



Upgrade

Open in app



Published in CodeX

This is your **last** free member-only story this month. [Upgrade for unlimited access.](#)



Andrew Zhu

Follow

May 11, 2021 · 3 min read ★ · Listen



# Set up a local Spark cluster step by step in 10 minutes

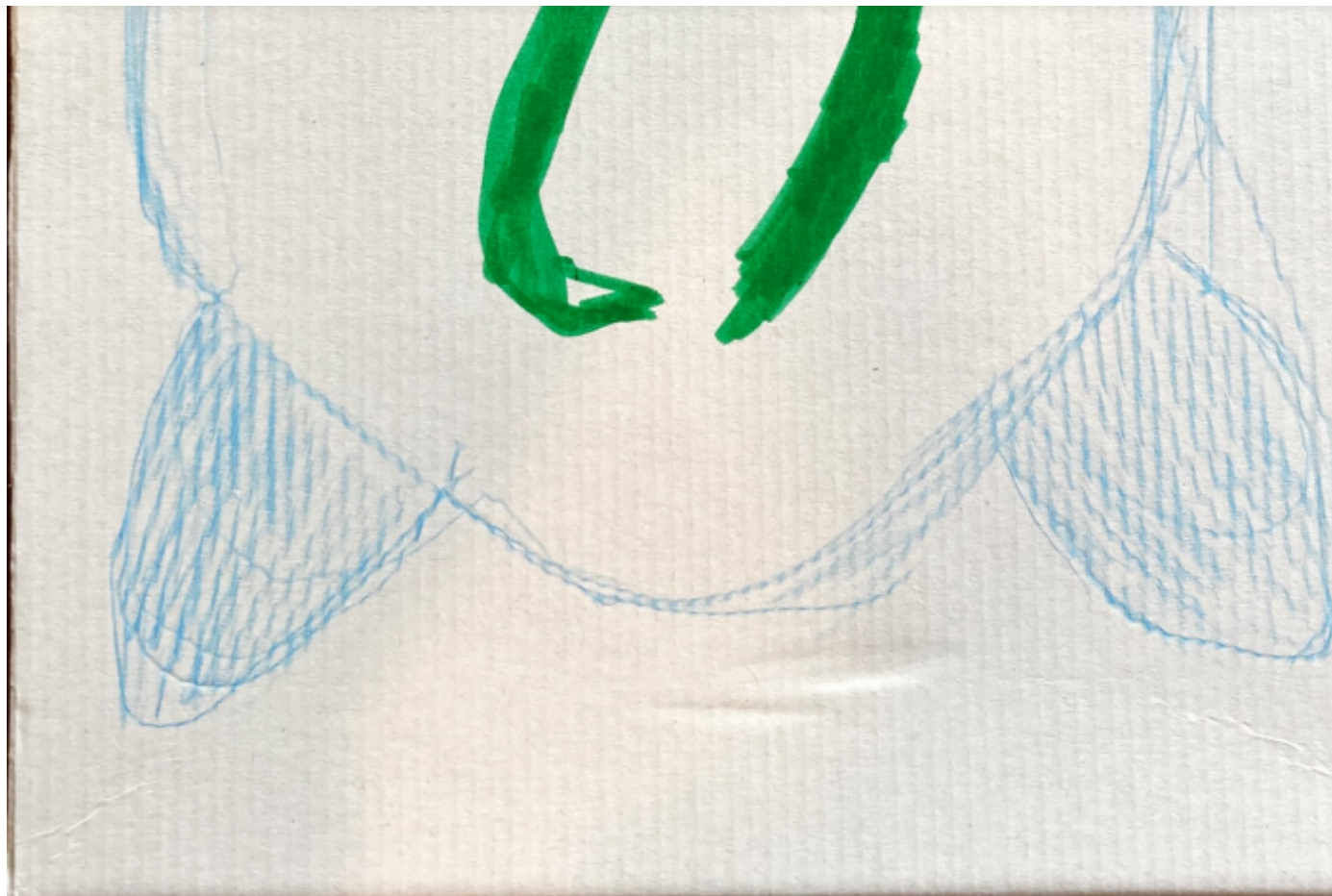
Set up a local Spark cluster with one master node and one worker node in Ubuntu from scratch completely, and for free.





Upgrade

Open in app



Turtle with 4 legs by Charles Zhu, my 6 yo son

This is an action list to install the open-sourced Spark master(or driver) and worker in local Ubuntu completely for free. (in contrast to Databricks for \$\$\$)

The following setup runs in a home intranet. On one Linux(Ubuntu) physical machine(Jetson Nano) and one WSL2(Ubuntu) inside of Windows 10.

## Step 1. Prepare environment

Make sure you have Java installed

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

Check if you get Java installed





Upgrade

Open in app

If you are going to use PySpark, go get Python installed

```
sudo apt install python3
```

Check if you get Python installed

```
python3 --version
```

## Step 2. Download and install Spark in the Driver machine

From the [Spark download](#) page, select your version, I select the newest. (in any directory)

```
curl -O https://apache.claz.org/spark/spark-3.1.1/spark-3.1.1-bin-hadoop3.2.tgz
```

Unpack it.

```
tar xvf spark-3.1.1-bin-hadoop3.2.tgz
```

Now, you should see a new folder `spark-3.1.1-bin-hadoop3.2`. Move this folder to

`/opt/spark`.

```
sudo mv spark-3.1.1-bin-hadoop3.2/ /opt/spark
```

To run Spark related bash script from anywhere, add related PATH to `~/.bashrc`





Upgrade

Open in app

Press `i` to shift to write mode, and append the following to the end of the file.

```
export SPARK_HOME=/opt/spark
export PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin
```

Press `esc` and `:`, then type `wq` to save and quit vim.

Run the following command to activate the changes.

```
source ~/.bashrc
```

### Step 3. Configure the master node, give IP address instead of the computer name

Navigate to the folder that holds configuration files.

```
cd /opt/spark/conf
```

Add a new file named `spark-env.sh`

```
sudo vim spark-env.sh
```

Add one variable

```
SPARK_MASTER_HOST=192.168.0.1
```





Upgrade

Open in app

```
start-master.sh
```

To stop it, run the following script. This will be useful for troubleshooting if you don't get the worker node connected, you will need to stop and start the master node.

```
stop-master.sh
```

The `start-master.sh` command will also start a spark node web page to show basic information include both and later connected worker nodes. do remember use 8080 port by default.

```
http://<master_ip_address>:8080
```

**Spark Master at spark://192.168.0.100:7077****URL:** spark://192.168.0.100:7077**Alive Workers:** 0**Cores in use:** 0 Total, 0 Used**Memory in use:** 0.0 B Total, 0.0 B Used**Resources in use:****Applications:** 0 Running, 0 Completed**Drivers:** 0 Running, 0 Completed**Status:** ALIVE

## Step 4. Setup Spark worker node in another Linux(Ubuntu) machine

Go open another Linux(Ubuntu) machine and repeat **step 2**. No need to take **Step 3** in the worker node.





Upgrade

Open in app

```
start-worker.sh spark://192.168.0.123:7077
```

Replace IP address 192.168.0.123 with yours.

Now, refresh the master UI web, you should see Alive Workers:1

**Spark Master at spark://192.168.0.123:7077**

URL: spark://192.168.0.123:7077

**Alive Workers: 1** ✓

Cores in use: 12 Total, 0 Used

Memory in use: 23.9 GiB Total, 0.0 B Used

Resources in use:

Applications: 0 Running, 0 Completed

Drivers: 0 Running, 0 Completed

Status: ALIVE

If you see this, you have a good Spark cluster full under your control.

## Connect Spark from remote PySpark

You can also verify the Spark cluster by starting a connection from PySpark.

```
from pyspark.sql import SparkSession
spark = SparkSession.builder.\
    master("spark://192.168.0.123:7077").getOrCreate()
print("spark session created")
```

## Links that are helpful

- [How to SSH into WSL2 on Windows 10 from an external machine](#) So that I don't "really" need to shift machines, I sit in one chair and get everything done.





[Upgrade](#)[Open in app](#)

- [spark worker not connecting to master](#)
- [Install Apache Spark on Ubuntu 20.04/18.04 & Debian 10/9](#)

---

## Sign up for CrunchX

By CodeX

A weekly newsletter on what's going on around the tech and programming space [Take a look.](#)

[Get this newsletter](#)

Emails will be sent to satvikj5@gmail.com.  
[Not you?](#)

