# How to Setup / Install an Apache Spark 3.1.1 Cluster on Ubuntu

Published on April 6, 2022



**Dr. Virendra Kumar Shrivastava**
Professor (Computer Science & Engineering) ||
Researcher || Blogger || Mentor || Alliance...
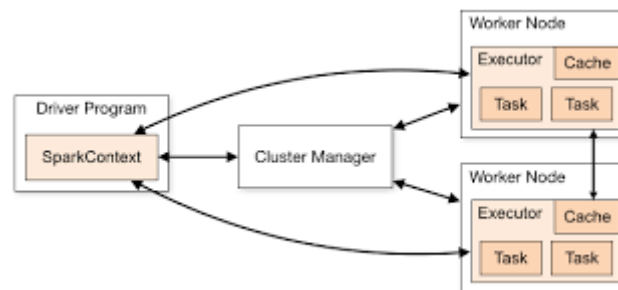
**49 articles**   ➕ Follow

Hey guys,

In this article, I will explain how to setup **Apache Spark** 3.1.1 on a multi-node cluster which includes installing spark master and workers. I will provide step by step instructions to setup spark on Ubuntu 16.04. So, if you are you are looking to make your hand dirty on Apache Spark cluster, this article can be a steppingstone for you.

**What is Apache Spark?**

It is an open-source and distributed processing system used for big data workloads. Spark is a fast, general engine and powerful engine for big data processing. Apache Spark follows a master/worker architecture (two main daemons) and a cluster manager.

- Master Daemon (Master/Driver Process)

- Worker Daemon (Slave Process)

- Cluster Manager



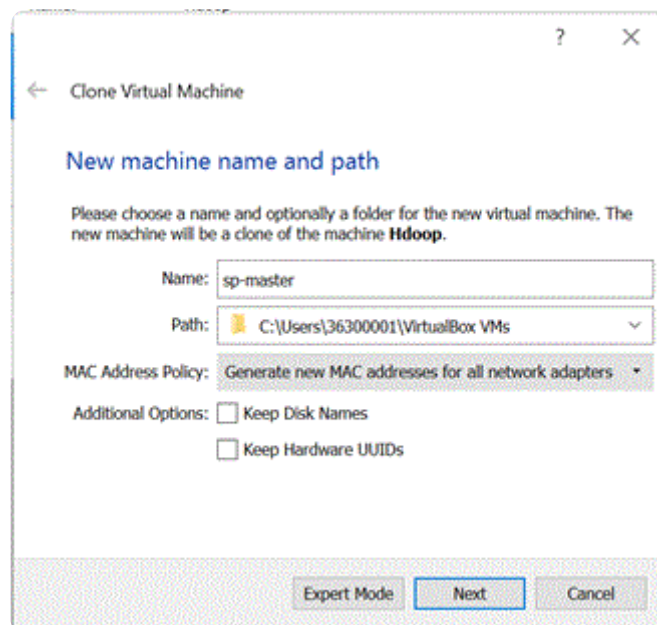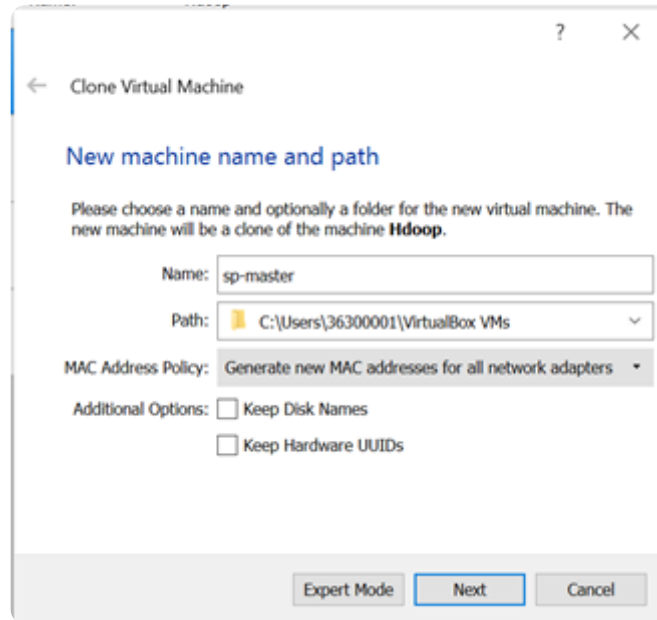An Apache Spark cluster used to have a Master and one or more than one Workers /Slaves.

## Pre-requisite:

Ubuntu 16.04 or higher installed on a virtual machine.

## Steps for installation of Apache Spark 3.1.1 Cluster on Hadoop 3.2

## Step 1.

for all network adapters" in MAC Address Policy. And also choose the option "Full Clone" in clone type.





## Step 2.

Go to settings option of virtual machines and make following network configuration on Adapter 2.

## Step 3.

You need to set the hostname of each virtual machine. Open the /etc/hostname file and type the name of the machine in it and save. Run the following command on each virtual machine:

$ **sudo nano /etc/hostname**

## Step 4.

To figure out IP address of the virtual machines run the following command:

$ **ip addr**

I run the above-mentioned command on master and workers (slaves). On my system, I found following IP addresses:
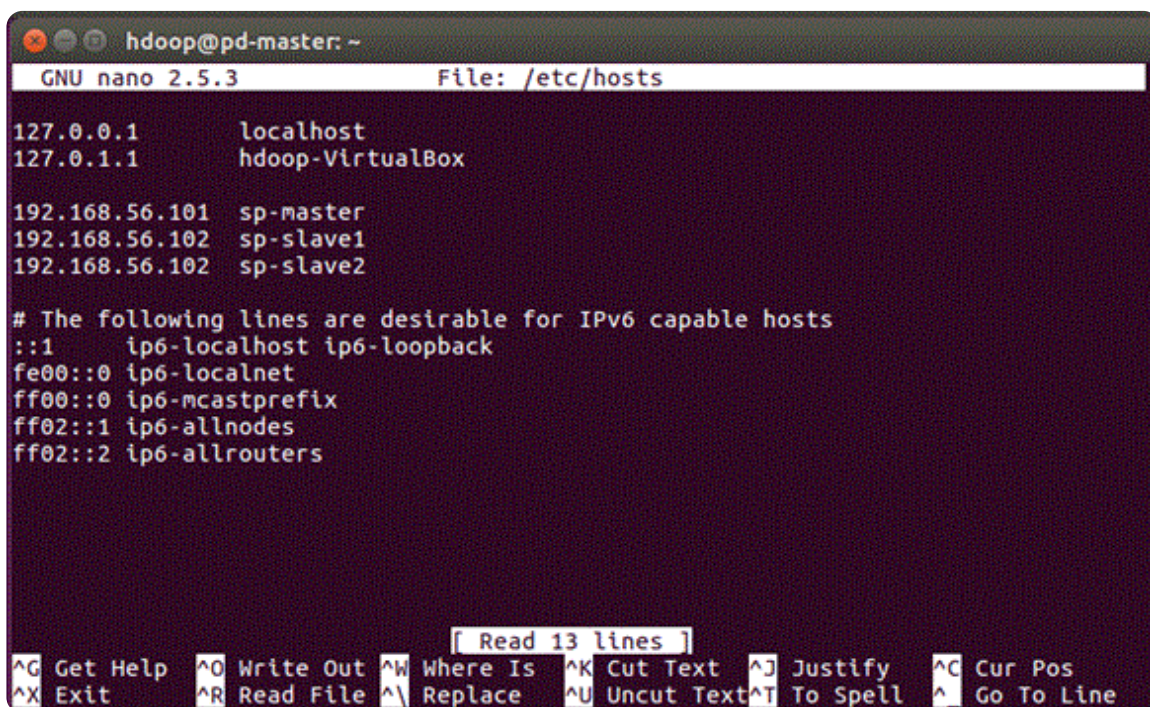
sp-master  192.168.56.101

sp-slave1   192.168.56.102

sp-slave2   192.168.56.103

# Step 5.

In this step, I edit **hosts** file and added IP address and hostname information, saved it, and reboot the machine. Run the following command on all the machines.

$ **sudo nano /etc/hosts**



$ **sudo reboot**

# Step 6.

 Run the following commands on all the Machines (master and workers / slaves).

 $ **sudo apt-get update**

#to check version of java, run the following command.

 $ **java -version**

## Step 7. (On Master and workers)

Install Scala on the all the machines (master and the worker / slaves). Run the following command:

$ **sudo apt-get install scala**

To check the version of Scala, run the following command:

$ **scala -version**

## Step 8. (On Master only)

Now configure Open SSH server-client on master. To configure Open SSH server-client, run the following command:

$ **sudo apt-get install openssh-server openssh-client**

Next step is to generate key pairs. For this purpose, run the following command:

$ **ssh-keygen -t rsa -P ""**

Run the following command to authorize the key:

$ **cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys**

ssh/authorized_keys (all the workers/slaves as well as master). Run the following commands:

$ **ssh-copy-id user@192.168.56.101**

$ **ssh-copy-id user@192.168.56.102**

$ **ssh-copy-id user@192.168.56.103**

Note: user name and IP will be different of your machines. So, use accordingly.

Now it's time to check if everything installed properly. Run the following command on master to connect to the slaves / workers:

$ **ssh 192.168.56.102**

$ **ssh 192.168.56.103**

You can exit from slave machine by type the command:

$ **exit**

## Step 9. (On all the Virtual Machines – Master and workers)

Download the stable version of Apache Spark. I will install spark-3.1.1 with Hadoop-3.2. To download spark-3.1.1 with Hadoop-3.2, run the following command:

Now run the following command to untar the spark tar file:

$ **tar xvf spark-3.1.1-bin-hadoop3.2.tgz**

Run the following command to move the spark files to the spark directory (*/usr/local/bin*):

$ **sudo mv spark-3.1.1-bin-hadoop3.2 /usr/local/spark**

**To set up the environment for Apache Spark, we need to e**dit the *.bashrc* file. Run the following command to edit .bashrc file:

$ **sudo nano ~/.bashrc**

Add the following line to the file and save.

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

The above line sets the location (Path) where the spark software file is located to the PATH variable.

Run the following command to make effective changes in the .bashrc file:

$ **source ~/.bashrc**

## Step 10. (On Master only)

Traverse to the spark/conf folder and make a copy of the spark-env.sh.template file as a spark-env.sh

$ **cd /usr/local/spark/conf**

$ **cp spark-env.sh.template spark-env.sh**

Now, to edit the *spark-env.sh* configuration file, run the following command:

$ **sudo nano *spark-env.sh***

Add the following parameters (line of code) at the end of file, save and exit.

export SPARK_MASTER_HOST='<Master-IP>'export JAVA_HOME=<Path_of_JAVA_installation>

Note: In my system, MASTER IP is 192.168.5.101 and JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64. You need to use these parameters as per your system.

**Add Workers or Slaves**

Now, to edit the configuration file *usr/local/spark/conf/slaves, run the following command on master:*

*$* **sudo nano /usr/local/spark/conf/slaves**

And add the master and workers/slaves name (given below) in the above-mentioned file, save and exit.

sh-slave1

sh-slave2

## Step 11.

**Now, our Apache Spark Cluster is ready.** To start the Apache Spark cluster, run the following command on master:

$ **cd /usr/local/spark**

$ **./sbin/start-all.sh**

Now, to check the services started by spark, run the following command:

$ **jps**



## Step 12.

In case of my system, my master IP is 192.168.56.101.

 you can see in the above snippet; we have two alive Workers are running.

## Step 13.

To stop services of Apache Spark cluster, run the following command:

**$ ./sbin/stop-all.sh**