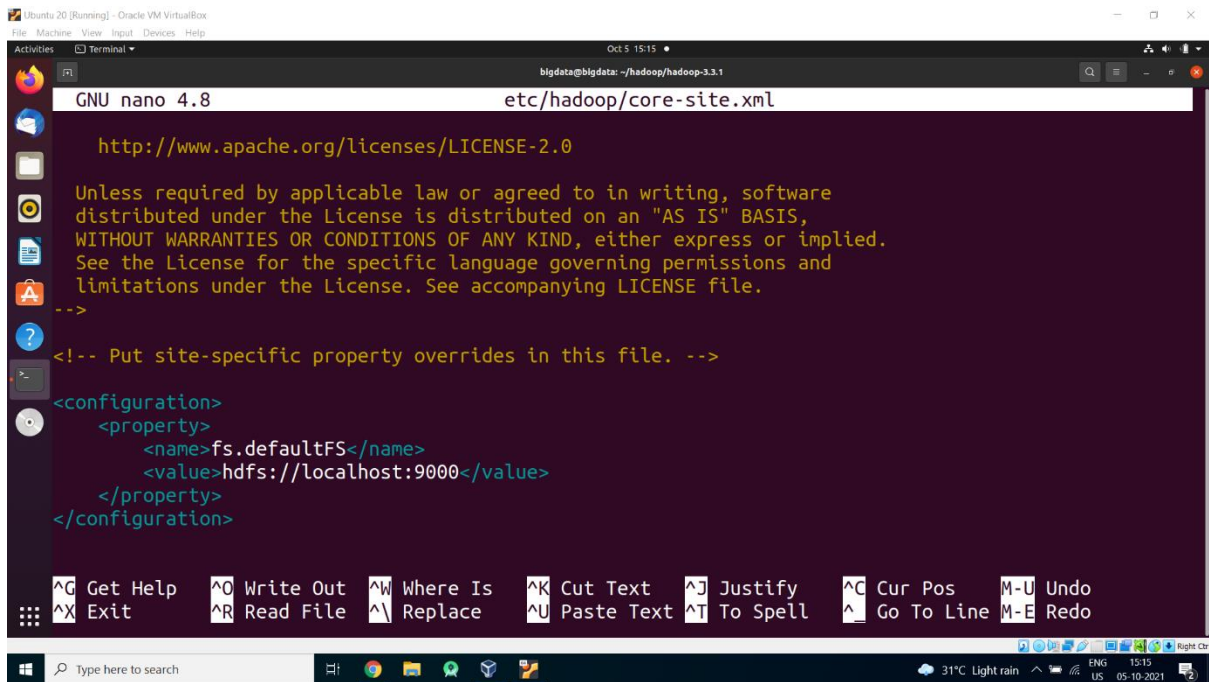Set the Java Path in hadoop-env.sh as shown in the image.

# Apache Hadoop 3.3.0 installation on Ubuntu Part 2

**Apache Hadoop 3.3.0 installation on Ubuntu Part 2**

Use the following property in the respective files

**File: nano etc/hadoop/core-site.xml:**

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



**File: nano etc/hadoop/hdfs-site.xml**

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

**File: nano etc/hadoop/mapred-site.xml**

<configuration>

<property>

<name>mapreduce.framework.name</name>

<value>yarn</value>

</property>

<property>

<name>mapreduce.application.classpath</name>

<value>
$HADOOP_MAPRED_HOME/share/hadoop/mapreduce/*:$HADOOP_MAPRED_HOM
E/share/hadoop/mapreduce/lib/*</value>

</property>

</configuration>

GNU nano 4.8                    etc/hadoop/mapred-site.xml

```
  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>mapreduce.application.classpath</name>
        <value>$HADOOP_MAPRED_HOME/share/hadoop/mapreduce/*:$HADOOP_MAPRED_HOME/share/hadoop/mapredu>
    </property>
</configuration>
```

**File: nano etc/hadoop/yarn-site.xml**

```
<configuration>
<property>
<name>yarn.nodemanager.aux-services</name>
<value>mapreduce_shuffle</value>
</property>
<property>
<name>yarn.nodemanager.env-whitelist</name>
<value>
JAVA_HOME,HADOOP_COMMON_HOME,HADOOP_HDFS_HOME,HADOOP_CONF_DIR,
CLASSPATH_PREPEND_DISTCACHE,
HADOOP_YARN_HOME,HADOOP_HOME,PATH,LANG,TZ,HADOOP_MAPRED_HOME*<
/value>
</property>
</configuration>
```

Now check that you can ssh to the localhost without a passphrase:

```
$ ssh localhost
```

If you cannot ssh to localhost without a passphrase, execute the following commands:
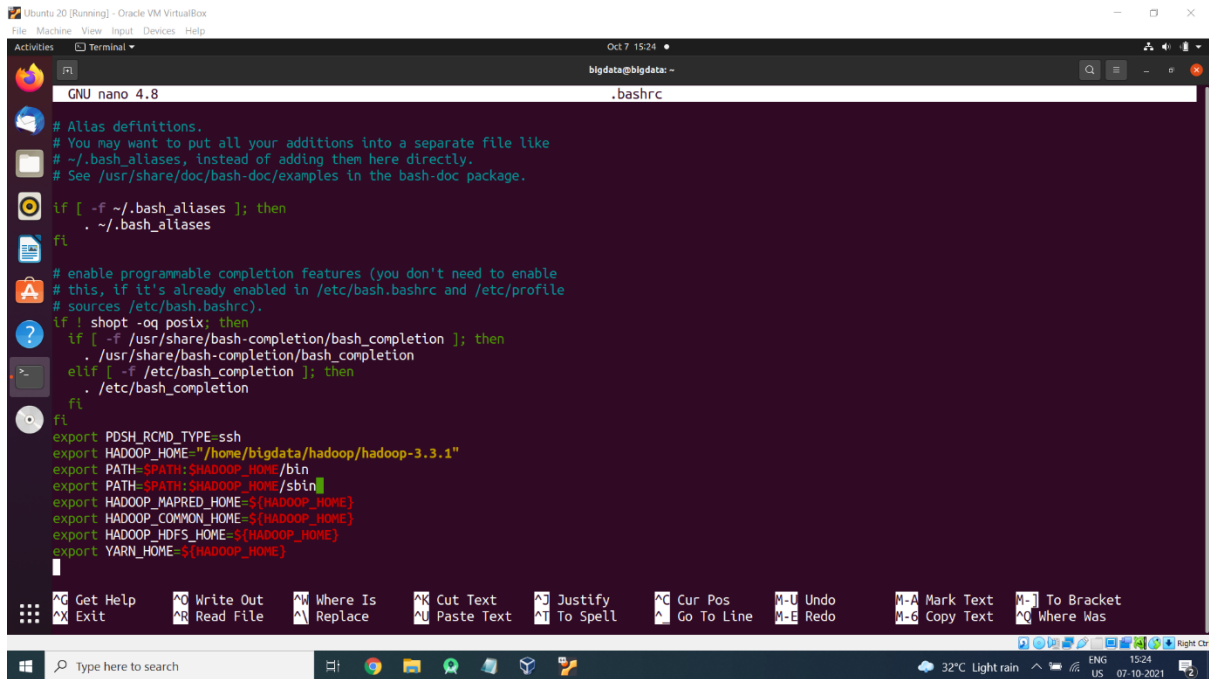
1. $ ssh-keygen -t rsa -P '' -f ~/.ssh/id_rsa
2. $ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
3. $ chmod 0600 ~/.ssh/authorized_keys

Open the bashrc files in the nano editor using the following command:

```
nano .bashrc
```

edit .bashrc file located in the user's home directory and add the following parameters:

1. export HADOOP_HOME="/home/bigdata/hadoop/hadoop-3.3.1"
2. export PATH=$PATH:$HADOOP_HOME/bin
3. export PATH=$PATH:$HADOOP_HOME/sbin
4. export HADOOP_MAPRED_HOME=${HADOOP_HOME}
5. export HADOOP_COMMON_HOME=${HADOOP_HOME}
6. export HADOOP_HDFS_HOME=${HADOOP_HOME}
7. export YARN_HOME=${HADOOP_HOME}

To save the changes you've made, press Ctrl+O. To exit the nano editor, press Ctrl+X and then press 'Y' to exit the editor.

Now, source the bashrc file so that the changes will come into effect:
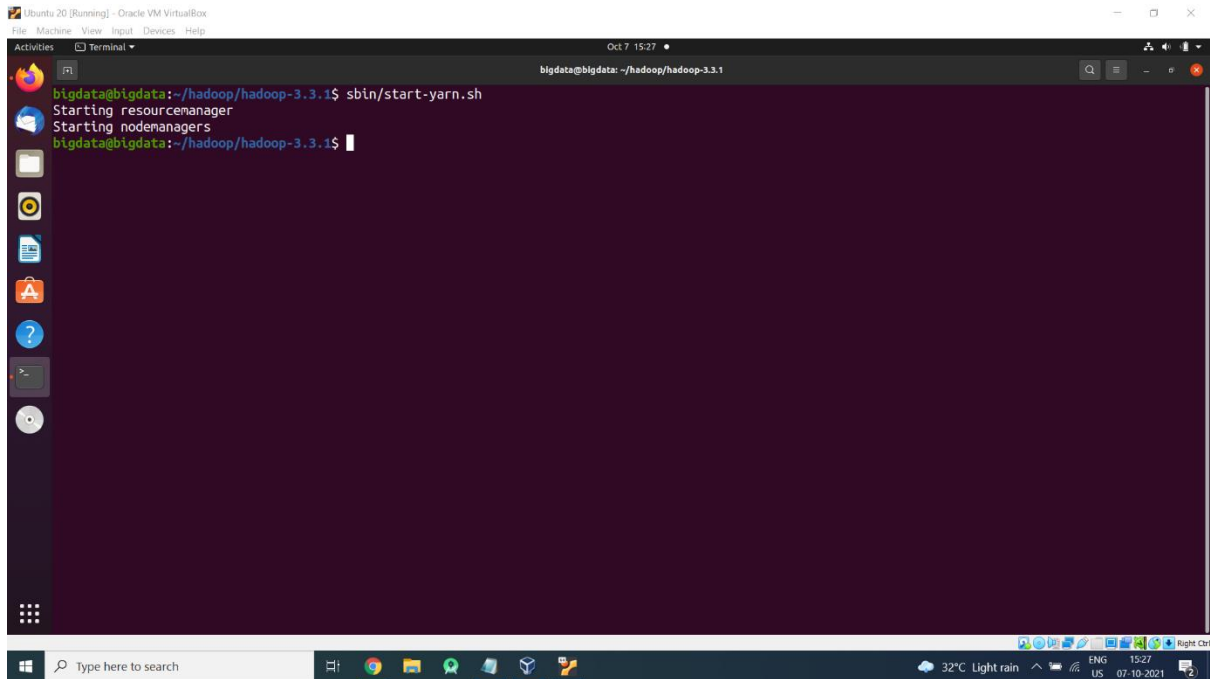
```
source ~/.bashrc
```

Format the filesystem:

```
$ bin/hdfs namenode -format
```

Start NameNode daemon and DataNode daemon:

```
$ sbin/start-dfs.sh
```

The hadoop daemon log output is written to the $HADOOP_LOG_DIR directory (defaults to $HADOOP_HOME/logs).

Browse the web interface for the NameNode; by default it is available at:

NameNode – http://localhost:9870/



Start ResourceManager daemon and NodeManager daemon:

```
$ sbin/start-yarn.sh
```

Browse the web interface for the ResourceManager; by default it is available at:

ResourceManager – http://localhost:8088/