# Hadoop Cluster Setup

Sorayut Glomglome

# Content

1. Create Instances from image

2. Configure Hadoop cluster

3. Run MapReduce

# Create 2 new instances from Hadoop image

# Choose m5.large

# Create 2 new instances from Hadoop image

# Name new instances as Slave 1 & Slave 2

# Private IPs

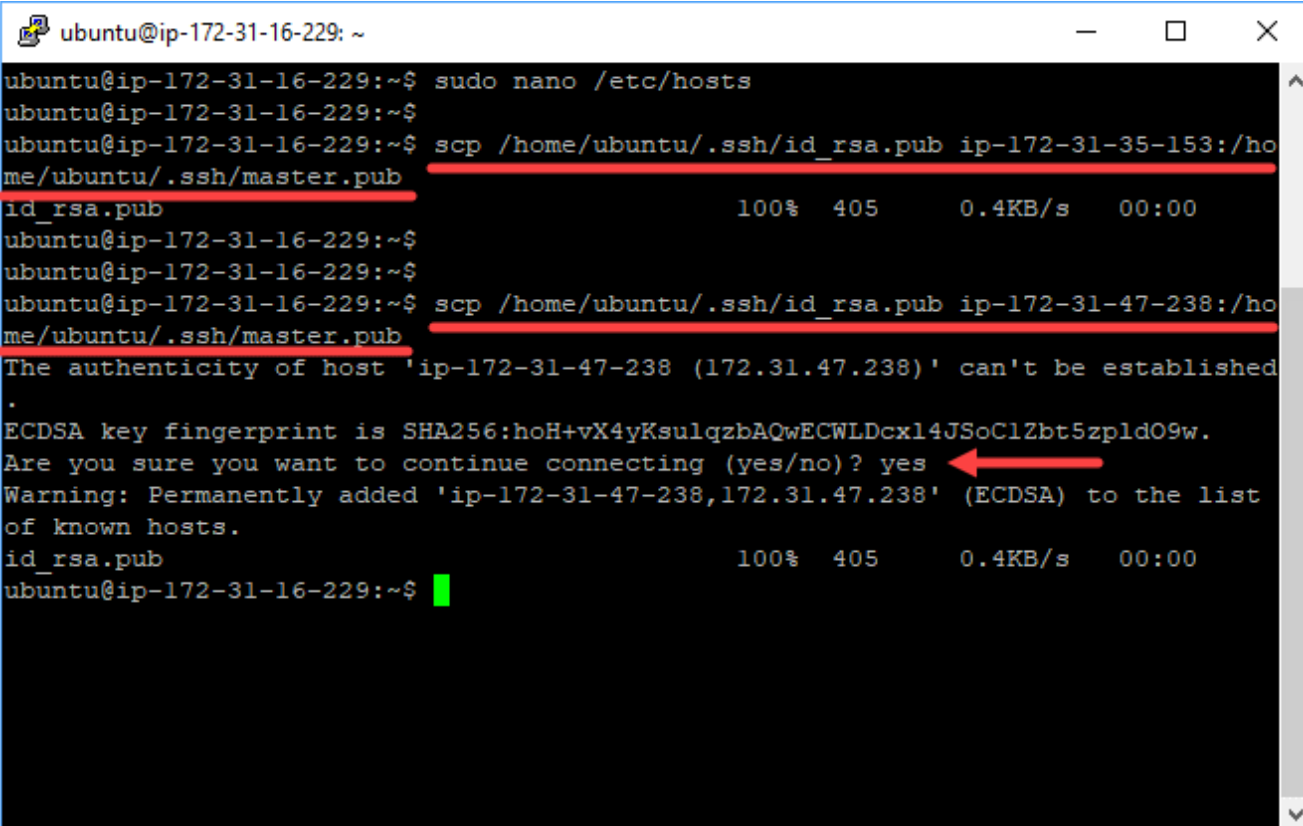| Name | Private IP | Private DNS |
|---|---|---|
| BDA-Hadoop Master | 172.31.16.229 | ip-172-31-16-229.us-west-2.compute.internal |
| BDA-Hadoop Slave 1 | 172.31.35.153 | ip-172-31-35-153.us-west-2.compute.internal |
| BDA-Hadoop Slave 2 | 172.31.47.238 | ip-172-31-47-238.us-west-2.compute.internal |

# At master node : Edit host file

At master node : Copy key file to Slave 1 & Slave 2
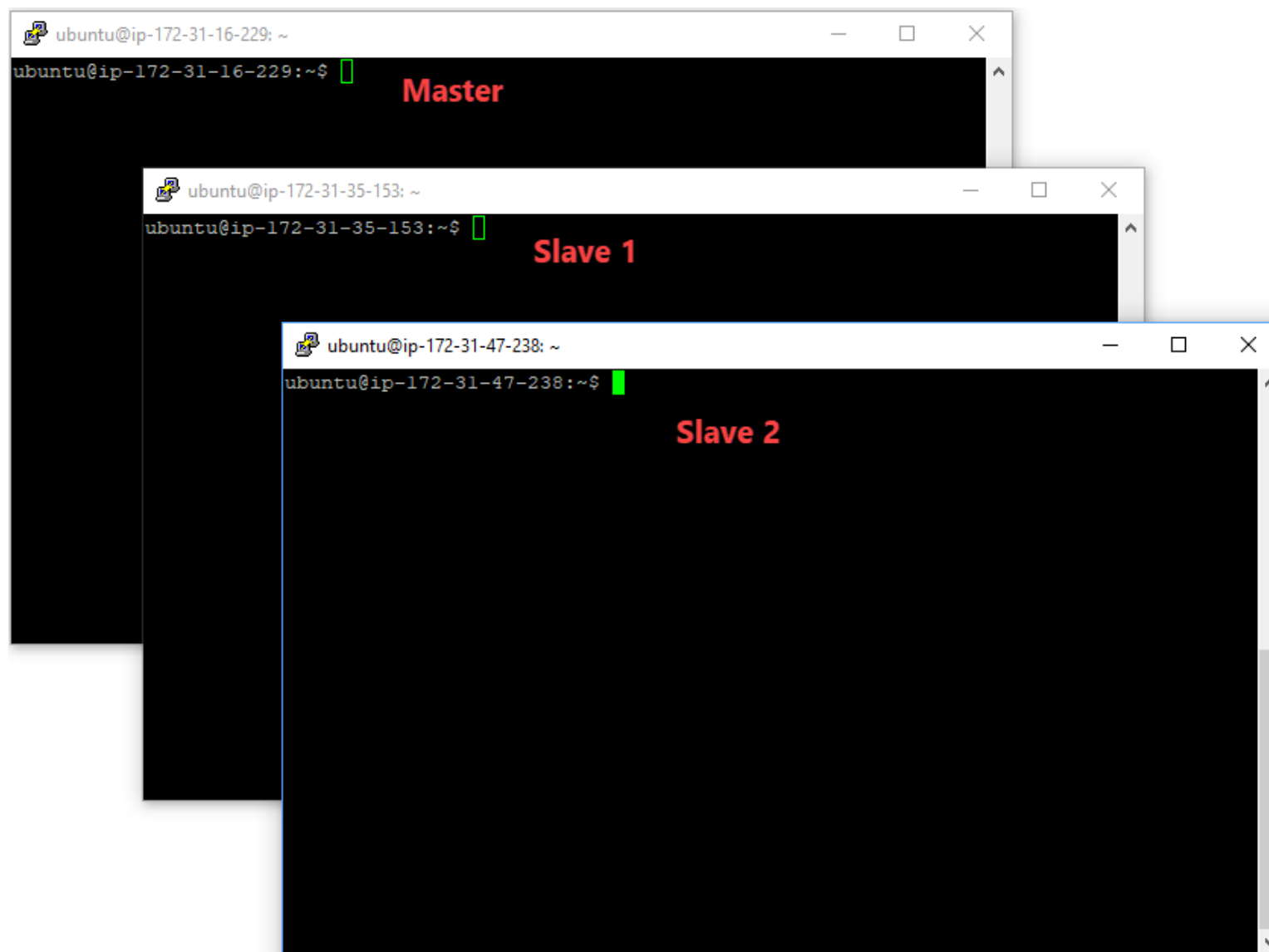
`$scp /home/ubuntu/.ssh/id_rsa.pub ip-172-31-35-153:/home/ubuntu/.ssh/master.pub`

`$scp /home/ubuntu/.ssh/id_rsa.pub ip-172-31-47-238:/home/ubuntu/.ssh/master.pub`

# Arrange SSH terminals

# At Slave1 & Slave2 : Append new key to key file

`$cat /home/ubuntu/.ssh/master.pub >> /home/ubuntu/.ssh/authorized_keys`

At Master : Test ssh to Slave1 & Slave2

```
$ssh ip-172-31-35-153
$exit
$ssh ip-172-31-47-238
$exit
```

# At Master : Test ssh to Slave1

# At Master : Test ssh to Slave2

# At Master : Add Slave1 & Slave2 to Hadoop slave file

`$nano /usr/local/hadoop/etc/hadoop/slaves`

# At Master : Add Slave1 & Slave2 to Hadoop slave file

**Add slave private DNS**
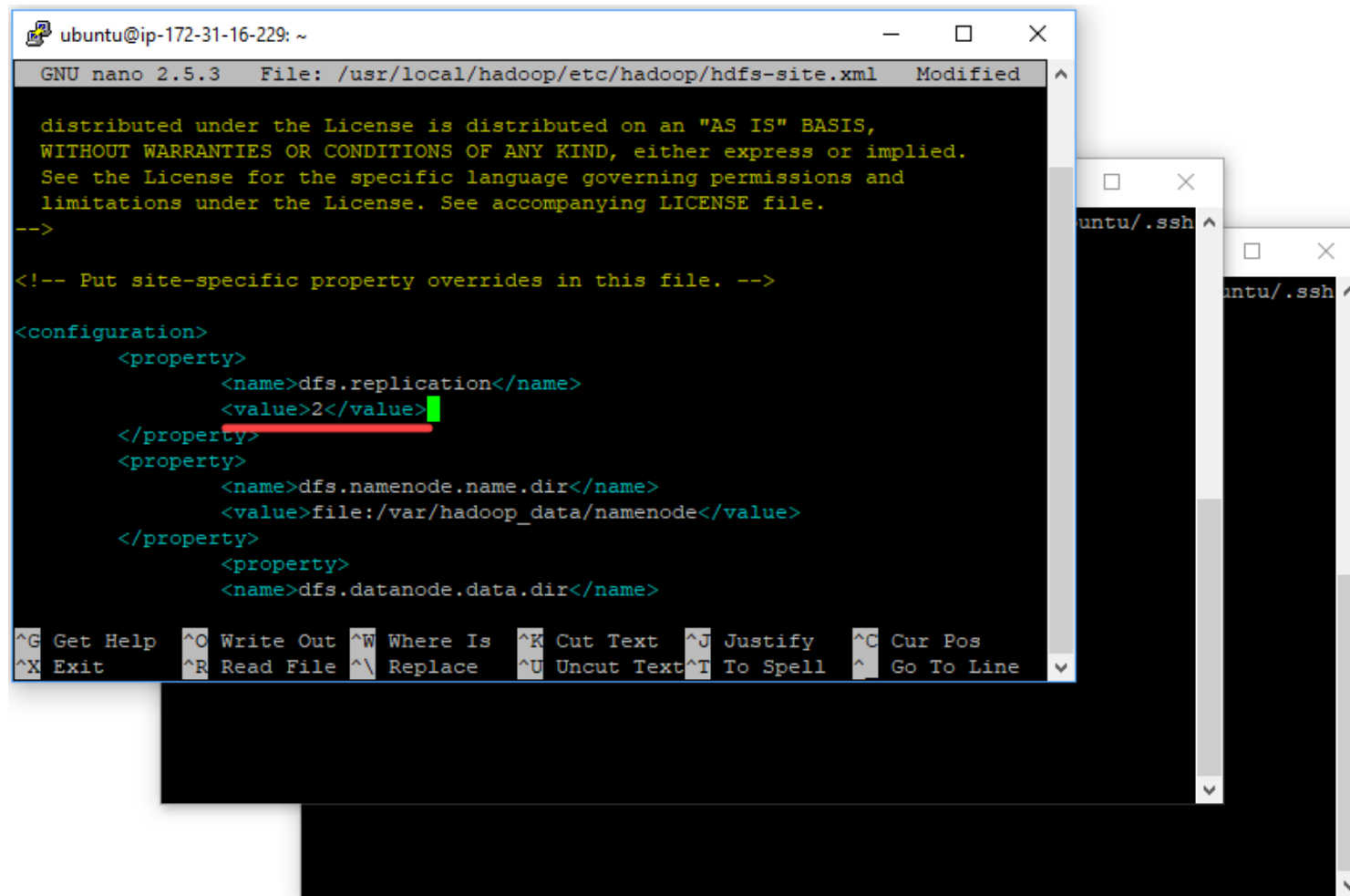
# At Master : Edit hdfs-site.xml

**$nano /usr/local/hadoop/etc/hadoop/hdfs-site.xml**

# At Master : Edit replication to 2

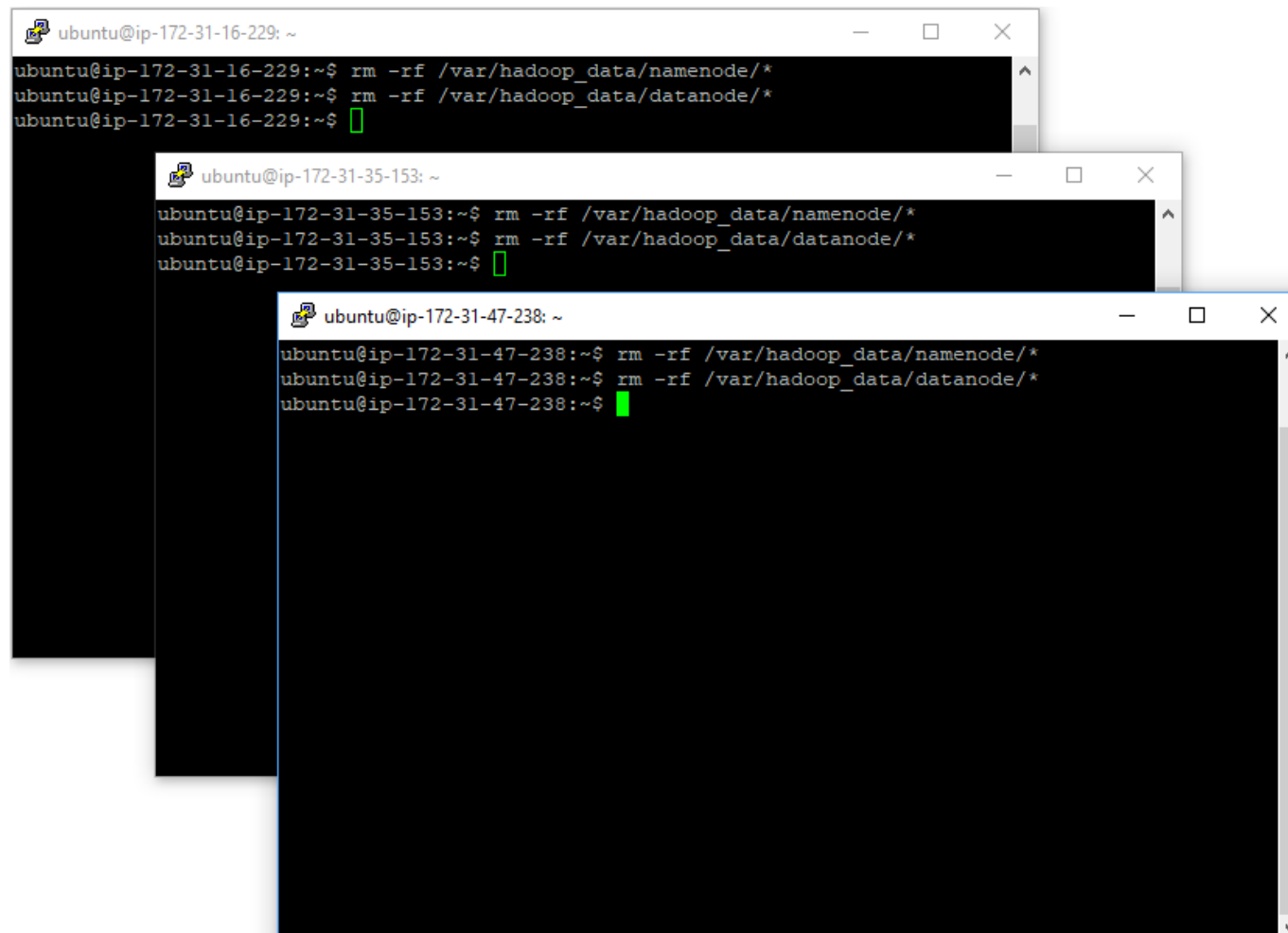# At all nodes : Remove directories of namenode and datanode

```
$rm -rf /var/hadoop_data/namenode/*
$rm -rf /var/hadoop_data/datanode/*
```

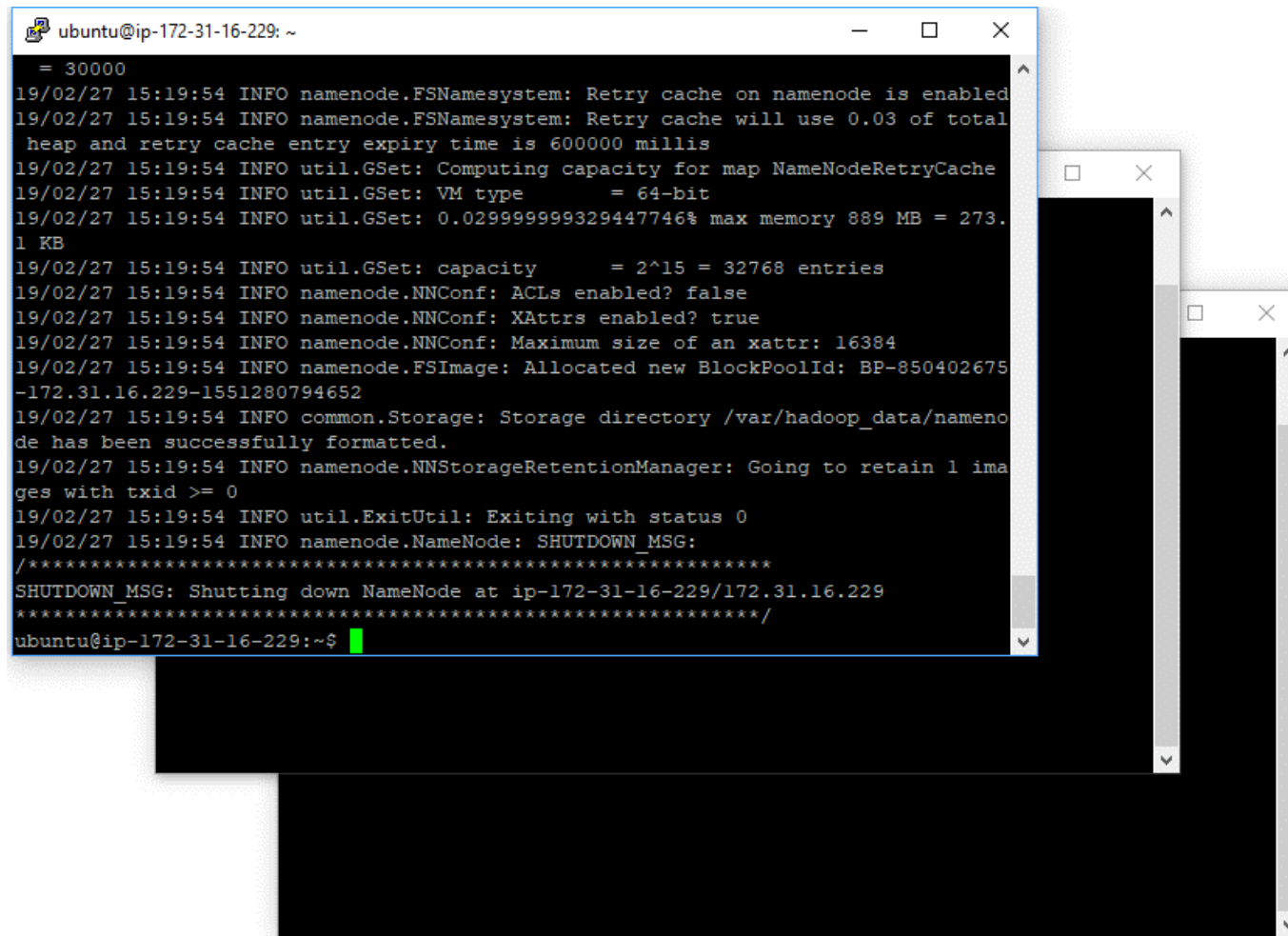# At Master: Format namenode

**`$hdfs namenode -format`**

# At Master: Execute start-dfs.sh

**`$start-dfs.sh`**

At all nodes: Use `jps` to check result, should see <u>NameNode</u> started on Master and <u>DataNode</u> started on both Slave1&Slave2

# At Master: Execute start-yarn.sh

## $start-yarn.sh

# At all nodes: Use `jps` to check result, should see <u>NameManager</u> started on both Slave1&Slave2

# Namenode Information:

`http://`<span style="color:red">`52.26.15.54`</span>`:50070`

| Cluster ID: | CID-3cc3d898-c601-4111-a36e-21754d6cf52d |
|---|---|
| Block Pool ID: | BP-850402675-172.31.16.229-1551280794652 |

## Summary

Security is off.

Safemode is off.

4 files and directories, 1 blocks = 5 total filesystem object(s).

Heap Memory used 122.49 MB of 207.5 MB Heap Memory. Max Heap Memory is 889 MB.

Non Heap Memory used 39.93 MB of 40.69 MB Commited Non Heap Memory. Max Non Heap Memory is -1 B.

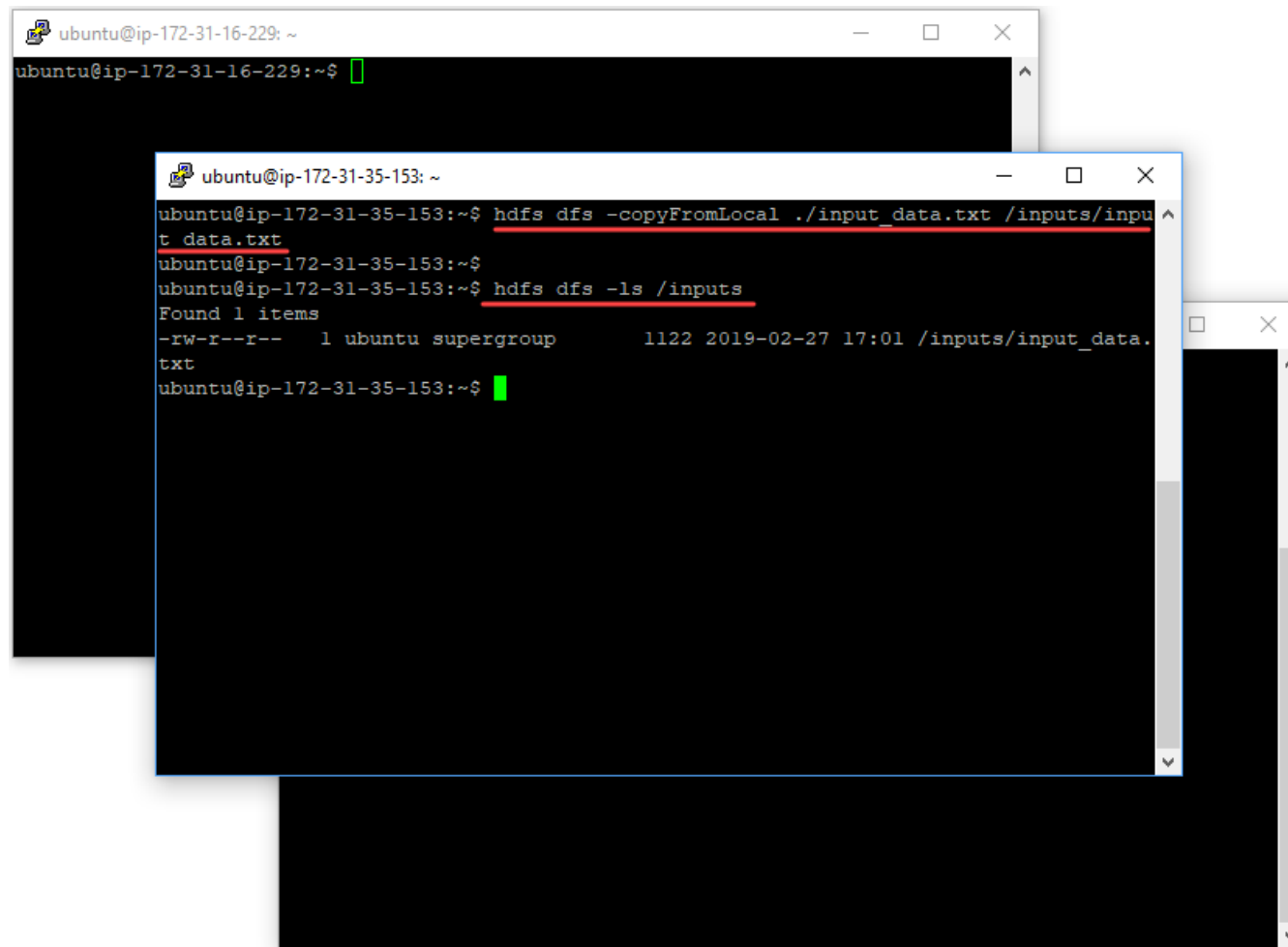| | |
|---|---|
| Configured Capacity: | 38.65 GB |
| DFS Used: | 48 KB |
| Non DFS Used: | 4.74 GB |
| DFS Remaining: | 33.91 GB |
| DFS Used%: | 0% |
| DFS Remaining%: | 87.74% |
| Block Pool Used: | 48 KB |
| Block Pool Used%: | 0% |
| DataNodes usages% (Min/Median/Max/stdDev): | 0.00% / 0.00% / 0.00% / 0.00% |
| Live Nodes | 2 (Decommissioned: 0) |
| Dead Nodes | 0 (Decommissioned: 0) |
| Decommissioning Nodes | 0 |

# Datanode Information

# Log Files

# Next Steps

- Import data to Hadoop cluster

- Execute MapReduce

- Compare result

- Stop Yarn & Stop DFS

Use the same commands as single node Hadoop

# At Slave1: Import data to Hadoop cluster

```
$hdfs dfs -copyFromLocal ./input_data.txt /inputs/input_data.txt
$hdfs dfs -ls /inputs
```
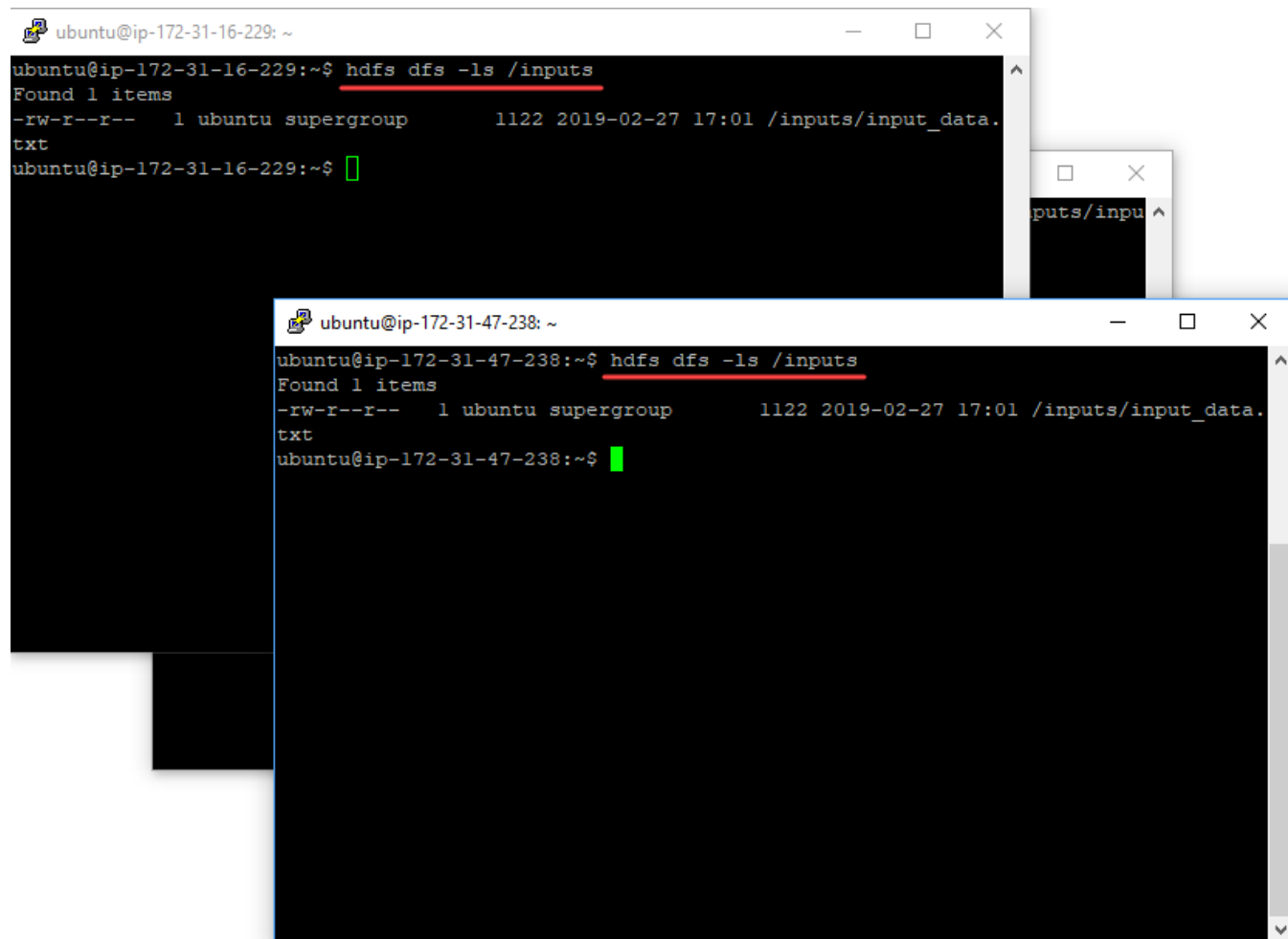
# At Master & Slave2: Check imported file

`$hdfs dfs –ls /inputs`

# At Slave1 & Slave2: Use `jps` to see Application Master and Yarn Child Container

# Don't Forget to SUSPEND/TERMINATE Instance !!!!