

# Running Kmeans Mapreduce code on Amazon AWS

## Pseudo Code

**Input:** Data points D, Number of clusters K

**Step 1:** Copy K and D into memory. Initialize each centroid with 0 as data points.

**Step 2:** Mapper:

Each map task computes the distance of each point with the centroid array. Assign data points to its nearest centroid.

**Step 3:** Mapper Output:

Output the key-value pair with key as centroid and value as the data points array.

**Step 4:** Reducer:

Combine all the values for each key (centroid) and compute the new centroid.

**Step 5:** Reducer Output:

Write the new centroids.

**Step 6:** Repeat steps 1 to 5 until the centroid converges

**Step 7:** Repeat steps 1 to 3 and write the mapper output.

**Output:** Data points with cluster membership.

## How to run the code

### Tools Required

1. [Amazon AWS Account](#)
2. [PuTTY](#) Windows Client (to connect to Amazon EC2 instance)
3. [PuTTYgen](#) (to generate private key – this will be used in putty to connect to EC2 instance)
4. [WinSCP](#) (secure copy)

### 1. Setting up Amazon EC2 Instances

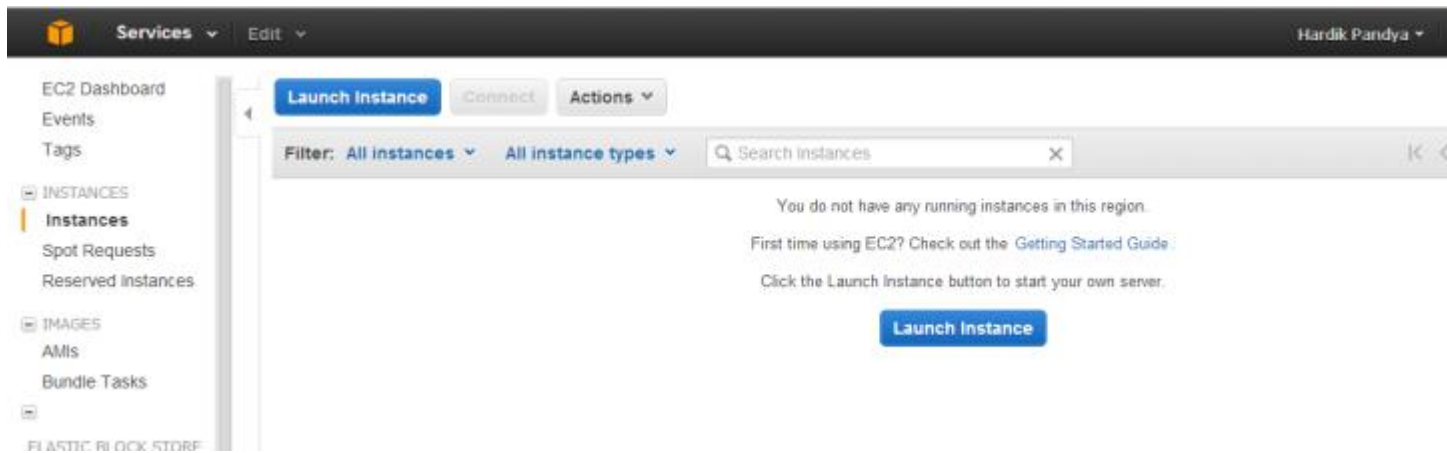
With 4 node clusters and minimum volume size of 8GB there would be an average \$2 of charge per day with all 4 running instances. You can stop the instance anytime to avoid the charge, but you will lose the public IP and host and restarting the instance will create new ones. You can also terminate your Amazon EC2 instance anytime and by default it will delete your instance upon termination, so just be careful what you are doing.

## 1.1 Get Amazon AWS Account

If you do not already have an account, please create a new one. Amazon EC2 comes with eligible free-tier instances.

## 1.2 Launch Instance

Once you have signed up for Amazon account. Login to Amazon Web Services, click on My Account and navigate to Amazon EC2 Console. Click on 'Launch Instance'.

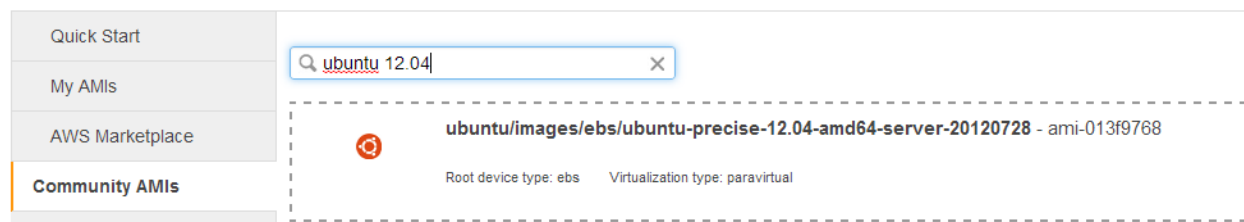


## 1.3 Select AMI

Under the Quick Start column on the left, click on Community AMIs. Search and select Ubuntu Server 12.04 Server 64-bit OS

### Step 1: Choose an Amazon Machine Image (AMI)

An AMI is a template that contains the software configuration (operating system, application server, and applications) required to launch your instance. Yo



## 1.4 Select Instance Type

Select the micro instance and click on 'Next: Configure Instance Details'; on bottom right.

## Step 2: Choose an Instance Type

Amazon EC2 provides a wide selection of instance types optimized to fit different use cases. Instances are virtual servers that can run your applications. [Learn more](#) about instance types and how they can meet your computing needs.

Filter by: All instances Current generation [Show/Hide Columns](#)

Currently selected: t1.micro (up to 2 ECUs, 1 vCPUs, 0.613 GiB memory, EBS only)

	Family	Type	ECUs	vCPUs
<input checked="" type="checkbox"/>	Micro instances Free tier eligible	t1.micro	up to 2	1

### 1.5 Configure Number of Instances

We are setting up 4 node Hadoop cluster, so please enter 4 as number of instances. Please check Amazon EC2 free-tier requirements, you may setup 3 node cluster with < 30GB storage size to avoid any charges. In production environment you want to have SecondNameNode as separate machine. Click on 'Next: Add Storage'; at bottom right.

Configure the instance to suit your requirements. You can launch multiple instances from the same AMI, request Spot Instances to take

**Number of instances**

**Purchasing option**  Request Spot Instances

**Network**  [Create new VPC](#)

**Subnet**  [Create new subnet](#)

**Public IP**  Automatically assign a public IP address to your instances

**IAM role**

**Shutdown behavior**

**Enable termination protection**  Protect against accidental termination

**Monitoring**  Enable CloudWatch detailed monitoring  
[Additional charges apply.](#)

**Tenancy**   
[Additional charges will apply for dedicated tenancy.](#)

### 1.6 Add Storage

Minimum volume size is 8GB. Change it to 20GB (since we can add upto 30GB in free tier) and also change the volume type to "General Purpose (SSD)". Click on 'Next: Tag Instance'; at bottom right.

## Step 4: Add Storage

Your instance will be launched with the following storage device settings. You can attach additional EBS volumes and instance store volumes to your instance, or edit the settings of the root volume. You can also attach additional EBS volumes after launching an instance, but not instance store volumes. [Learn more](#) about storage options in Amazon EC2.

Type <i>i</i>	Device <i>i</i>	Snapshot <i>i</i>	Size (GiB) <i>i</i>	Volume Type <i>i</i>	IOPS <i>i</i>	Delete on Termination <i>i</i>	Encrypted <i>i</i>
Root	/dev/sda1	snap-49967039	20	General Purpose (SSD)	60 / 3000	<input checked="" type="checkbox"/>	Not Encrypted

## 1.7 Instance Description

Give your instance name as “HadoopEC2MultiNodeCluster” and click on ‘Next: Configure Security group’; at bottom right.

## Step 5: Tag Instance

A tag consists of a case-sensitive key-value pair. For example, you could define a tag with key = Name and value = Webserver. [Learn more](#) about tags.

Key <small>(127 characters maximum)</small>	Value <small>(255 characters maximum)</small>
<input type="text" value="Name"/>	<input type="text" value="HadoopEC2MultiNodeCluster"/>
<input type="button" value="Create Tag"/> <small>(Up to 10 tags maximum)</small>	

## 1.8 Define a Security Group

Create a new security group, later on we are going to modify the security group with security rules. Name it ‘HadoopEC2SecurityGroup’. Click ‘Review and Launch’; at bottom right.

## Step 6: Configure Security Group


A security group is a set of firewall rules that control the traffic for your instance. On this page, you can add rules to allow specific traffic to reach your instance. For example, if you want to allow access to the HTTP and HTTPS ports. You can create a new security group or select from an existing one below. [Learn more](#) about Amazon EC2 security groups.

Assign a security group:  Create a new security group  
 Select an existing security group

Security group name:

Description:

Type <i>i</i>	Protocol <i>i</i>	Port Range <i>i</i>
SSH	TCP	22

 **Warning**  
Rules with source of 0.0.0.0/0 allow all IP addresses to access your instance. We recommend setting security group rules to allow access from known IP addresses only.

## 1.9 Launch Instance and Create Security Pair

Review and Launch Instance.

Amazon EC2 uses public–key cryptography to encrypt and decrypt login information. Public–key cryptography uses a public key to encrypt a piece of data, such as a password, then the recipient uses the private key to decrypt the data. The public and private keys are known as a *key pair*.

Create a new keypair and give it a name “hadoopec2cluster” and download the keypair (.pem) file to your local machine. Click Launch Instance


### Select an existing key pair or create a new key pair ✕

A key pair consists of a **public key** that AWS stores, and a **private key file** that you store. Together, they allow you to connect to your instance securely. For Windows AMIs, the private key file is required to obtain the password used to log into your instance. For Linux AMIs, the private key file allows you to securely SSH into your instance.

Create a new key pair ▼

**Key pair name**  
hadoopec2cluster

**Download Key Pair**

 You have to download the **private key file** (\*.pem file) before you can continue. **Store it in a secure and accessible location.** You will not be able to download the file again after it's created.

**Cancel** **Launch Instances**

## 1.10 Launching Instances

Once you click “Launch Instance” 4 instance should be launched with “pending” state

Launch Instance Connect Actions

Filter: All Instances All instance types Search Instances 1 to 4 of

<input type="checkbox"/>	Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks	Alarm Status
	HadoopEC2Cluster	i-a8d5b388	t1.micro	us-east-1b	pending	Initializing	None
	HadoopEC2Cluster	i-a9d5b389	t1.micro	us-east-1b	pending	Initializing	None
	HadoopEC2Cluster	i-aad5b38a	t1.micro	us-east-1b	pending	Initializing	None
	HadoopEC2Cluster	i-abd5b38b	t1.micro	us-east-1b	pending	Initializing	None

Once in “running” state we are now going to rename the instance name as below.

1. HadoopNameNode (Master)
2. HadoopSecondaryNameNode
3. HadoopSlave1 (data node will reside here)
4. HadoopSlave2 (data node will reside here)

You can rename the instance by clicking by hovering on the name and clicking on the pen icon showed next to it. Once renamed click on the tick mark.

Filter: All Instances All instance types Search Instances 1 to 4 of

<input type="checkbox"/>	Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks	Alarm Status
<input type="checkbox"/>	HadoopNameNode	i-a8d5b388	t1.micro	us-east-1b	running	2/2 check...	None
<input type="checkbox"/>	HadoopSecondaryNameNode	i-a9d5b389	t1.micro	us-east-1b	running	2/2 check...	None
<input type="checkbox"/>	HadoopSlave1	i-aad5b38a	t1.micro	us-east-1b	running	2/2 check...	None
<input type="checkbox"/>	HadoopSlave2	i-abd5b38b	t1.micro	us-east-1b	running	2/2 check...	None

Please note down the Instance ID, Public DNS/URL like (ec2-54-209-221-112.compute-1.amazonaws.com) and Public IP for each instance for your reference. We will need it later on to connect from Putty client. Also notice we are using “HadoopEC2SecurityGroup”.

Filter: All instances ▾ All instance types ▾  X 1 to

<input type="checkbox"/>	Name	Instance ID	Instance Type	Availability Zone	Instance State	Status Checks	Alarm Status
<input checked="" type="checkbox"/>	HadoopNameNonde	i-a8d5b388	t1.micro	us-east-1b	running	2/2 check...	None
<input type="checkbox"/>	HadoopSecondaryNameNode	i-a9d5b389	t1.micro	us-east-1b	running	2/2 check...	None
<input type="checkbox"/>	HadoopSlave1	i-aad5b38a	t1.micro	us-east-1b	running	2/2 check...	None
<input type="checkbox"/>	HadoopSlave2	i-abd5b38b	t1.micro	us-east-1b	running	2/2 check...	None

Instance: i-a8d5b388 (HadoopNameNonde) Public DNS: ec2-54-209-221-112.compute-1.amazonaws.com

**Description** | Status Checks | Monitoring | Tags

<b>Instance ID</b>	i-a8d5b388	<b>Public DNS</b>	ec2-54-209-221-112.compute-1.amazonaws.com
<b>Instance state</b>	running	<b>Public IP</b>	54.209.221.112
<b>Instance type</b>	t1.micro	<b>Elastic IP</b>	-
<b>Private DNS</b>	ip-172-31-35-98.ec2.internal	<b>Availability zone</b>	us-east-1b
<b>Private IPs</b>	172.31.35.98	<b>Security groups</b>	HadoopEC2SecurityGroup. view
<b>Secondary private IPs</b>		<b>Scheduled events</b>	No scheduled events
<b>VPC ID</b>	vpc-120e1470	<b>AMI ID</b>	ubuntu-precise-12.04-amd64-se (ami-a73264ce)

You can use the existing group or create a new one. When you create a group with default options it add a rule for SSH at port 22.In order to have TCP and ICMP access we need to add 2 additional security rules. Add 'All TCP', 'All ICMP' and 'SSH (22)' under the inbound rules to "HadoopEC2SecurityGroup". This will allow ping, SSH, and other similar commands among servers and from any other machine on internet. Make sure to "Apply Rule changes" to save your changes.

These protocols and ports are also required to enable communication among cluster servers. As this is a test setup we are allowing access to all for TCP, ICMP and SSH and not bothering about the details of individual server port and security.

Viewing: All Security Groups

	Group ID	Name	VPC ID	Description
<input type="checkbox"/>	sg-0436c861	default	vpc-120e1470	default VPC security group
<input checked="" type="checkbox"/>	sg-0811ef6d	HadoopEC2SecurityGroup	vpc-120e1470	security group for hadoop ec2 security group

**1 Security Group selected**

**Security Group: HadoopEC2SecurityGroup**

Create a new rule: Custom TCP rule

Port range:

(e.g., 80 or 49152-65535)

Source:

(e.g., 192.168.2.0/24, sg-47ad482e, or 1234567890/default)

Your changes have not been applied yet.

ICMP	
Port (Service)	Source
ALL	0.0.0.0/0
TCP	
Port (Service)	Source
22 (SSH)	0.0.0.0/0
0 - 65535	0.0.0.0/0

## 2. Setting up client access to Amazon Instances

Now, let's make sure we can connect to all 4 instances. For that we are going to use Putty client. We are going to setup password-less SSH access among servers to setup the cluster. This allows remote access from Master Server to Slave Servers so Master Server can remotely start the Data Node and Task Tracker services on Slave servers.

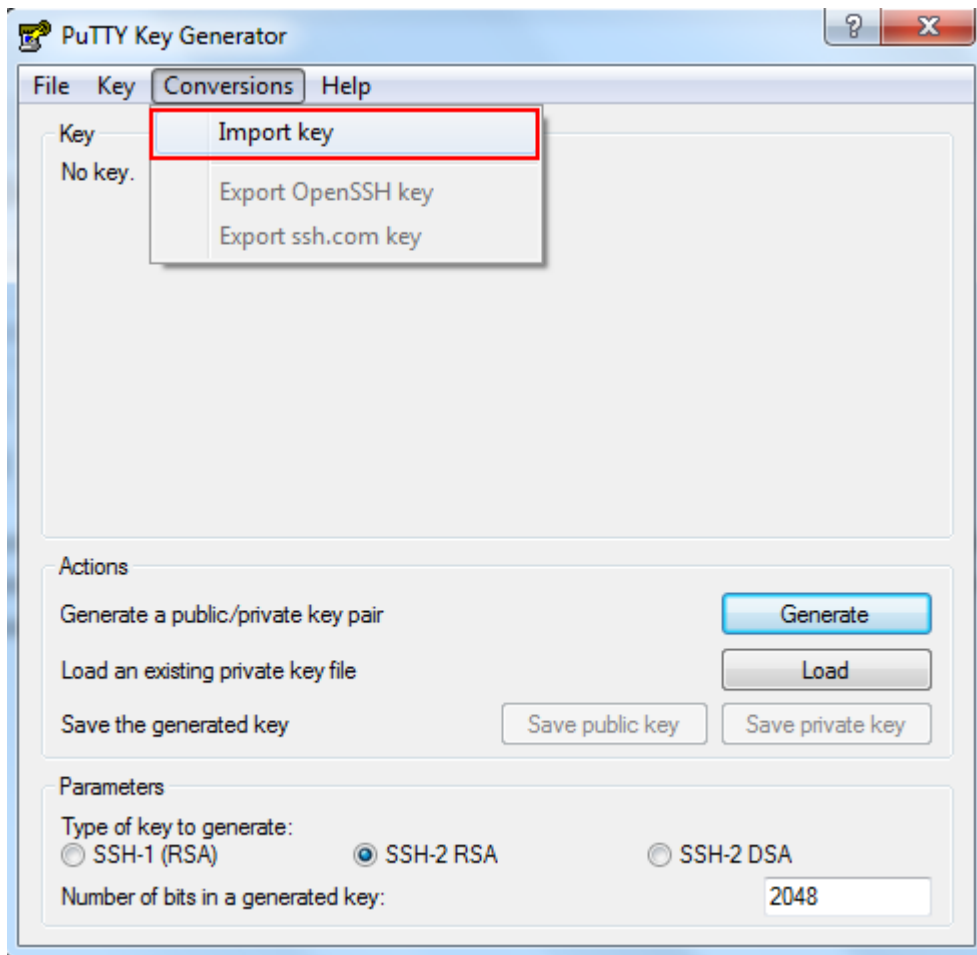
We are going to use downloaded hadoop-ec2-cluster.pem file to generate the private key (.ppk). In order to generate the private key we need Puttygen client. You can download the putty and puttygen and various utilities in zip from [here](#).

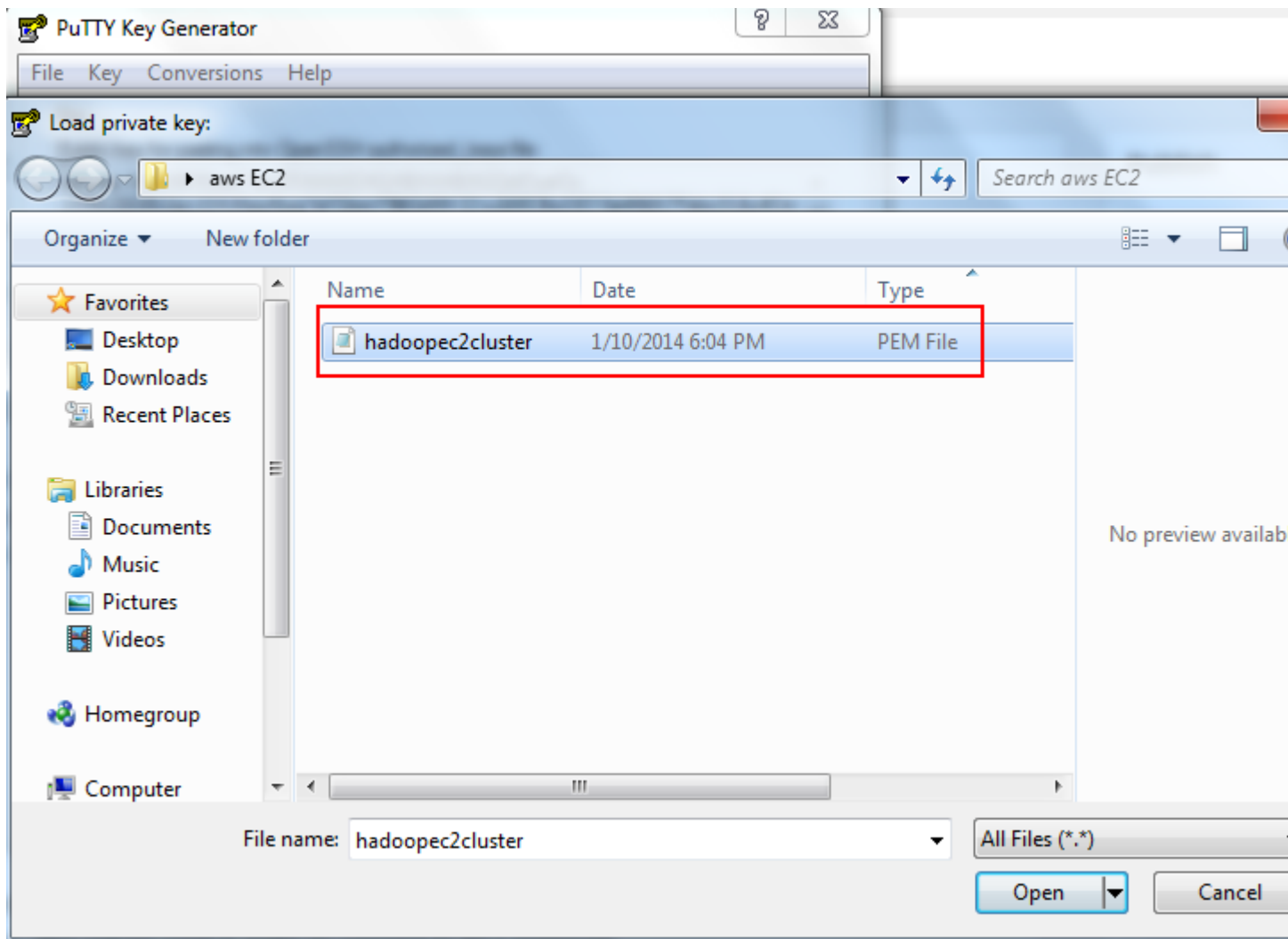
### 2.1 Generating Private Key

Let's launch PUTTYGEN client and import the key pair we created during launch instance step – "hadoop-ec2-cluster.pem"



Navigate to Conversions and "Import Key"



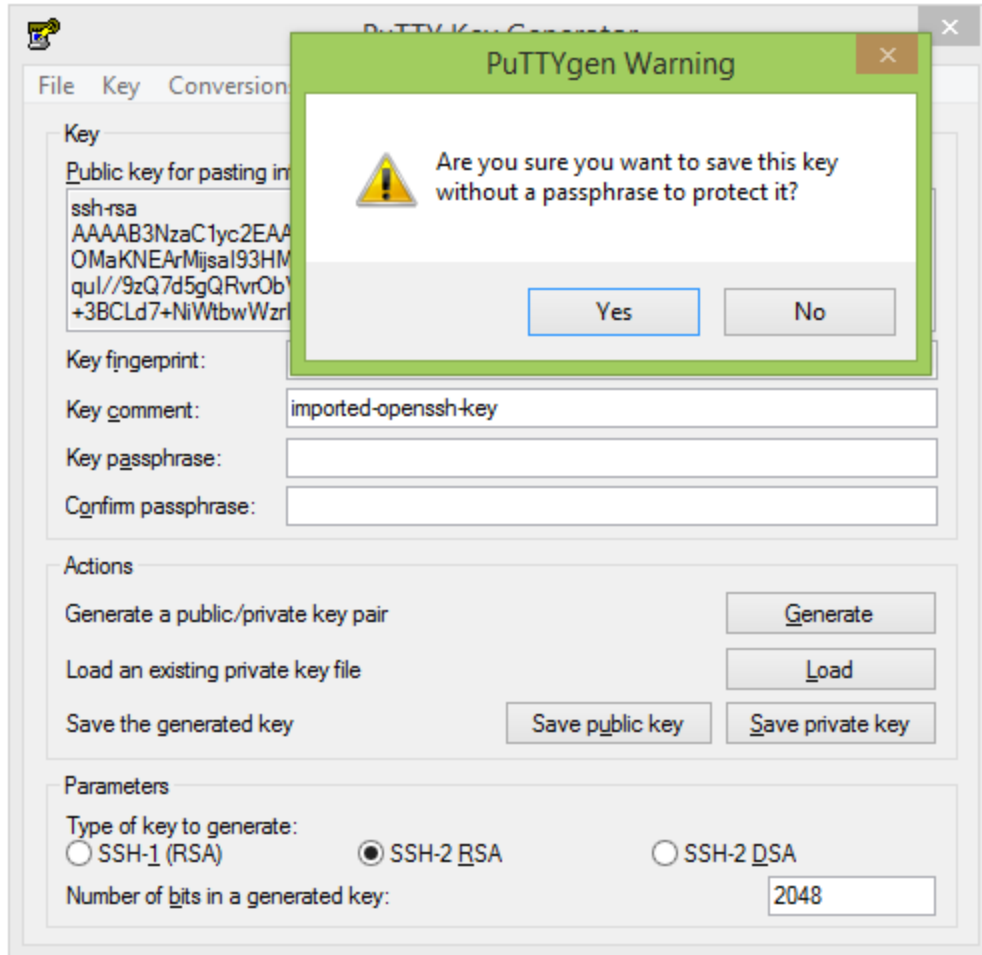


Once you import the key you can enter passphrase to protect your private key or leave the passphrase fields blank to use the private key without any passphrase. But for now leave it **blank**. Passphrase protects the private key from any unauthorized access to servers using your machine and your private key.

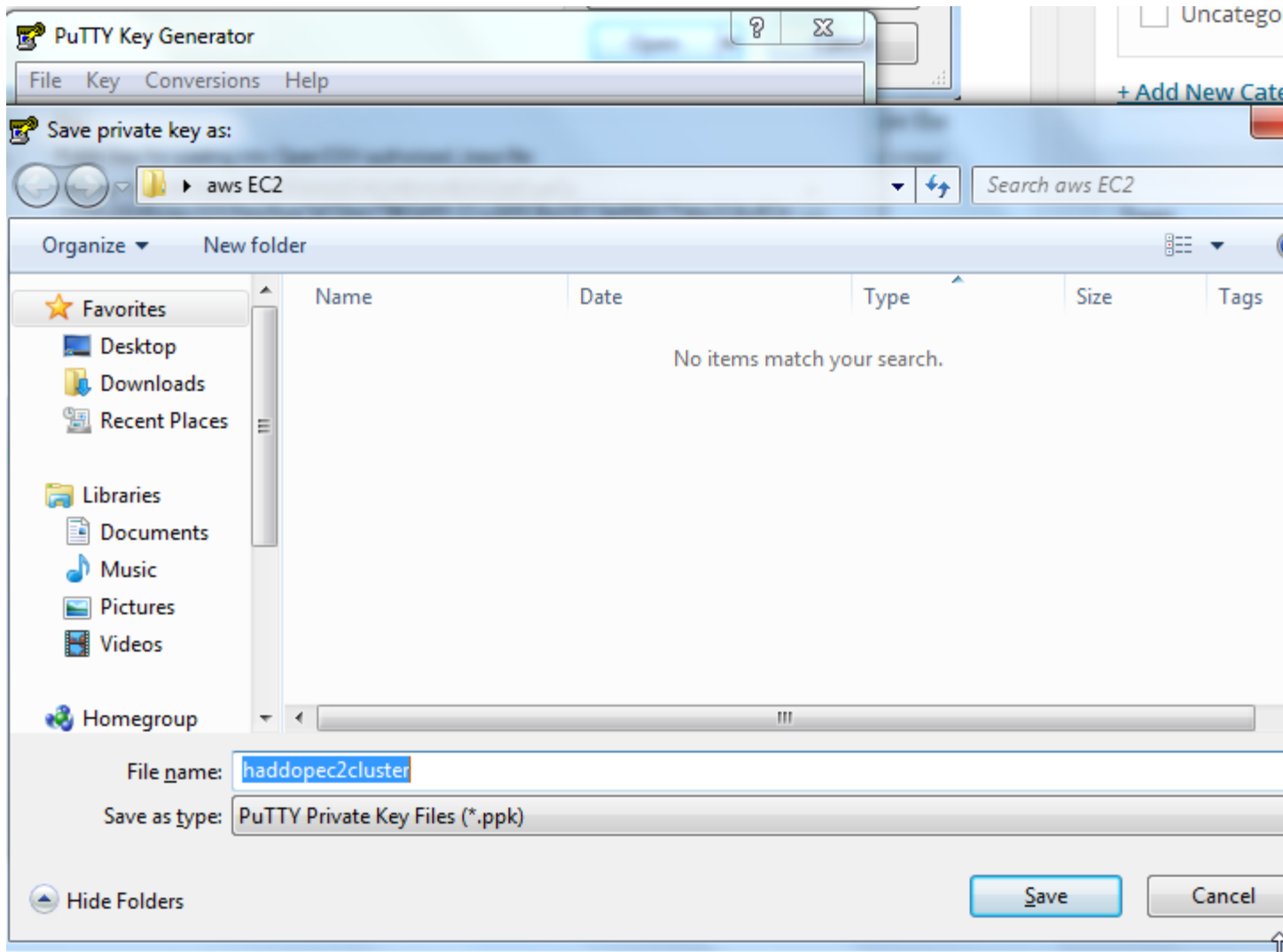
Any access to server using passphrase protected private key will require the user to enter the passphrase to enable the private key enabled access to AWS EC2 server.

## 2.2 Save Private Key

Now save the private key by clicking on "Save Private Key" and click "Yes" as we are going to leave passphrase empty.



Save the .ppk file and give it the same name.

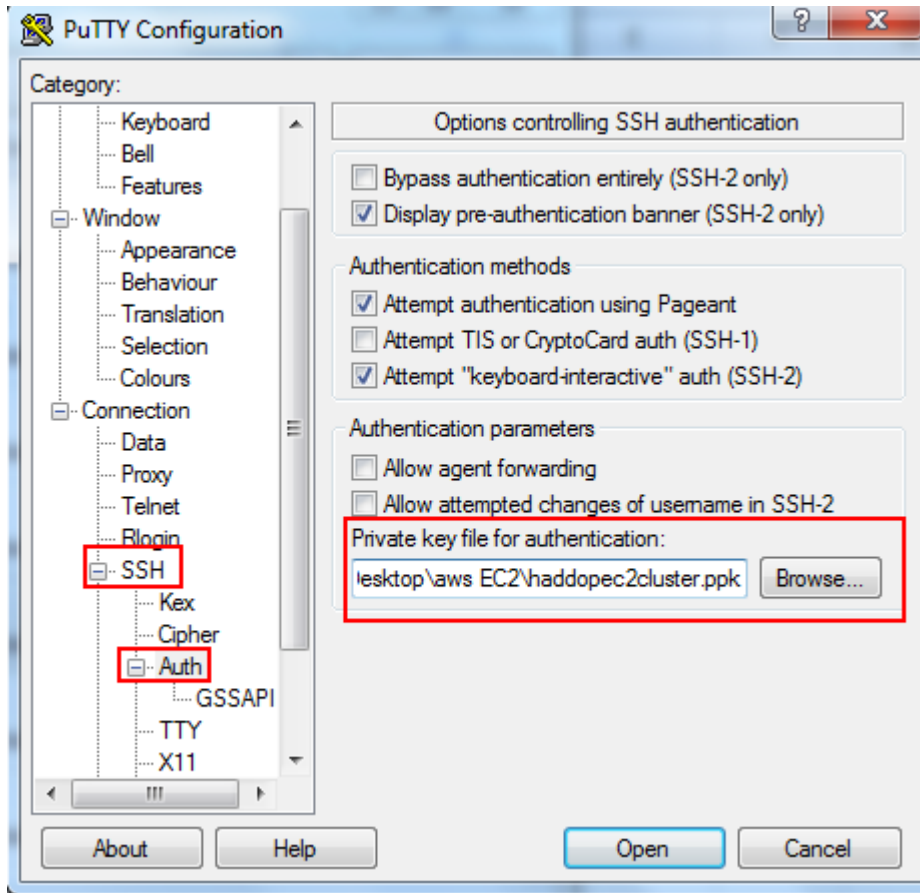


Now we are ready to connect to our Amazon Instance Machine for the first time.

## 2.3 Connect to Amazon Instance

Let's connect to HadoopNameNode first. Launch Putty client, grab the public URL (the DNS ec2-....-amazonaws.com from the console step 1.10), import the .ppk private key that we just created for password-less SSH access. As per [amazon documentation](#), for Ubuntu machines username is "ubuntu"

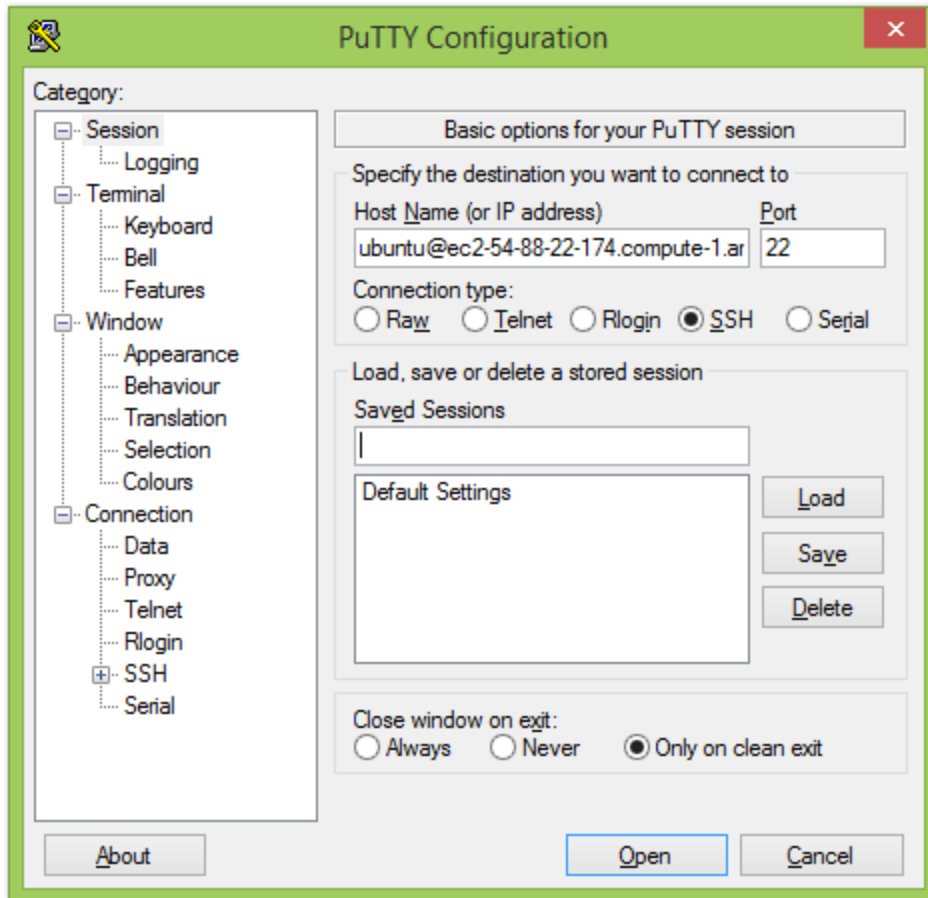
### 2.3.1 Provide private key for authentication



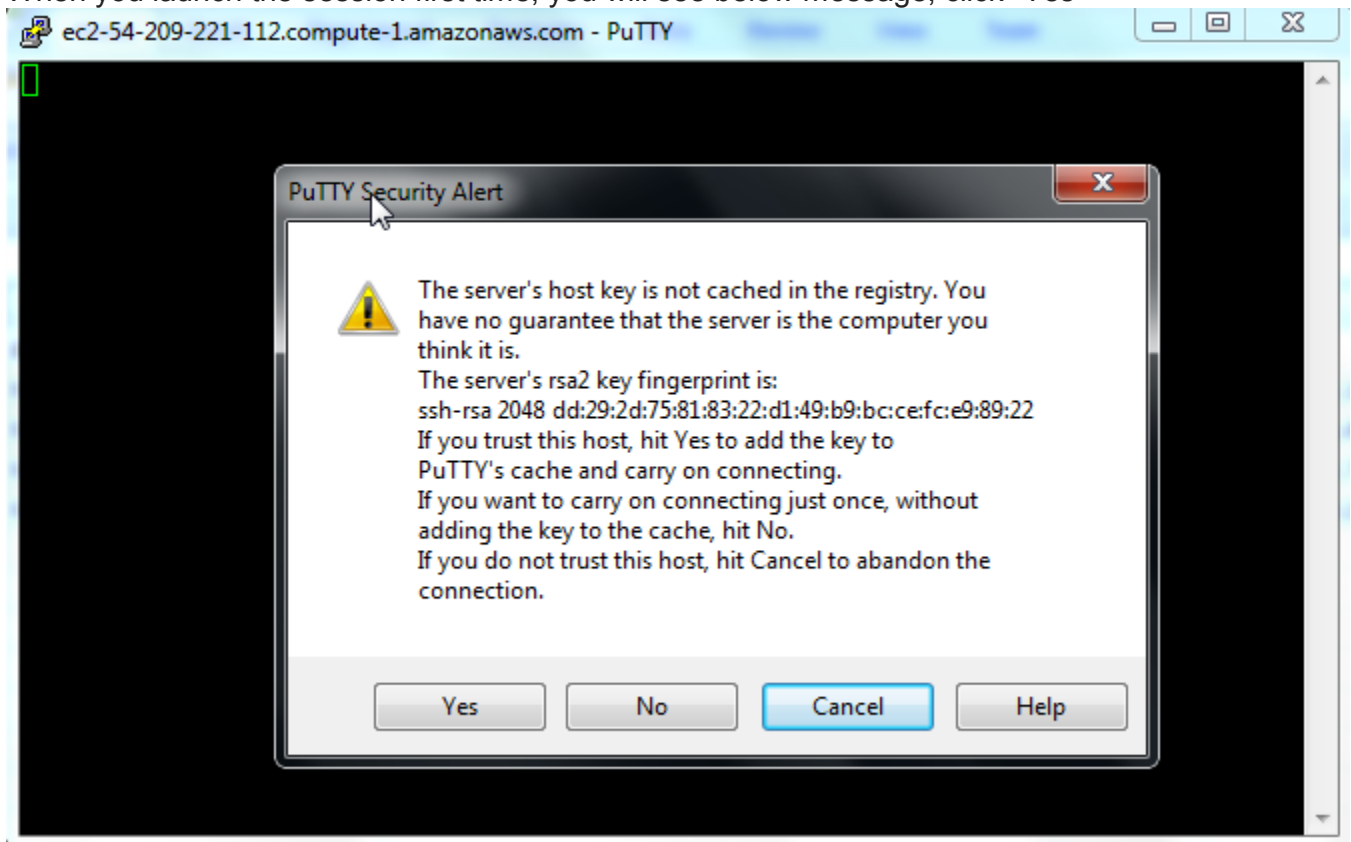
### 2.3.2 Hostname and Port and Connection Type

Host name will be like “**Ubuntu**@ec2-.....compute-1.amazonaws.com”

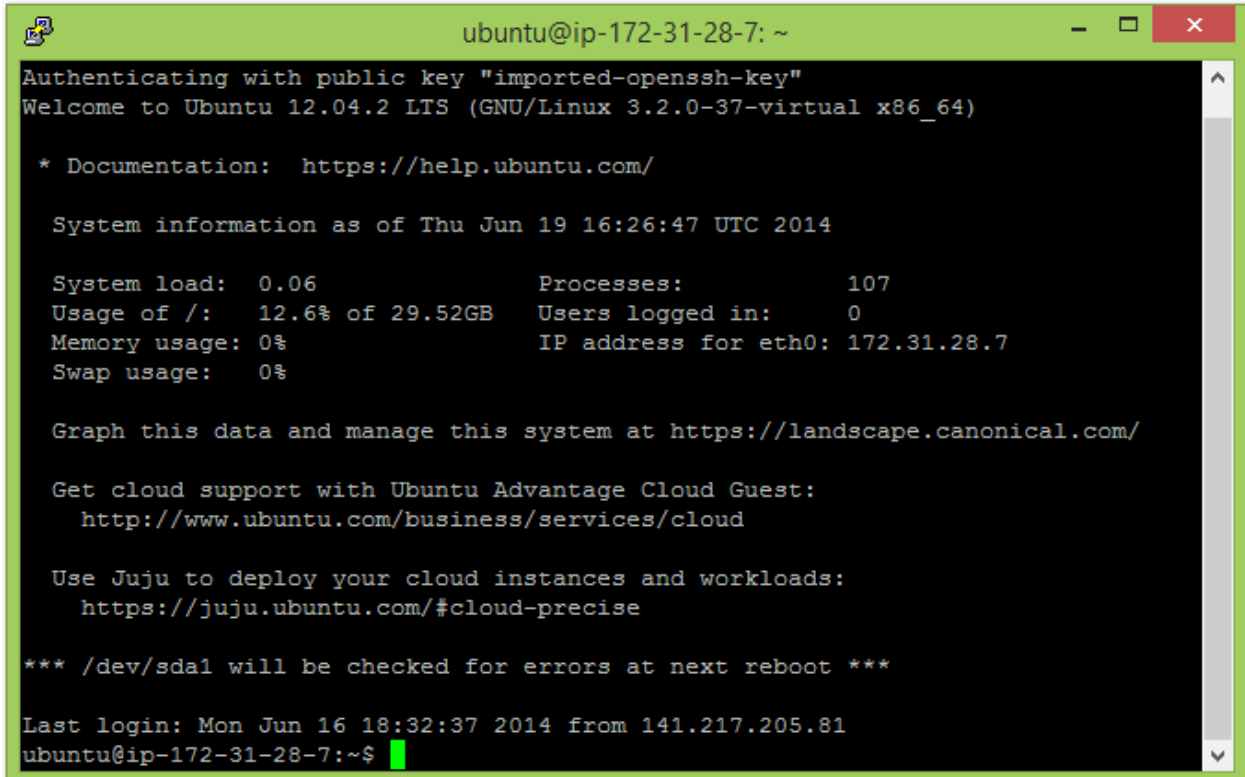
Click “Open” to launch putty session.



When you launch the session first time, you will see below message, click "Yes"



If everything goes well you will be presented welcome message with Unix shell at the end.



```
ubuntu@ip-172-31-28-7: ~
Authenticating with public key "imported-openssh-key"
Welcome to Ubuntu 12.04.2 LTS (GNU/Linux 3.2.0-37-virtual x86_64)

* Documentation:  https://help.ubuntu.com/

System information as of Thu Jun 19 16:26:47 UTC 2014

System load:  0.06          Processes:            107
Usage of /:   12.6% of 29.52GB Users logged in:     0
Memory usage: 0%          IP address for eth0: 172.31.28.7
Swap usage:   0%

Graph this data and manage this system at https://landscape.canonical.com/

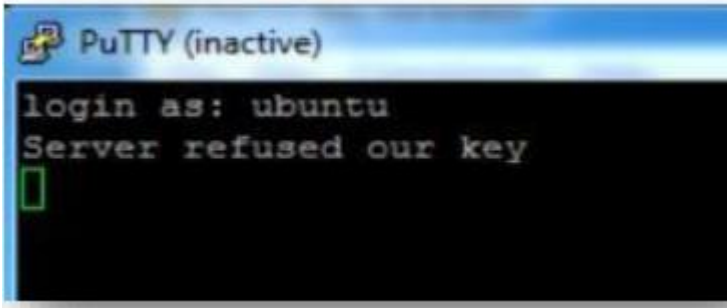
Get cloud support with Ubuntu Advantage Cloud Guest:
http://www.ubuntu.com/business/services/cloud

Use Juju to deploy your cloud instances and workloads:
https://juju.ubuntu.com/#cloud-precise

*** /dev/sda1 will be checked for errors at next reboot ***

Last login: Mon Jun 16 18:32:37 2014 from 141.217.205.81
ubuntu@ip-172-31-28-7:~$
```

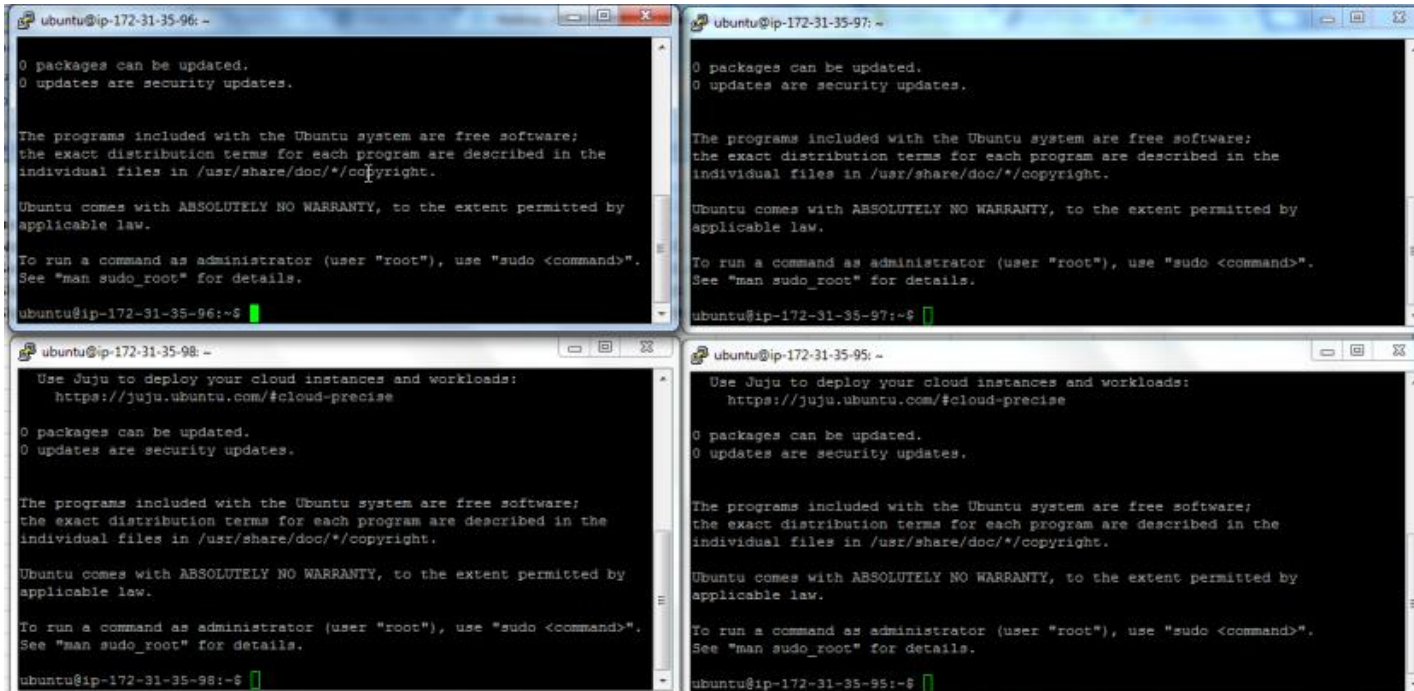
If there is a problem with your key, you may receive below error message. Check manual again and eradicate the error.



```
PuTTY (inactive)
login as: ubuntu
Server refused our key
█
```

Similarly connect to remaining 3 machines HadoopSecondaryNameNode, HadoopSlave1 and HadoopSlave2 respectively to make sure you can connect successfully.





## 2.4 Enable Public Access

Issue `ifconfig` command in the terminal and note down the ip address. Next, we are going to update the hostname with ec2 public URL and finally we are going to update `/etc/hosts` file to map the ec2 public URL with ip address. This will help us to configure master and slaves nodes with hostname instead of ip address.

Following is the output on HadoopNameNode `ifconfig`

```
ubuntu@ip-172-31-35-98: ~
To run a command as administrator (user "root"), use "sudo <command>".
See "man sudo_root" for details.

ubuntu@ip-172-31-35-98:~$ ifconfig
eth0      Link encap:Ethernet  HWaddr 12:cb:69:95:84:79
          inet addr:172.31.35.98  Bcast:172.31.47.255  Mask:255.255.240.0
          inet6 addr: fe80::10cb:69ff:fe95:8479/64 Scope:Link
          UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1
          RX packets:378  errors:0  dropped:0  overruns:0  frame:0
          TX packets:374  errors:0  dropped:0  overruns:0  carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:37541 (37.5 KB)  TX bytes:40496 (40.4 KB)
          Interrupt:25

lo        Link encap:Local Loopback
          inet addr:127.0.0.1  Mask:255.0.0.0
          inet6 addr: ::1/128 Scope:Host
          UP LOOPBACK RUNNING  MTU:16436  Metric:1
          RX packets:0  errors:0  dropped:0  overruns:0  frame:0
          TX packets:0  errors:0  dropped:0  overruns:0  carrier:0
          collisions:0 txqueuelen:0
          RX bytes:0 (0.0 B)  TX bytes:0 (0.0 B)

ubuntu@ip-172-31-35-98:~$
```

This IP is same as what we have in the console

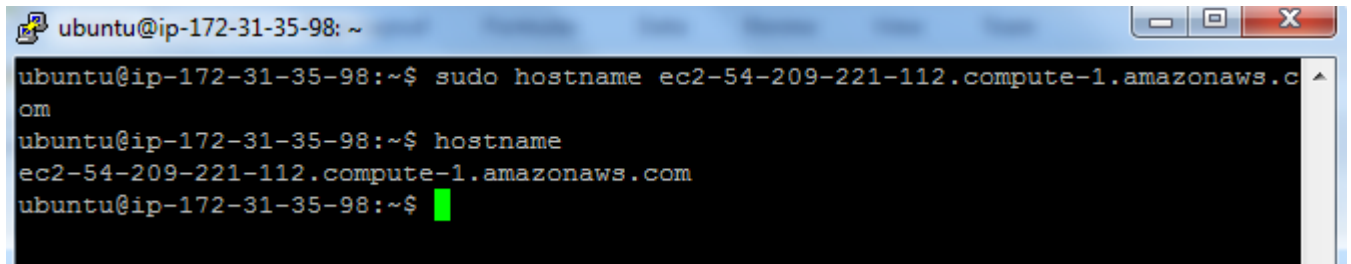
A	B	E
AMI name	public dns	inet address
HadoopNameNode	ec2-54-209-221-112.compute-1.amazonaws.com	172.31.35.98

Now, issue the hostname command, it will display the ip address same as inet address from ifconfig command.

```
ubuntu@ip-172-31-35-98: ~
ubuntu@ip-172-31-35-98:~$ hostname
ip-172-31-35-98
ubuntu@ip-172-31-35-98:~$
```

We need to modify the hostname to ec2 public URL with below command

~\$ sudo hostname ec2.....compute-1.amazonaws.com (Please put the URL which you got)

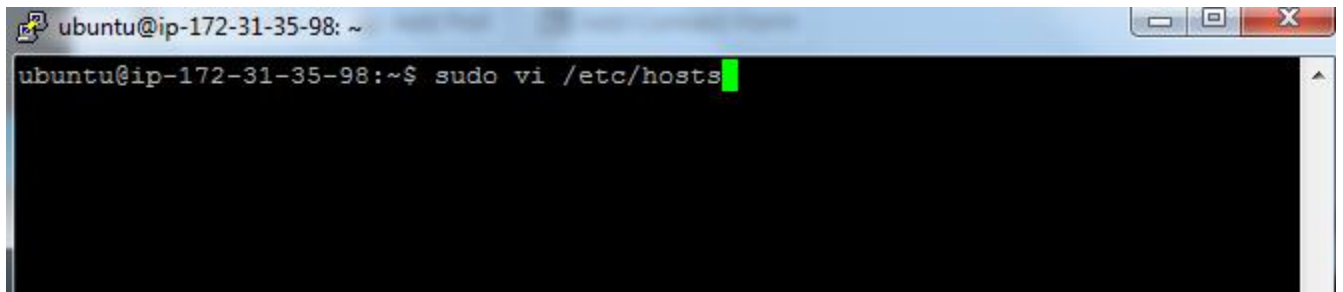
A terminal window titled 'ubuntu@ip-172-31-35-98: ~' showing the execution of the 'sudo hostname' command. The command is 'sudo hostname ec2-54-209-221-112.compute-1.amazonaws.com'. The output shows the hostname is set to 'ec2-54-209-221-112.compute-1.amazonaws.com'.

```
ubuntu@ip-172-31-35-98: ~  
ubuntu@ip-172-31-35-98:~$ sudo hostname ec2-54-209-221-112.compute-1.amazonaws.com  
om  
ubuntu@ip-172-31-35-98:~$ hostname  
ec2-54-209-221-112.compute-1.amazonaws.com  
ubuntu@ip-172-31-35-98:~$
```

## 2.5 Modify /etc/hosts

Let's change the host to EC2 public IP and hostname.

Open the /etc/hosts in 'VI' with "sudo vi /etc/hosts" command, in a very first line it will show 127.0.0.1 localhost, we need to replace that with amazon ec2 hostname and ip address we just collected.

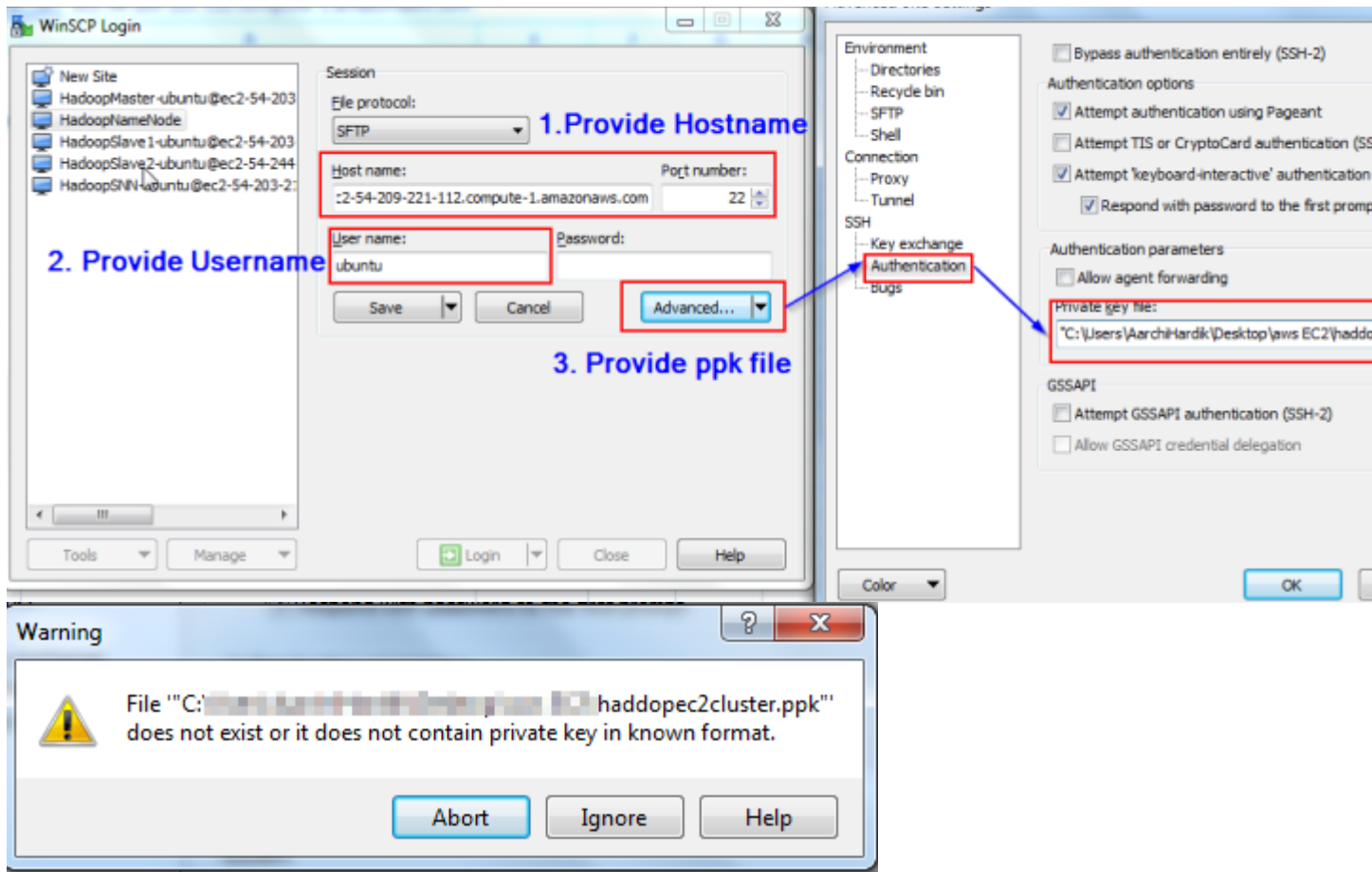
A terminal window titled 'ubuntu@ip-172-31-35-98: ~' showing the execution of the 'sudo vi /etc/hosts' command. The cursor is at the end of the command.

```
ubuntu@ip-172-31-35-98: ~  
ubuntu@ip-172-31-35-98:~$ sudo vi /etc/hosts
```

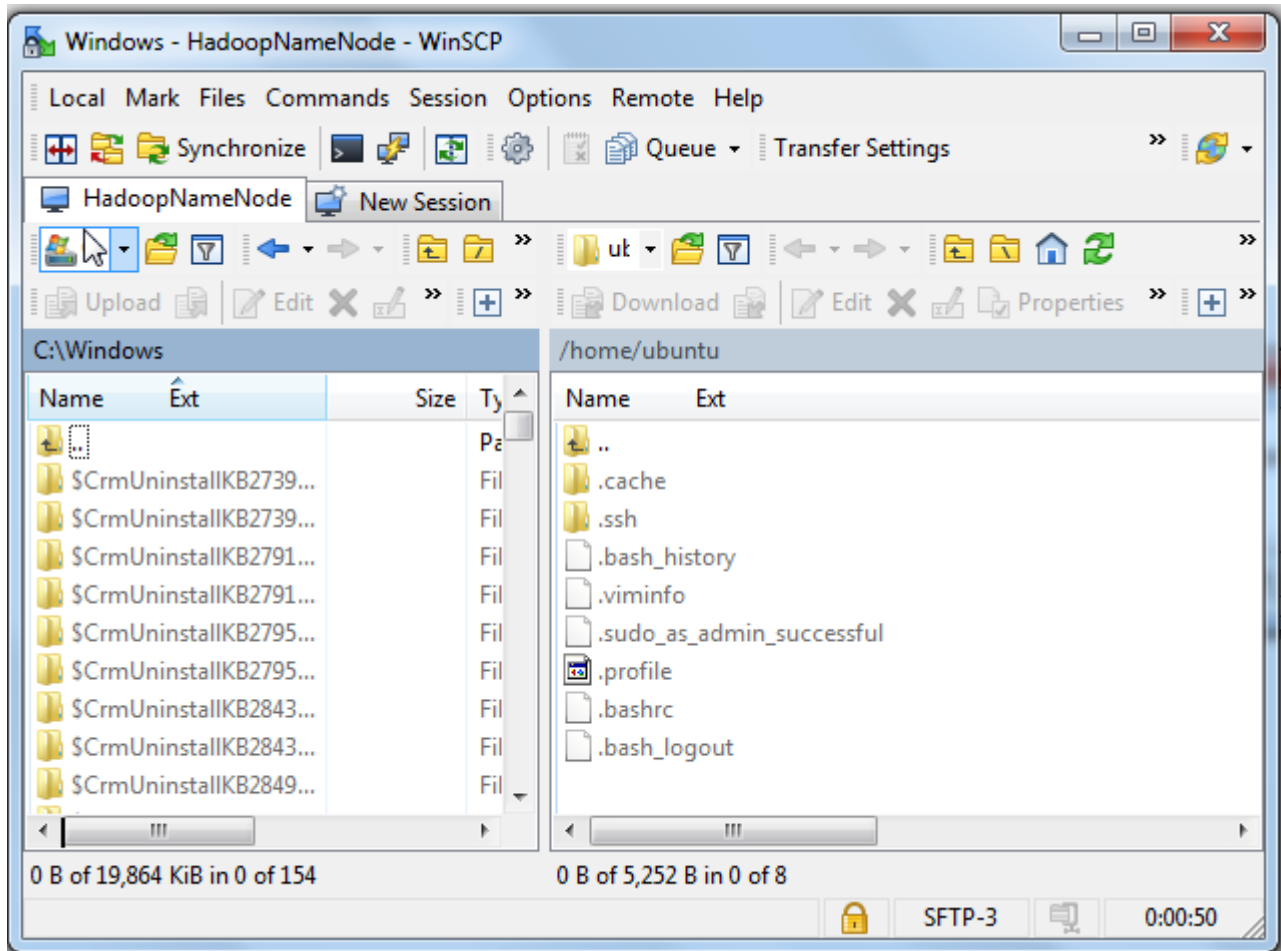
Modify the file and save your changes. To save a file you need to type the following sequence

1. `Esc`, to quit edit mode and fallback to command mode
2. `:wq`, start with a colon and then press w and q to write and quit
3. Then press `Enter` to validate.





If you see above error, just ignore and you upon successful login you will see unix file system of a logged in user /home/ubuntu your Amazon EC2 Ubuntu machine.



Upload the .pem file to master machine (HadoopNameNode). It will be used while connecting to slave nodes during hadoop startup daemons.

## 1. Apache Hadoop Installation and Cluster Setup

1.1 Update the packages and dependencies.

Let's update the packages, I will start with master, **repeat this for SecondaryNameNode and 2 slaves.**

Open the connection to the MasterNode using the steps provided in '2.3 Connect to Amazon Instance'

Type the following in the opened terminal

```
$ sudo apt-get update
```

Once it's complete, let's install java

## 1.2 Install Java

Add following PPA and install the latest Oracle Java (JDK) 7 in Ubuntu

```
$ sudo add-apt-repository ppa:webupd8team/java
```

Then type

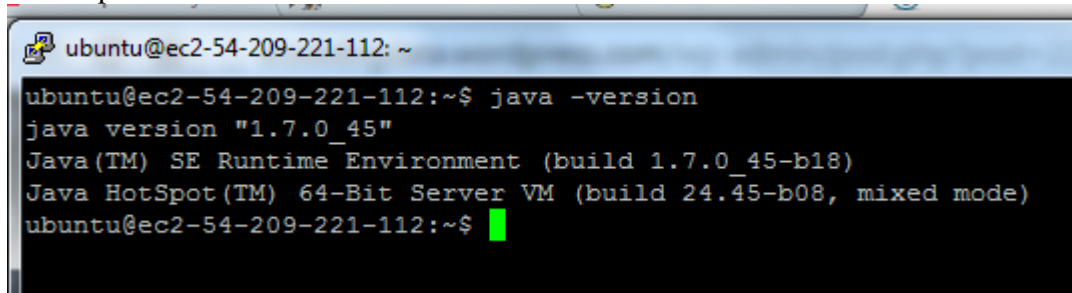
```
$ sudo apt-get update && sudo apt-get install oracle-jdk7-installer
```

Check if Ubuntu uses JDK 7

Type:

```
$ java -version
```

The response should be like this

A terminal window screenshot showing the command 'java -version' and its output. The terminal title is 'ubuntu@ec2-54-209-221-112: ~'. The output is: 'java version "1.7.0\_45"', 'Java(TM) SE Runtime Environment (build 1.7.0\_45-b18)', and 'Java HotSpot(TM) 64-Bit Server VM (build 24.45-b08, mixed mode)'. The prompt returns to 'ubuntu@ec2-54-209-221-112:~\$' with a green cursor.

```
ubuntu@ec2-54-209-221-112: ~  
ubuntu@ec2-54-209-221-112:~$ java -version  
java version "1.7.0_45"  
Java(TM) SE Runtime Environment (build 1.7.0_45-b18)  
Java HotSpot(TM) 64-Bit Server VM (build 24.45-b08, mixed mode)  
ubuntu@ec2-54-209-221-112:~$
```

**Repeat this for SNN and 2 slaves.**

## 1.3 Download Hadoop

I am going to use haddop 1.2.1 stable version from apache [download](#) page and here is the [1.2.1 mirror](#)

Issue wget command from shell

```
$ wget http://apache.mirror.gtcomm.net/hadoop/common/hadoop-1.2.1/hadoop-1.2.1.tar.gz
```

```
ubuntu@ec2-54-209-221-112: ~  
ubuntu@ec2-54-209-221-112:~$ wget http://apache.mirror.gtcomm.net/hadoop/common/  
hadoop-1.2.1/hadoop-1.2.1.tar.gz  
--2014-01-11 22:23:10-- http://apache.mirror.gtcomm.net/hadoop/common/hadoop-1.  
2.1/hadoop-1.2.1.tar.gz  
Resolving apache.mirror.gtcomm.net (apache.mirror.gtcomm.net)... 67.215.8.196  
Connecting to apache.mirror.gtcomm.net (apache.mirror.gtcomm.net)|67.215.8.196|:  
80... connected.  
HTTP request sent, awaiting response... 200 OK  
Length: 63851630 (61M) [application/x-gzip]  
Saving to: `hadoop-1.2.1.tar.gz'  
  
54% [=====>] 35,072,702 5.25M/s eta 6s
```

Unzip the files and review the package content and configuration files.

```
$ tar -xzvf hadoop-1.2.1.tar.gz
```

```
ubuntu@ec2-54-209-221-112:~$ ls  
hadoop-1.2.1 hadoop-1.2.1.tar.gz hadooppec2cluster.pem  
ubuntu@ec2-54-209-221-112:~$
```

For simplicity, rename the 'hadoop-1.2.1' directory to 'hadoop' for ease of operation and maintenance.

```
$ mv hadoop-1.2.1 hadoop
```

```
ubuntu@ec2-54-209-221-112: ~  
ubuntu@ec2-54-209-221-112:~$ ls  
hadoop hadoop-1.2.1.tar.gz hadooppec2cluster.pem  
ubuntu@ec2-54-209-221-112:~$
```

#### 1.4 Setup Environment Variable

Setup Environment Variable for 'ubuntu' user

Update the .bashrc file to add important Hadoop paths and directories.

Navigate to home directory



```
$ cd
```

Open .bashrc file in vi edit

```
$ vi .bashrc
```

Add following at the end of file

```
export HADOOP_CONF=/home/ubuntu/hadoop/conf
export HADOOP_PREFIX=/home/ubuntu/hadoop
```

```
#Set JAVA_HOME
```

```
export JAVA_HOME=/usr/lib/jvm/java-7-oracle
```

```
# Add Hadoop bin/ directory to path
```

```
export PATH=$PATH:$HADOOP_PREFIX/bin
```

Save and Exit by pressing escape and typing ':wq' and enter to validate.

To check whether it has been updated correctly or not, reload bash profile, use following commands

```
$ source ~/.bashrc
```

By typing the following two commands, there should be some value which should come up

```
$ echo $HADOOP_PREFIX
```

```
$ echo $HADOOP_CONF
```

**Repeat 1.3 and 1.4 for remaining 3 machines (SNN and 2 slaves).**

## 1.5 Setup Password-less SSH on Servers

Master server remotely starts services on slave nodes, which requires password-less access to Slave Servers. AWS Ubuntu server comes with pre-installed OpenSSH server.

### Quick Note:

The public part of the key loaded into the agent must be put on the target system in ~/.ssh/authorized\_keys. This has been taken care of by the AWS Server creation process

Now we need to add the AWS EC2 Key Pair identity 'HadoopEc2cluster.pem' to SSH profile. In order to do that we will need to use following ssh utilities

- 'ssh-agent' is a background program that handles passwords for SSH private keys.
- 'ssh-add' command prompts the user for a private key password and adds it to the list maintained by ssh-agent. Once you add a password to ssh-agent, you will not be asked to provide the key when using SSH or SCP to connect to hosts with your public key. Amazon EC2 Instance has already taken care of 'authorized\_keys' on master server, execute following commands to allow password-less SSH access to slave servers.

First of all we need to protect our keypair files, if the file permissions are too open (see below) you will get an error

```
ubuntu@ubuntu:~$ ls -l hadoop-1.2.1.tar.gz hadoop-ec2cluster.pem
-rw-rw-r-- 1 ubuntu ubuntu 1696 Jan 10 23:04 hadoop-1.2.1.tar.gz
-rw-rw-r-- 1 ubuntu ubuntu 1696 Jan 10 23:04 hadoop-ec2cluster.pem
ubuntu@ec2-54-209-221-112:~$ ls
hadoop hadoop-1.2.1.tar.gz hadoop-ec2cluster.pem
ubuntu@ec2-54-209-221-112:~$ eval `ssh-agent -s`
Agent pid 5645
ubuntu@ec2-54-209-221-112:~$ ssh-add hadoop-ec2cluster.pem
@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
@                WARNING: UNPROTECTED PRIVATE KEY FILE!                @
@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
Permissions 0664 for 'hadoop-ec2cluster.pem' are too open.
It is required that your private key files are NOT accessible by others.
This private key will be ignored.
ubuntu@ec2-54-209-221-112:~$ ls
```

To fix this problem, we need to issue following commands

```
$ chmod 644 authorized_keys
```

**Quick Tip:** If you set the permissions to 'chmod 644', you get a file that can be written by you, but can only be read by the rest of the world.

```
$ chmod 400 hadoop-ec2cluster.pem
```

**Quick Tip:** chmod 400 is a very restrictive setting giving only the file owner read-only access. No write / execute capabilities for the owner, and no permissions what-so-ever for anyone else.

To use ssh-agent and ssh-add, follow the steps below:

1. At the Unix prompt, enter: eval `ssh-agent`  
**Note:** Make sure you use the backquote (`), located under the tilde (~), rather than the single quote (').

2. Enter the command: 'ssh-add hadoopec2cluster.pem'. Make sure you are in the directory where this .pem file is.

It should work this time.

```
ubuntu@ec2-54-209-221-112:~$ chmod 644 .ssh/authorized_keys
ubuntu@ec2-54-209-221-112:~$ chmod 400 hadoopec2cluster.pem
ubuntu@ec2-54-209-221-112:~$ eval `ssh-agent -s`
Agent pid 5652
ubuntu@ec2-54-209-221-112:~$ ssh-add hadoopec2cluster.pem
Identity added: hadoopec2cluster.pem (hadoopec2cluster.pem)
ubuntu@ec2-54-209-221-112:~$ ls -l
total 62364
drwxr-xr-x 15 ubuntu ubuntu 4096 Jul 22 22:26 hadoop
-rw-rw-r-- 1 ubuntu ubuntu 63851630 Jul 22 22:27 hadoop-1.2.1.tar.gz
-r----- 1 ubuntu ubuntu 1696 Jan 10 23:04 hadoopec2cluster.pem
ubuntu@ec2-54-209-221-112:~$
```

Keep in mind ssh session will be lost upon shell exit and you have repeat ssh-agent and ssh-add commands.

### Remote SSH

Let's verify that we can connect into SNN and slave nodes from master

```
ubuntu@ec2-54-209-221-112:~$ ssh ubuntu@ec2-54-209-219-2.compute-1.amazonaws.com
Welcome to Ubuntu 12.04.3 LTS (GNU/Linux 3.2.0-54-virtual x86_64)

* Documentation:  https://help.ubuntu.com/

System information as of Sun Jan 12 01:00:03 UTC 2014

System load:  0.0          Processes:           61
Usage of /:   20.1% of 7.87GB  Users logged in:   1
Memory usage: 27%          IP address for eth0: 172.31.35.95
Swap usage:   0%

Graph this data and manage this system at https://landscape.canonical.com/

Get cloud support with Ubuntu Advantage Cloud Guest:
http://www.ubuntu.com/business/services/cloud

Use Juju to deploy your cloud instances and workloads:
https://juju.ubuntu.com/#cloud-precise

Last login: Sun Jan 12 00:59:39 2014 from ip-172-31-35-98.ec2.internal
ubuntu@ec2-54-209-219-2:~$
```

\$ ssh ubuntu@<your-amazon-ec2-public URL for SNN or any of your slave nodes>

On successful login the IP address on the shell will change.

Type the following to exit from SNN or other nodes and to come back to the master node.

```
$ exit
```

## 1.6 Hadoop Cluster Setup

This section will cover the hadoop cluster configuration. We will have to modify

- **hadoop-env.sh** - This file contains some environment variable settings used by Hadoop. You can use these to affect some aspects of Hadoop daemon behavior, such as where log files are stored, the maximum amount of heap used etc. The only variable you should need to change at this point is in this file is `JAVA_HOME`, which specifies the path to the Java 1.7.x installation used by Hadoop.
- **core-site.xml** – key property `fs.default.name` – for namenode configuration for e.g `hdfs://namenode/`
- **hdfs-site.xml** – key property - `dfs.replication` – by default 3
- **mapred-site.xml** - key property `mapred.job.tracker` for jobtracker configuration for e.g `jobtracker:8021`

We will first start with master (NameNode) and then copy above xml changes to remaining 3 nodes (SNN and slaves)

Finally, in section 1.6.2 we will have to configure `conf/masters` and `conf/slaves`.

- **masters** - defines on which machines Hadoop will start secondary NameNodes in our multi-node cluster.
- **slaves** - defines the lists of hosts, one per line, where the Hadoop slave daemons (datanodes and tasktrackers) will run.

**Lets go over one by one. Start with masters (namenode).**

**Perform the following**

**hadoop-env.sh**

```
$ vi $HADOOP_CONF/hadoop-env.sh and add JAVA_HOME shown below and save changes.
```

```
ubuntu@ec2-54-209-221-112: ~/hadoop/conf
# Set Hadoop-specific environment variables here.

# The only required environment variable is JAVA_HOME. All others are
# optional. When running a distributed configuration it is best to
# set JAVA_HOME in this file, so that it is correctly defined on
# remote nodes.

# The java implementation to use. Required.
export JAVA_HOME=/usr/lib/jvm/java-7-oracle
```

### core-site.xml

This file contains configuration settings for Hadoop Core (for e.g I/O) that are common to HDFS and MapReduce Default file system configuration property – fs.default.name goes here it could for e.g hdfs / s3 which will be used by clients.

```
$ sudo vi $HADOOP_CONF/core-site.xml
```

We are going to add two properties

- fs.default.name will point to NameNode URL and port (usually 8020)
- hadoop.tmp.dir - A base for other temporary directories. Its important to note that every node needs hadoop tmp directory. I am going to create a new directory “hdfstmp” as below in all 4 nodes. Ideally you can write a shell script to do this for you, but for now going the manual way.
- **Perform the following**

### Exit from core-site.xml

Then

```
$ cd
```

```
$ mkdir hdfstmp
```

**Quick Tip:** Some of the important directories are dfs.name.dir, dfs.data.dir in hdfs-site.xml. The default value for the dfs.name.dir is \${hadoop.tmp.dir}/dfs/data and dfs.data.dir is \${hadoop.tmp.dir}/dfs/data. It is **critical** that you choose your directory location wisely in production environment.

Fill the following in the \$HADOOP\_CONF/core-site.xml

```
<configuration>
<property>
<name>fs.default.name</name>
<value>hdfs://ec2-54-209-221-112.compute-1.amazonaws.com:8020</value>
</property>
<property>
<name>hadoop.tmp.dir</name>
<value>/home/ubuntu/hdfstmp</value>
</property>
</configuration>
```

## Save and Exit

### hdfs-site.xml

This file contains the configuration for HDFS daemons, the NameNode, SecondaryNameNode and data nodes.

We are going to add 2 properties

- **dfs.permissions.enabled** with value *false*, This means that any user, not just the “hdfs” user, can do anything they want to HDFS so do not do this in production unless you have a very good reason. If “true”, enable permission checking in HDFS. If “false”, permission checking is turned off, but all other behavior is unchanged. Switching from one parameter value to the other does not change the mode, owner or group of files or directories. Be very careful before you set this
- **dfs.replication** – Default block replication is 3. The actual number of replications can be specified when the file is created. The default is used if replication is not specified in create time. Since we have 2 slave nodes we will set this value to 2.

Perform the following

```
$ sudo vi $HADOOP_CONF/ hdfs-site.xml
```

Fill it with the following

```
<configuration>
<property>
<name>dfs.replication</name>
<value>2</value>
</property>
```

```
<property>
<name>dfs.permissions</name>
<value>>false</value>
</property>
</configuration>
```

### Save and exit

```
<configuration>
<property>
<name>dfs.replication</name>
<value>2</value>
</property>
<property>
<name>dfs.permissions</name>
<value>>false</value>
</property>
</configuration>
```

### mapred-site.xml

This file contains the configuration settings for MapReduce daemons; the job tracker and the task-trackers.

The mapred.job.tracker parameter is a hostname (or IP address) and port pair on which the Job Tracker listens for RPC communication. This parameter specifies the location of the Job Tracker for Task Trackers and MapReduce clients.

JobTracker will be running on **master** (NameNode)

Perform the following

```
$ sudo vi $HADOOP_CONF/mapred-site.xml
```

Fill it with the following

```
<configuration>
<property>
<name>mapred.job.tracker</name>
<value>hdfs://ec2-54-209-221-112.compute-1.amazonaws.com:8021</value>
</property>
</configuration>
```

### Save and Exit

## 1.6.1 Move configuration files to Slaves

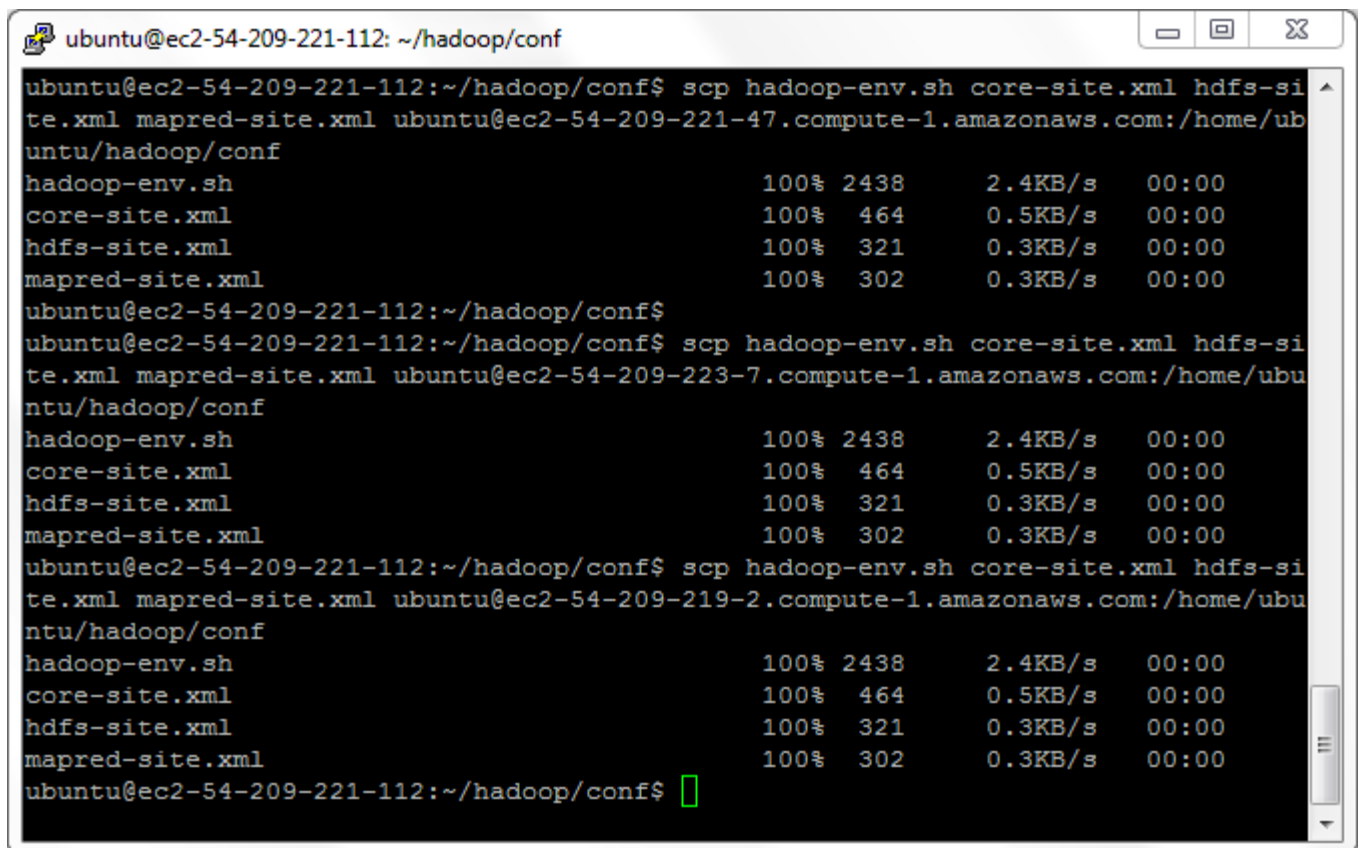
Now, we are done with hadoop xml files configuration master, lets copy the files to remaining 3 nodes using secure copy (scp)

start with SNN, if you are starting a new session, follow ssh-add as per section 1.5

from master's unix shell issue below command

```
$ scp hadoop-env.sh core-site.xml hdfs-site.xml mapred-site.xml ubuntu@<URL of your Secondary Name node>:/home/ubuntu/hadoop/conf
```

**Repeat this for slave nodes and check if they got copied in all the nodes**



```
ubuntu@ec2-54-209-221-112: ~/hadoop/conf
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ scp hadoop-env.sh core-site.xml hdfs-site.xml mapred-site.xml ubuntu@ec2-54-209-221-47.compute-1.amazonaws.com:/home/ubuntu/hadoop/conf
hadoop-env.sh          100% 2438    2.4KB/s   00:00
core-site.xml          100%  464    0.5KB/s   00:00
hdfs-site.xml          100%  321    0.3KB/s   00:00
mapred-site.xml        100%  302    0.3KB/s   00:00
ubuntu@ec2-54-209-221-112:~/hadoop/conf$
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ scp hadoop-env.sh core-site.xml hdfs-site.xml mapred-site.xml ubuntu@ec2-54-209-223-7.compute-1.amazonaws.com:/home/ubuntu/hadoop/conf
hadoop-env.sh          100% 2438    2.4KB/s   00:00
core-site.xml          100%  464    0.5KB/s   00:00
hdfs-site.xml          100%  321    0.3KB/s   00:00
mapred-site.xml        100%  302    0.3KB/s   00:00
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ scp hadoop-env.sh core-site.xml hdfs-site.xml mapred-site.xml ubuntu@ec2-54-209-219-2.compute-1.amazonaws.com:/home/ubuntu/hadoop/conf
hadoop-env.sh          100% 2438    2.4KB/s   00:00
core-site.xml          100%  464    0.5KB/s   00:00
hdfs-site.xml          100%  321    0.3KB/s   00:00
mapred-site.xml        100%  302    0.3KB/s   00:00
ubuntu@ec2-54-209-221-112:~/hadoop/conf$
```

## 1.6.2 Configure Master and Slaves

Every hadoop distribution comes with master and slaves files. By default it contains one entry for localhost, we have to modify these 2 files on both “masters” (HadoopNameNode) and “slaves” (HadoopSlave1 and HadoopSlave2) machines – we have a dedicated machine for HadoopSecondaryNameNode.



```

ubuntu@ec2-54-209-221-112:~/hadoop/conf$ ls
capacity-scheduler.xml      hadoop-policy.xml          slaves
configuration.xsl          hdfs-site.xml             ssl-client.xml.example
core-site.xml               log4j.properties         ssl-server.xml.example
fair-scheduler.xml         mapred-queue-acls.xml     taskcontroller.cfg
hadoop-env.sh              mapred-site.xml           task-log4j.properties
hadoop-metrics2.properties masters

```

### 1.6.3 Modify masters file on Master machine

conf/masters file defines on which machines Hadoop will start *Secondary NameNodes* in our multi-node cluster. In our case, there will be two machines HadoopNameNode and HadoopSecondaryNameNode

[Hadoop HDFS user guide](#) : “The secondary NameNode merges the fsimage and the edits log files periodically and keeps edits log size within a limit. It is usually run on a different machine than the primary NameNode since its memory requirements are on the same order as the primary NameNode. The secondary NameNode is started by “bin/start-dfs.sh“ on the nodes specified in “conf/masters“ file.”

#### Perform the following

\$ vi \$HADOOP\_CONF/masters and provide an entry for the hostname where you want to run SecondaryNameNode daemon. In our case HadoopNameNode and HadoopSecondaryNameNode

```

ec2-54-209-221-112.compute-1.amazonaws.com
ec2-54-209-221-47.compute-1.amazonaws.com
~
~HadoopSecondaryNameNode      HadoopNameNode

```

### 1.6.4 MODIFY THE SLAVES FILE ON MASTER MACHINE

The slaves file is used for starting DataNodes and TaskTrackers

\$ vi \$HADOOP\_CONF/slaves

```

ec2-54-209-223-7.compute-1.amazonaws.com
ec2-54-209-219-2.compute-1.amazonaws.com

```

### 1.6.5 Copy masters and slaves to SecondaryNameNode

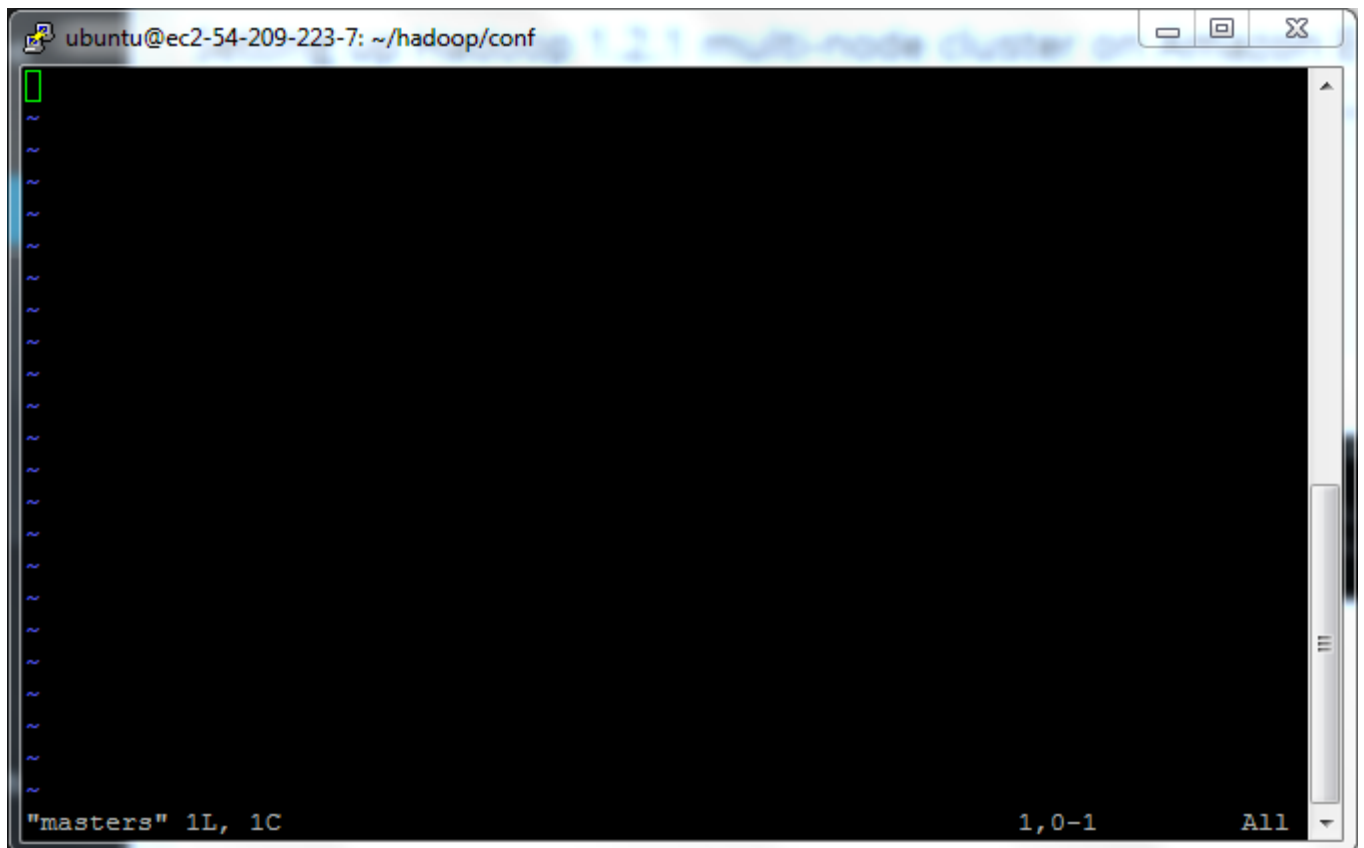
Since SecondaryNameNode configuration will be same as NameNode, we need to copy master and slaves to HadoopSecondaryNameNode.

```
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ scp masters slaves ubuntu@ec2-54-209-221-47.compute-1.amazonaws.com:~/hadoop/conf
masters          100%  83    0.1KB/s   00:00
slaves           100%  82    0.1KB/s   00:00
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ █
```

### 1.6.7 CONFIGURE MASTER AND SLAVES ON “SLAVES” NODE

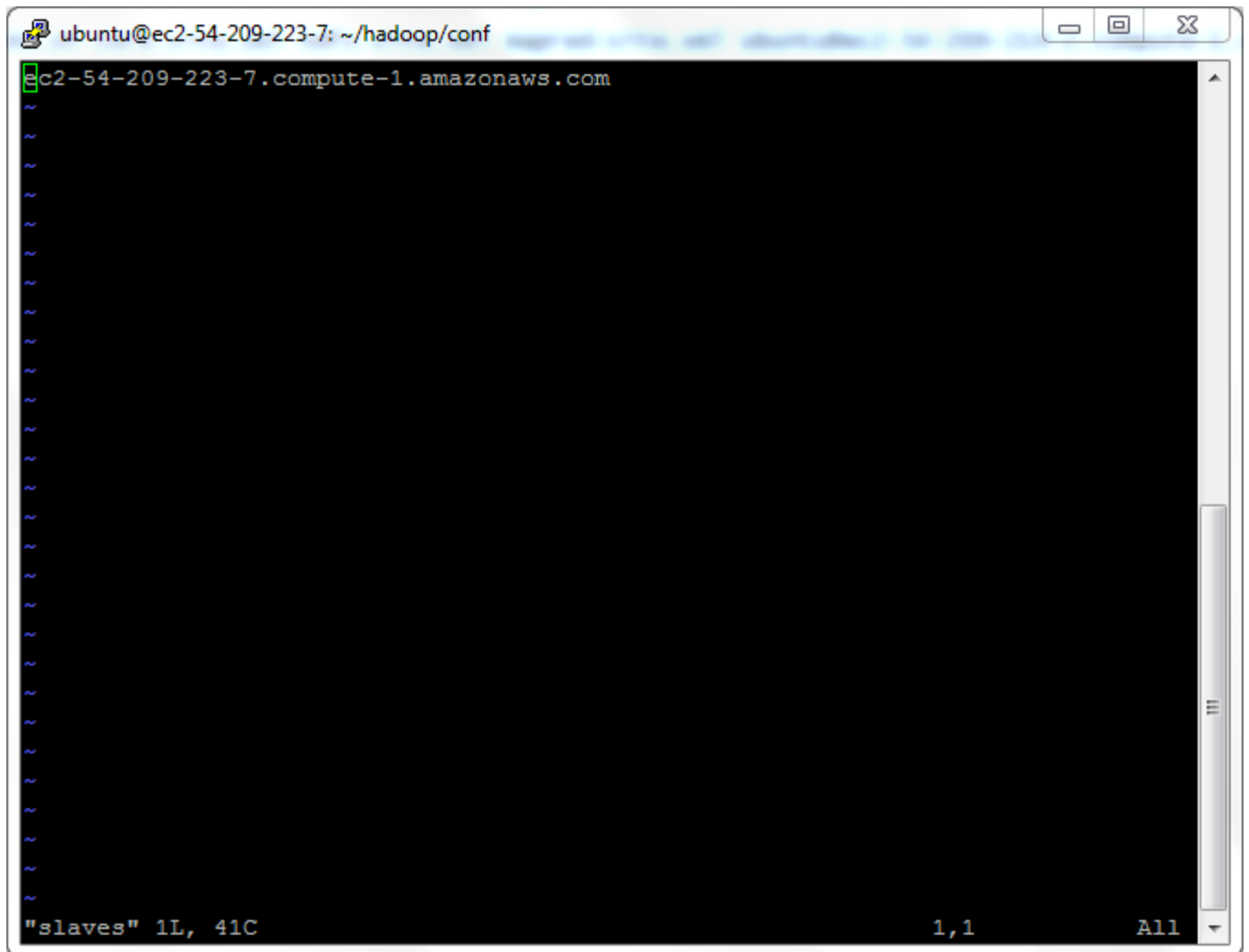
Since we are configuring slaves (HadoopSlave1 & HadoopSlave2), masters file on slave machine is going to be empty

```
$ vi $HADOOP_CONF/masters
```



Next, update the 'slaves' file on Slave server (HadoopSlave1) with the IP address of the *slave node*. Notice that the 'slaves' file at Slave node contains only its own IP address and not of any other Data Node in the cluster.

```
$ vi $HADOOP_CONF/slaves
```



The image shows a terminal window titled 'ubuntu@ec2-54-209-223-7: ~/hadoop/conf'. The prompt is 'ec2-54-209-223-7.compute-1.amazonaws.com'. The terminal displays a list of IP addresses, one per line, representing slave nodes. At the bottom of the terminal, the command '"slaves" 1L, 41C' is shown, along with '1,1' and 'All'.

Similarly update masters and slaves for HadoopSlave2

## 1.7 Hadoop Daemon Startup

The first step to starting up your Hadoop installation is formatting the Hadoop filesystem which is implemented on top of the local filesystems of your cluster. You need to do this the first time you set up a Hadoop installation. **Do not format a running Hadoop filesystem**, this will cause all your data to be erased.

To format the namenode

Goto Namenode(master node) and perform the following

```
$ hadoop namenode -format
```

```

ubuntu@ec2-54-209-221-112:~/hadoop/conf$ hadoop namenode -format
14/01/13 02:36:27 INFO namenode.NameNode: STARTUP_MSG:
/*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = ec2-54-209-221-112.compute-1.amazonaws.com/172.31.35.98
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 1.2.1
STARTUP_MSG: build = https://svn.apache.org/repos/asf/hadoop/common/branches/branch-1.2 -r 1503152; compiled by 'mattf' on Mon Jul 22
STARTUP_MSG: java = 1.7.0_45
*****/
14/01/13 02:36:28 INFO util.GSet: Computing capacity for map BlocksMap
14/01/13 02:36:28 INFO util.GSet: VM type = 64-bit
14/01/13 02:36:28 INFO util.GSet: 2.0% max memory = 1013645312
14/01/13 02:36:28 INFO util.GSet: capacity = 2^21 = 2097152 entries
14/01/13 02:36:28 INFO util.GSet: recommended=2097152, actual=2097152
14/01/13 02:36:28 INFO namenode.FSNamesystem: fsOwner=ubuntu
14/01/13 02:36:28 INFO namenode.FSNamesystem: supergroup=supergroup
14/01/13 02:36:28 INFO namenode.FSNamesystem: isPermissionEnabled=false
14/01/13 02:36:28 INFO namenode.FSNamesystem: dfs.block.invalidate.limit=100
14/01/13 02:36:28 INFO namenode.FSNamesystem: isAccessTokenEnabled=false accessKeyUpdateInterval=0 min(s), accessTokenLifetime=0 min(s)
14/01/13 02:36:28 INFO namenode.FSEditLog: dfs.namenode.edits.toleration.length = 0
14/01/13 02:36:28 INFO namenode.NameNode: Caching file names occurring more than 10 times
14/01/13 02:36:29 INFO common.Storage: Image file /home/ubuntu/hdfstmp/dfs/name/current/fsimage of size 112 bytes saved in 0 seconds.
14/01/13 02:36:29 INFO namenode.FSEditLog: closing edit log: position=4, editlog=/home/ubuntu/hdfstmp/dfs/name/current/edits
14/01/13 02:36:29 INFO namenode.FSEditLog: close success: truncate to 4, editlog=/home/ubuntu/hdfstmp/dfs/name/current/edits
14/01/13 02:36:30 INFO common.Storage: Storage directory /home/ubuntu/hdfstmp/dfs/name has been successfully formatted.
14/01/13 02:36:30 INFO namenode.NameNode: SHUTDOWN_MSG:
/*****
SHUTDOWN_MSG: Shutting down NameNode at ec2-54-209-221-112.compute-1.amazonaws.com/172.31.35.98
*****/
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ █

```

Lets start all hadoop daemons from HadoopNameNode

```
$ cd $HADOOP_CONF
```

```
$ start-all.sh
```

This will start

- NameNode,JobTracker and SecondaryNameNode daemons on **HadoopNameNode**

```

ubuntu@ec2-54-209-221-112: ~/hadoop/conf — ssh — 118x31
ubuntu@ec2.../conf — ssh  ubuntu@ec2...47: ~ — ssh  ubuntu@ec2...-7: ~ — ssh  ubuntu@ec2...-2: ~
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ clear

ubuntu@ec2-54-209-221-112:~/hadoop/conf$ start-all.sh
starting namenode, logging to /home/ubuntu/hadoop/libexec/./logs/hadoop-ubuntu-namenode-ec2-54-209-221-112.
amazonaws.com.out
ec2-54-209-219-2.compute-1.amazonaws.com: starting datanode logging to /home/ubuntu/hadoop/libexec/./logs/hadoop-ub
ntu-datanode-ec2-54-209-219-2.compute-1.amazonaws.com.out
ec2-54-209-223-7.compute-1.amazonaws.com: starting datanode logging to /home/ubuntu/hadoop/libexec/./logs/hadoop-ub
ntu-datanode-ec2-54-209-223-7.compute-1.amazonaws.com.out
The authenticity of host 'ec2-54-209-221-112.compute-1.amazonaws.com (172.31.35.98)' can't be established.
ECDSA key fingerprint is f3:90:74:77:31:5f:2f:f6:d7:5b:94:f1:0c:65:0d:df.
Are you sure you want to continue connecting (yes/no)? ec2-54-209-221-47.compute-1.amazonaws.com: starting
menode, logging to /home/ubuntu/hadoop/libexec/./logs/hadoop-ubuntu-secondarynamenode-ec2-54-209-221-47.c
zonaws.com.out
yes
ec2-54-209-221-112.compute-1.amazonaws.com: Warning: Permanently added 'ec2-54-209-221-112.compute-1.amazo
.31.35.98' (ECDSA) to the list of known hosts.
ec2-54-209-221-112.compute-1.amazonaws.com: starting secondarynamenode, logging to /home/ubuntu/hadoop/lib
/hadoop-ubuntu-secondarynamenode-ec2-54-209-221-112.compute-1.amazonaws.com.out
starting jobtracker, logging to /home/ubuntu/hadoop/libexec/./logs/hadoop-ubuntu-jobtracker-ec2-54-209-22
e-1.amazonaws.com.out
ec2-54-209-219-2.compute-1.amazonaws.com: starting tasktracker, logging to /home/ubuntu/hadoop/libexec/./
ubuntu-tasktracker-ec2-54-209-219-2.compute-1.amazonaws.com.out
ec2-54-209-223-7.compute-1.amazonaws.com: starting tasktracker, logging to /home/ubuntu/hadoop/libexec/./
ubuntu-tasktracker-ec2-54-209-223-7.compute-1.amazonaws.com.out
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ jps
15184 JobTracker
15100 SecondaryNameNode
15269 Jps
14857 NameNode
ubuntu@ec2-54-209-221-112:~/hadoop/conf$ █

```

Starting all hadoop daemons from HadoopNameNode

- SecondaryNameNode daemons on **HadoopSecondaryNameNode**

```

Last login: Mon Jan 13 03:15:02 2014 from bas1-malton23-1177880673.dsl.bell.ca
ubuntu@ec2-54-209-221-47:~$ ls
hadoop hadoop-1.2.1.tar.gz hdfstmp
ubuntu@ec2-54-209-221-47:~$ jps
10748 Jps
10522 SecondaryNameNode
ubuntu@ec2-54-209-221-47:~$ █

```

- and DataNode and TaskTracker daemons on slave nodes **HadoopSlave1** and **HadoopSlave2**

```

ubuntu@ec2-54-209-223-7:~$ jps
11816 DataNode
12211 Jps
11977 TaskTracker
ubuntu@ec2-54-209-223-7:~$ █

```

```
ubuntu@ec2-54-209-219-2:~$ jps
11210 Jps
10976 TaskTracker
10815 DataNode
ubuntu@ec2-54-209-219-2:~$ █
```

We can check the namenode status from <http://ec2-54-209-221-112.compute-1.amazonaws.com:50070/dfshealth.jsp>

← → ↻ ec2-54-209-221-112.compute-1.amazonaws.com:50070/dfshealth.jsp

## NameNode 'ec2-54-209-221-112.compute-1.amazonaws.com:8020'

**Started:** Mon Jan 13 14:23:22 UTC 2014  
**Version:** 1.2.1, r1503152  
**Compiled:** Mon Jul 22 15:23:09 PDT 2013 by mattf  
**Upgrades:** There are no upgrades in progress.

[Browse the filesystem](#)  
[Namenode Logs](#)

### Cluster Summary

7 files and directories, 1 blocks = 8 total. Heap Size is 32.65 MB / 966.69 MB (3%)

Configured Capacity	: 15.75 GB
DFS Used	: 56 KB
Non DFS Used	: 4.14 GB
DFS Remaining	: 11.61 GB
DFS Used%	: 0 %
DFS Remaining%	: 73.71 %
<a href="#">Live Nodes</a>	: 2
<a href="#">Dead Nodes</a>	: 0
<a href="#">Decommissioning Nodes</a>	: 0
Number of Under-Replicated Blocks	: 0

### NameNode Storage:

Storage Directory	Type	State
/home/ubuntu/hd/stmp/dfs/name	IMAGE_AND_EDITS	Active

This is [Apache Hadoop](#) release 1.2.1

Check Jobtracker status : <http://<Your AMAZON MASTER URL>:50030/jobtracker.jsp>

## ec2-54-209-221-112 Hadoop Map/Reduce Administration

**State:** RUNNING  
**Started:** Mon Jan 13 14:25:09 UTC 2014  
**Version:** 1.2.1, r1503152  
**Compiled:** Mon Jul 22 15:23:09 PDT 2013 by mattf  
**Identifier:** 201401131425  
**SafeMode:** OFF

### Cluster Summary (Heap Size is 9.31 MB/966.69 MB)

Running Map Tasks	Running Reduce Tasks	Total Submissions	Nodes	Occupied Map Slots	Occupied Reduce Slots	Reserved Map Slots	Reserved Reduce Slots	Map Task Capacity	Reduce Task Capacity	Avg. Tasks/Node	Blacklisted Nodes
0	0	0	2	0	0	0	0	4	4	4.00	0

### Scheduling Information

Queue Name	State	Scheduling Information
<a href="#">default</a>	running	N/A

Filter (Jobid, Priority, User, Name)

Example: 'user:smith 3200' will filter by 'smith' only in the user field and '3200' in all fields

### Running Jobs

none

### Retired Jobs

none

### Local Logs

[Log](#) directory, [Job Tracker History](#)

This is [Apache Hadoop](#) release 1.2.1

Slave Node Status for HadoopSlave1 : http://<YOUR AMAZON MASTER URL>:50060/tasktracker.jsp

## tracker\_ec2-54-209-223-7.compute-1.amazonaws.com:127.0.0.1/127.0.0.1:39269 Task Tracker Status



Version: 1.2.1, r1503152  
Compiled: Mon Jul 22 15:23:09 PDT 2013 by mattf

### Running tasks

Task Attempts	Status	Progress	Errors
---------------	--------	----------	--------

### Non-Running Tasks

Task Attempts	Status
---------------	--------

### Tasks from Running Jobs

Task Attempts	Status	Progress	Errors
---------------	--------	----------	--------

### Local Logs

[Log](#) directory

This is [Apache Hadoop](#) release 1.2.1

Slave Node Status for HadoopSlave2 : <http://ec2-54-209-219-2.compute-1.amazonaws.com:50060/tasktracker.jsp>

← → C ec2-54-209-219-2.compute-1.amazonaws.com:50060/tasktracker.jsp

## tracker\_ec2-54-209-219-2.compute-1.amazonaws.com:127.0.0.1/127.0.0.1:48106 Task Tracker Status



Version: 1.2.1, r1503152  
Compiled: Mon Jul 22 15:23:09 PDT 2013 by mattf

### Running tasks

Task Attempts	Status	Progress	Errors
---------------	--------	----------	--------

### Non-Running Tasks

Task Attempts	Status
---------------	--------

### Tasks from Running Jobs

Task Attempts	Status	Progress	Errors
---------------	--------	----------	--------

### Local Logs

[Log](#) directory

This is [Apache Hadoop](#) release 1.2.1



## Running the code

**Step 1:** Browse to the 'Mapreduce\_Kmeans' directory and check if we have the following:

- a) 'Input' directory
- b) JAR file named 'ProcessCorpus.jar'
- c) JAR file named 'GetCentroids.jar'
- d) JAR file named 'MapRedKMeans.jar'

**Step 2:** Type 'java -jar ProcessCorpus.jar'. When prompted, answer:

Enter the directory where the corpus is located: **Input**

Enter the name of the file to write the result to: **vectors**

Enter the max number of docs to use in each subdirectory: **100**

How many of those words do you want to use...? **10000**

This will create a file called "vectors" that has 20 \* 100 lines, each of which is a vectorized representation of a document. The dictionary size that is used is 10000 words

**Step 3:** Now type 'java -jar GetCentroids.jar'. When prompted, answer:

Enter the data file to select the clusters from: **vectors**

Enter the name of the file to write the result to: **clusters**

Enter the number of clusters to select: **20**

This will create a file called "clusters" that has 20 lines, each of which describes a cluster centroids. This initial set of cluster centroids is simply randomly sampled from the "vectors" file.

**Step 4:** Now we'll copy the "vectors" and "clusters" files into HDFS:

```
hadoop fs -mkdir /user/ubuntu/data
```

```
hadoop fs -mkdir /user/ubuntu/clusters
```

```
hadoop fs -copyFromLocal vectors /user/ubuntu/data
```

```
hadoop fs -copyFromLocal clusters /user/ubuntu/clusters
```

**Step 5:** Run the kmeans jar by typing

```
$ hadoop jar MapRedKMeans.jar KMeans /user/ubuntu/data
```

/user/ubuntu/clusters 3

This will run for 3 iterations.

**Step 6:** When done, check the '/user/ubuntu/clusters2/part-r-00000' file on HDFS to get the cluster distribution.

```
$ Hadoop fs -cat /user/ubuntu/clusters/clusters2/part-r-00000
```

Or you can copy the clusters file on your local machine

```
$ Hadoop fs -copyToLocal /user/ubuntu/clusters2 /home/Ubuntu/kmeans/
```

Then, look into the part file in the clusters2 directory for distribution.

## Stopping the Hadoop Daemon

```
$ cd $HADOOP_CONF
```

```
$ stop-all.sh
```

### **Cleanup (Important)**

**Step 1:** Logon to Amazon AWS and under Services select 'Ec2'.

**Step 2:** Under the 'Instances' tab in the left column; click on 'Instances'.

**Step 3:** Locate all your Hadoop instances and select them. On the top locate 'Actions' drop down button and click 'Stop' to stop the instances. You can start it and connect to the same settings whenever you want. If you terminate it, you have to create a new instance all together.

### **Caveats**

When you stop and restart the amazon instances, the Public IP and the URL of the instances changes. You have to make changes in the following with the new URLs

1. hostname
2. Step [2.3.2](#)
3. /etc/hosts
4. \$HADOOP\_CONF/ core-site.xml
5. \$HADOOP\_CONF/ core-site.xml
6. \$HADOOP\_CONF/ masters
7. \$HADOOP\_CONF/ slaves
8. Repeat Step [1.5](#)
9. No need to format the namenode

## 10. Start the Hadoop daemon