

Tutorial- Counting Words in File(s) using Elastic MapReduce (AWS)

Prepared by- Srinivasan Rajappa

Under the guidance of Prof. Bina Ramamurthy

1 OVERVIEW

This document serves as a tutorial to setup and run a simple application in Elastic MapReduce which is a service provided by Amazon Web Services.

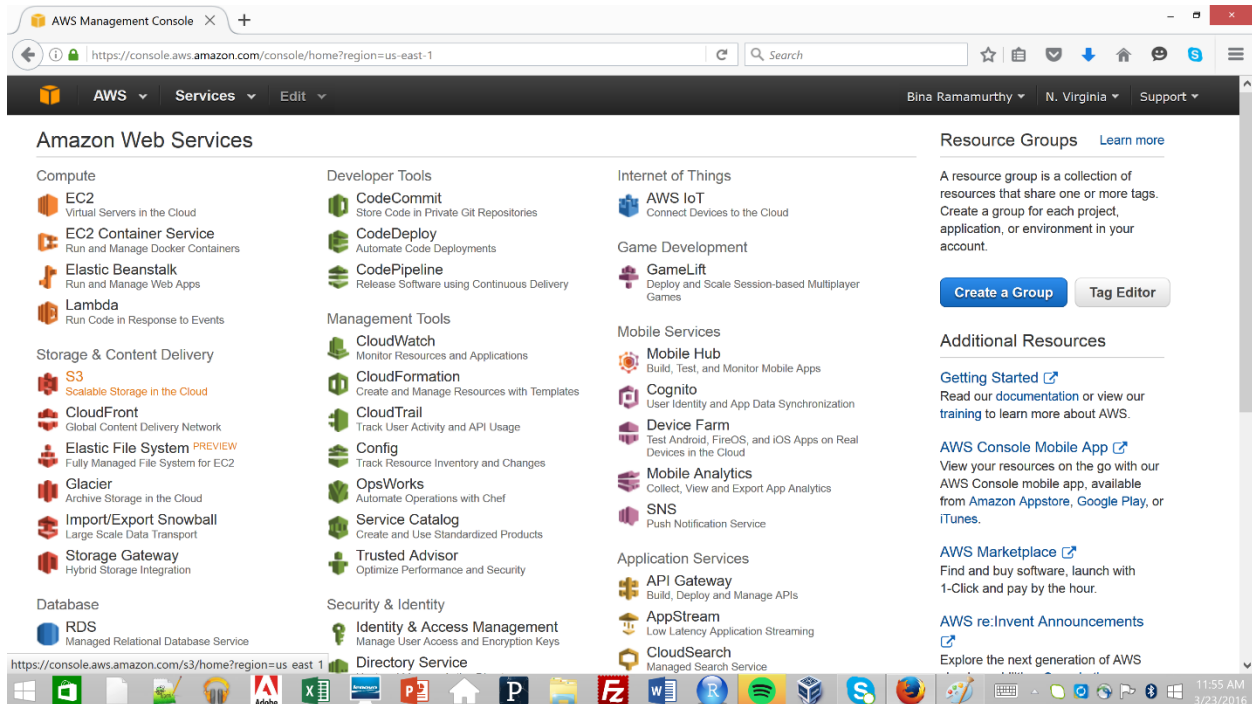
The process involves three phases viz. setting up an AWS Educate account, Creating Buckets within S3 and then to configure and run a cluster. In addition to covering these topics, additional notes and warnings are also be provided.

2 CREATING AWS EDUCATE ACCOUNT

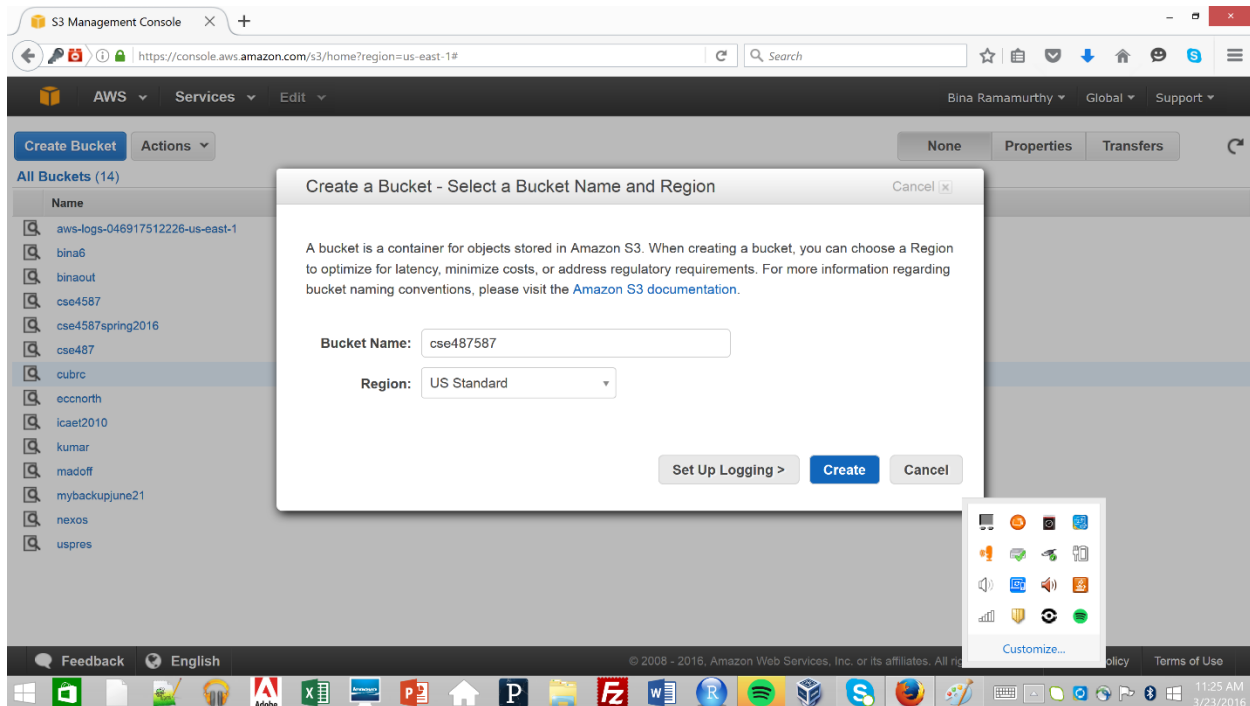
An account can be created by following this video link [here](#). Please ensure that you register using your [buffalo.edu](#) [\[why?\]](#) email account as it will provide various benefits.

3 S3 BUCKETS

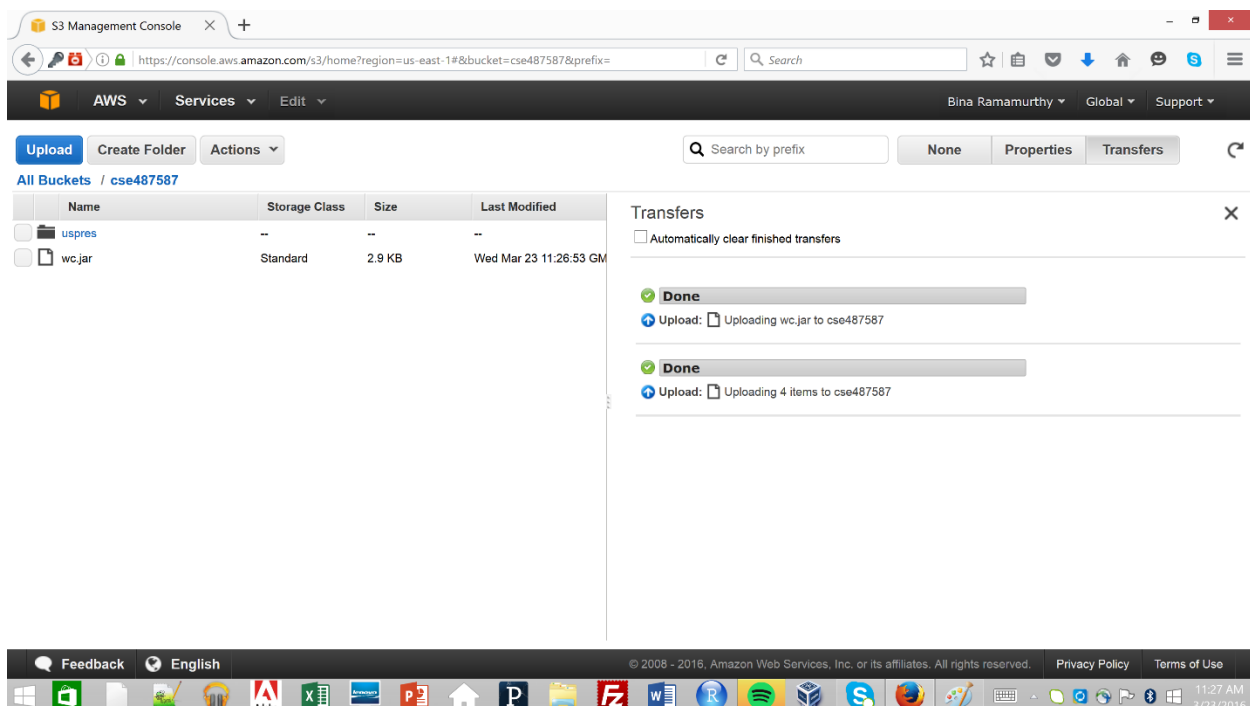
After successfully creating an account. Please login to this [AWS console](#). Once you reach the page, click on Simple Storage Service (S3). Please refer to the following screenshot:



After opening the S3 webpage, please click on create a new bucket. Soon a popup requesting details will appear. Provide a universally unique name to the bucket, it should only comprise of lower case alphabets and numbers^[why?]. Also, select the region as **US Standard**. Please refer to the following screenshot:

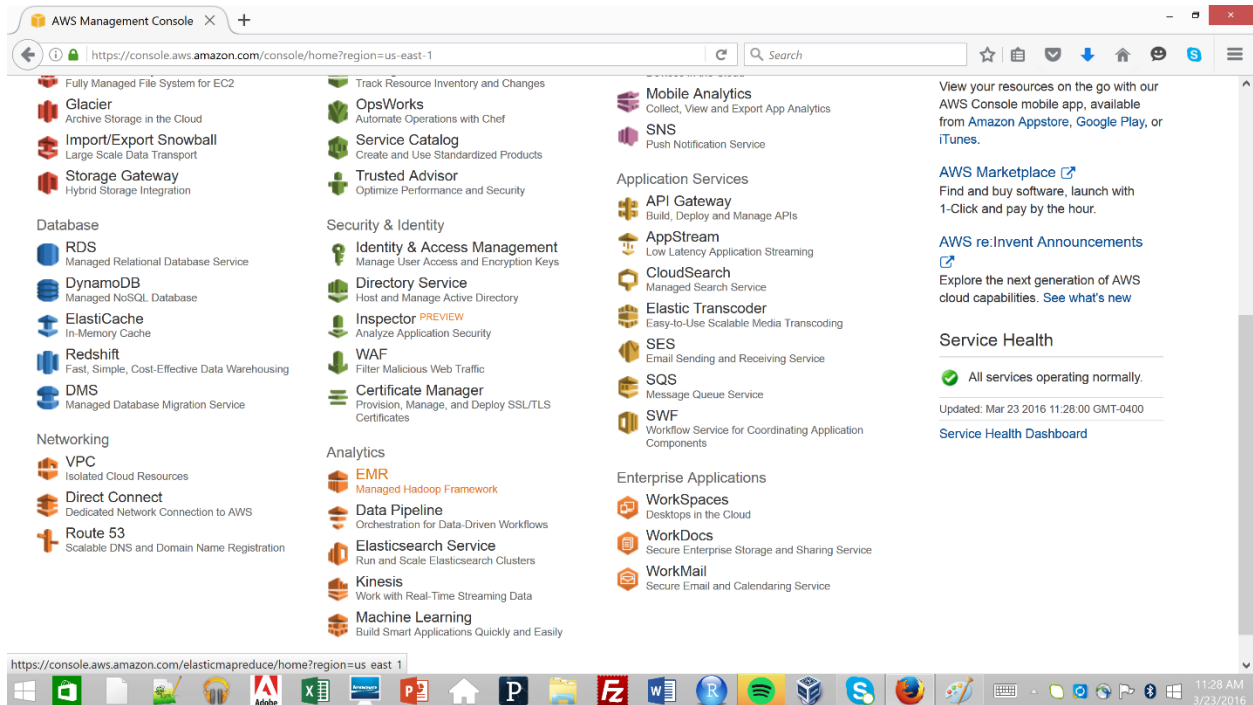


Once the bucket is created upload the `.jar` file and input files from local system or the VM. Click on Upload button and perform the upload of `wc.jar` and the input files. In this case, input files reside in **uspres** bucket file. Please refer to the following screenshot:

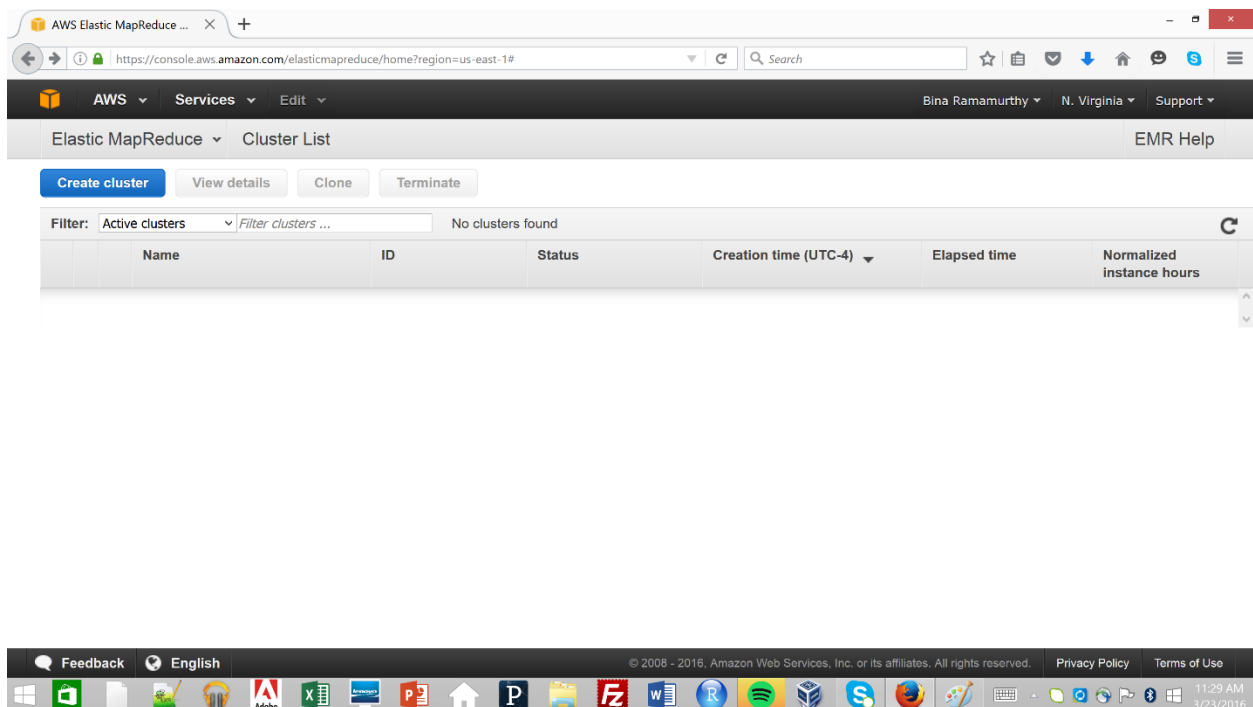


4 CONFIGURATIONS

For configuration click on the AWS Service Cube icon present on the top left of the page. Once the page opens click on **EMR** residing in **Analytics** section. Please refer to the following screenshot:



Here click on Create Cluster. Please refer to the following screenshot:



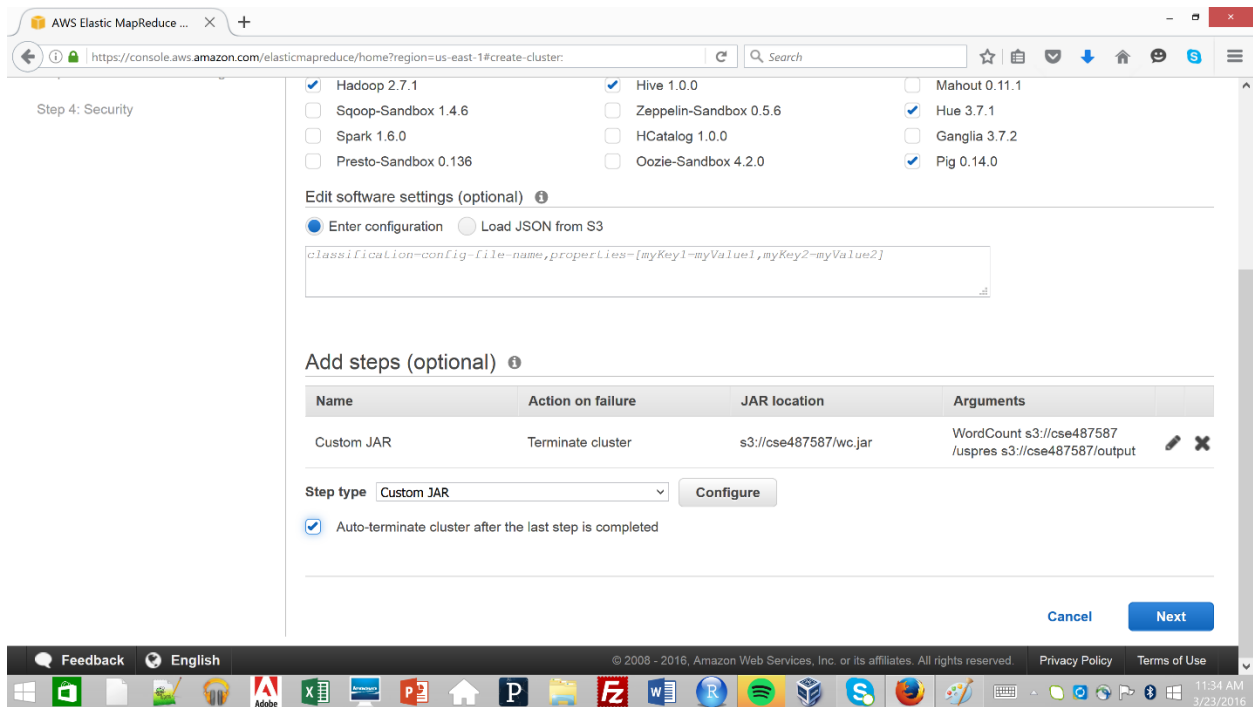
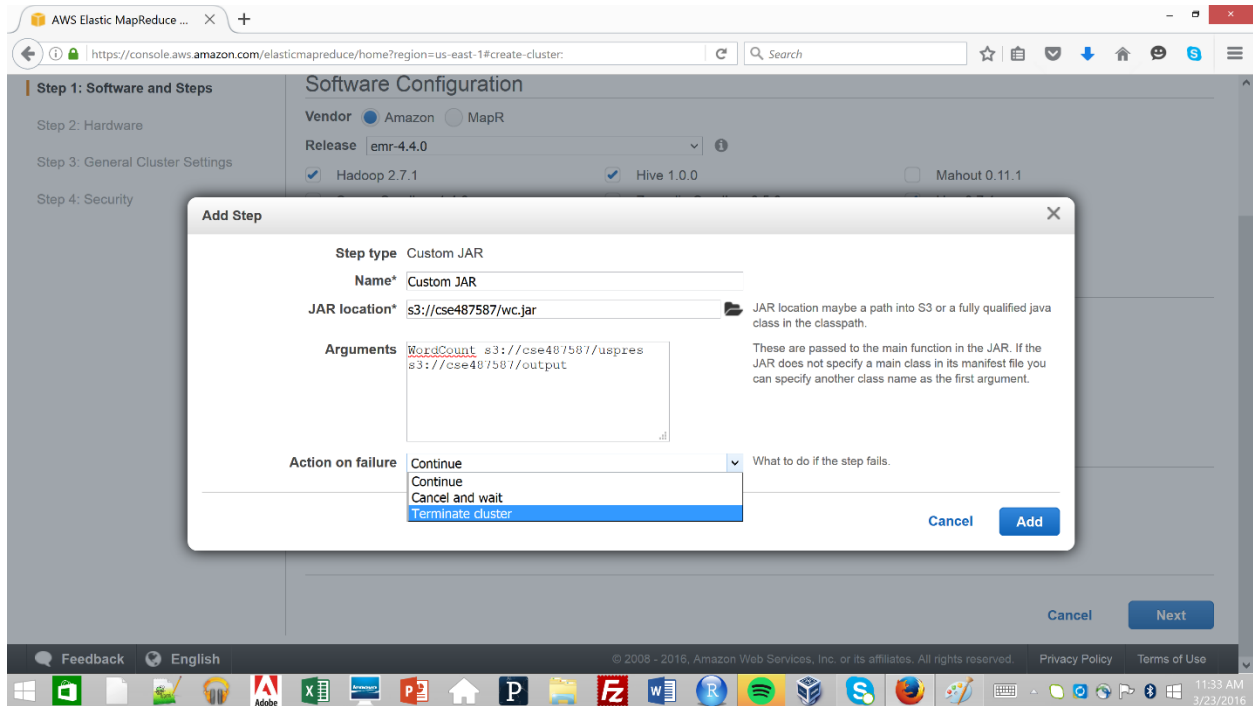
On opening a new web page, please click on the Advanced Options. Please refer to the following screenshot:

The screenshot shows the AWS Elastic MapReduce console's 'Create Cluster - Quick Options' page. The browser address bar shows the URL: `https://console.aws.amazon.com/elasticmapreduce/home?region=us-east-1#quick-create:`. The page has a navigation bar with 'AWS', 'Services', 'Edit', and user information 'Bina Ramamurthy', 'N. Virginia', and 'Support'. Below the navigation bar, there's a breadcrumb 'Elastic MapReduce > Create Cluster' and an 'EMR Help' link. The main heading is 'Create Cluster - Quick Options' with a link to 'Go to advanced options'. The 'General Configuration' section includes a 'Cluster name' field with 'My cluster', a checked 'Logging' checkbox, and an 'S3 folder' field with 's3://aws-logs-046917512226-us-east-1/elasticmapred'. The 'Launch mode' has 'Cluster' selected. The 'Software configuration' section shows 'Vendor' as 'Amazon', 'Release' as 'emr-4.4.0', and 'Applications' with 'All Applications' selected. The taskbar at the bottom shows various application icons and the system clock at 11:30 AM on 3/23/2016.

Once in the advanced options, go to section **Add Steps** and select **Step Type** to Custom JAR. Please refer to the following screenshot:

The screenshot shows the 'Software Configuration' page in the AWS Elastic MapReduce console. The left sidebar shows a progress bar with 'Step 1: Software and Steps' selected, followed by 'Step 2: Hardware', 'Step 3: General Cluster Settings', and 'Step 4: Security'. The main content area is titled 'Software Configuration'. It shows 'Vendor' as 'Amazon' and 'Release' as 'emr-4.4.0'. Under 'Hadoop 2.7.1', several checkboxes are visible: 'Hadoop 2.7.1' (checked), 'Sqoop-Sandbox 1.4.6', 'Spark 1.6.0', 'Presto-Sandbox 0.136', 'Hive 1.0.0' (checked), 'Zeppelin-Sandbox 0.5.6', 'HCatalog 1.0.0', 'Oozie-Sandbox 4.2.0', 'Mahout 0.11.1', 'Hue 3.7.1' (checked), 'Ganglia 3.7.2', and 'Pig 0.14.0' (checked). Below this is the 'Edit software settings (optional)' section with 'Enter configuration' selected and a text area containing 'classification-config-file-name,properties={myKey1-myValue1,myKey2-myValue2}'. The 'Add steps (optional)' section has a 'Step type' dropdown menu open, showing options: 'Custom JAR' (selected), 'Select a step', 'Auto-terminate', 'Streaming program', 'Hive program', 'Pig program', 'Spark application', and 'Custom JAR'. A 'Configure' button is next to the dropdown. At the bottom right, there are 'Cancel' and 'Next' buttons. The footer shows 'Feedback', 'English', copyright information '© 2008 - 2016, Amazon Web Services, Inc. or its affiliates. All rights reserved.', 'Privacy Policy', and 'Terms of Use'. The taskbar at the bottom shows various application icons and the system clock at 11:31 AM on 3/23/2016.

Later click on **Configure** button, it will open another popup. Fill in the details by providing the JAR location in S3 and corresponding arguments. Please ensure that **Action on Failure** has been selected to “*Terminate cluster*” [Why?]. Moreover, please check “Auto-terminate after the last step is complete”. Please refer to the following screenshots:



Click **Next** button and click **Next** button again skipping the Hardware configurations. Now after you reach **General Cluster Settings**, please uncheck **Debugging** and **Termination Protection**. Please refer to the following screenshot:

Step 1: Software and Steps
Step 2: Hardware
Step 3: General Cluster Settings
Step 4: Security

General Options

Cluster name

☒ Logging ⓘ
S3 folder

☐ Debugging ⓘ
☐ Termination protection ⓘ

Tags ⓘ

Key	Value (optional)
<input type="text" value="Add a key to create a tag"/>	<input type="text"/>

Additional Options

☐ EMRFS consistent view ⓘ

▶ Bootstrap Actions

Cancel Previous Next

After these steps click on **Next** button. In this **Security** section, please ensure the following configuration and then click on **Create Cluster**.

Create Cluster - Advanced Options [Go to quick options](#)

Step 1: Software and Steps
Step 2: Hardware
Step 3: General Cluster Settings
Step 4: Security

Security Options

EC2 key pair ⓘ

☒ Cluster visible to all IAM users in account ⓘ

Permissions ⓘ
☒ Default ☐ Custom
Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.

EMR role [EMR_DefaultRole](#) ⓘ

EC2 instance profile [EMR_EC2_DefaultRole](#) ⓘ

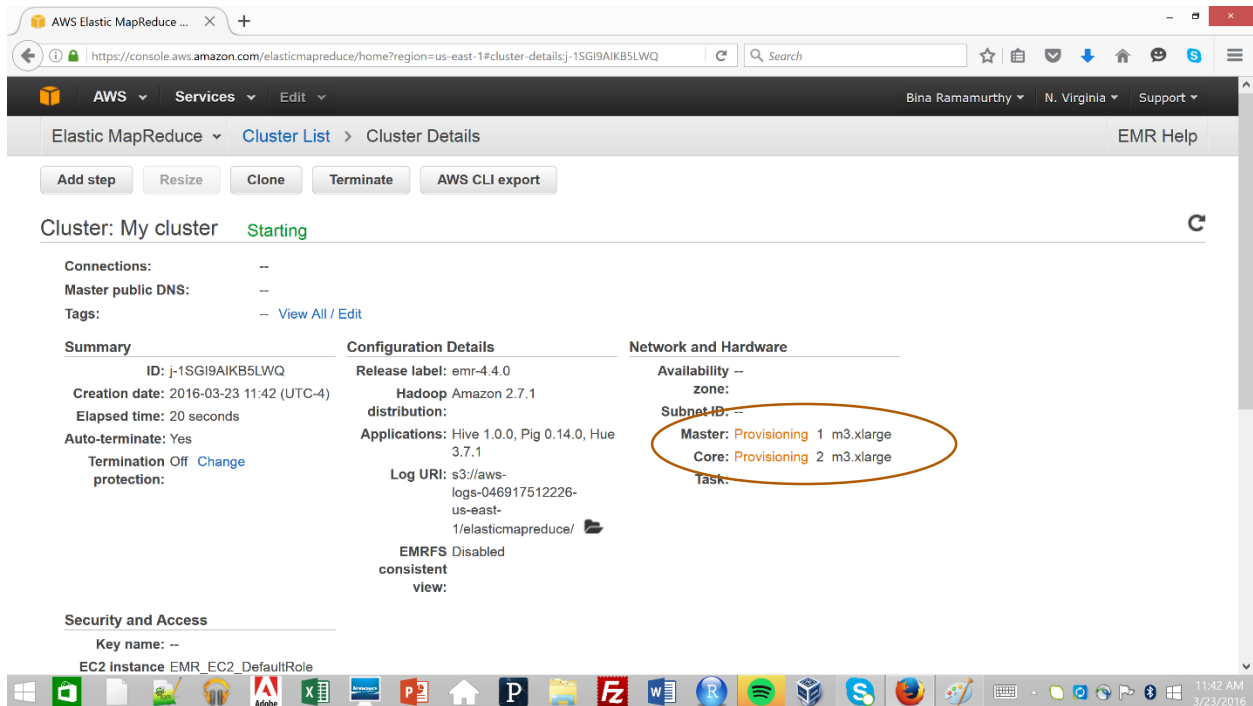
▶ EC2 Security Groups

▶ Encryption Options

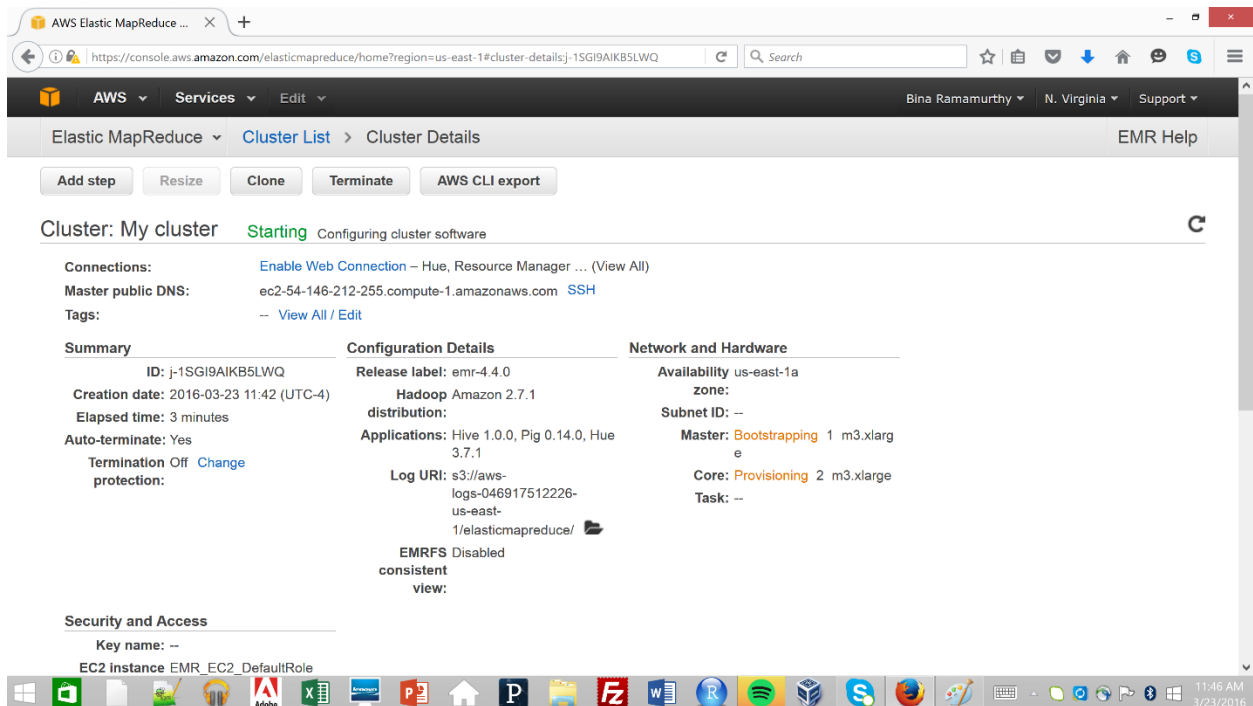
No EC2 key pair has been selected, so you will not be able to SSH to this cluster or connect to HUE (unless you are using a VPN). [Learn how to create an EC2 Key Pair.](#)

Cancel Previous Create cluster

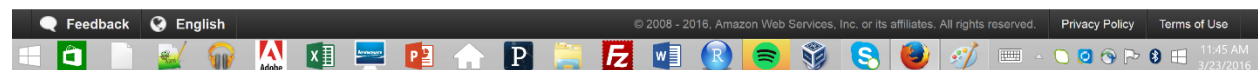
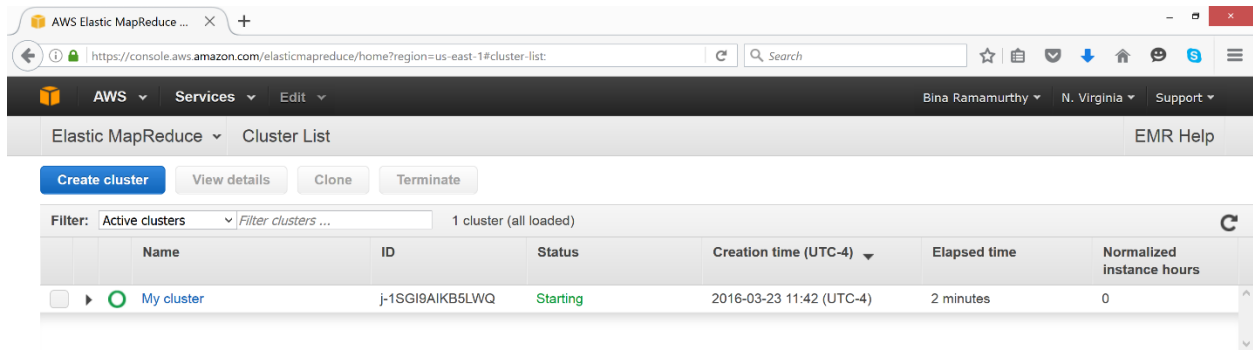
This will take you to a new page where details of the running cluster is shown. Please refer to the following image:



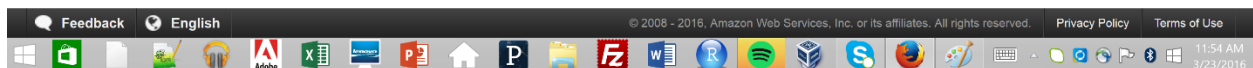
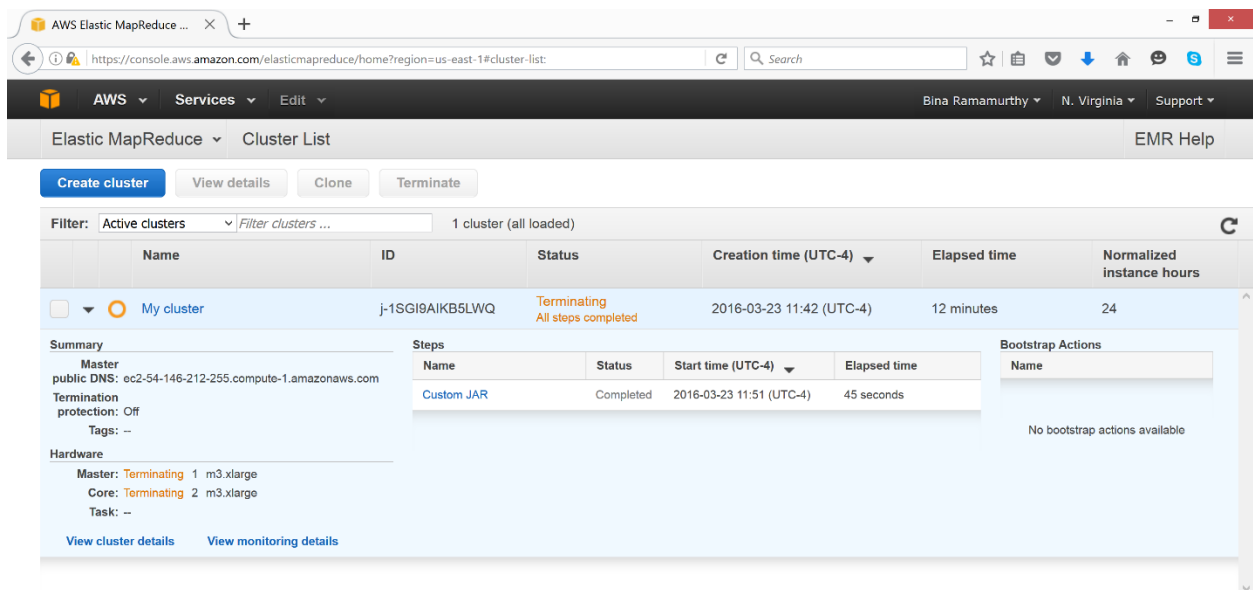
The subsequent processes typically take around 20 minutes. The next step involves the changing of the state of **Master** and **Core**, they change from *Provisioning* to *Bootstrap* to *Running* to *Terminated*. In the following screenshot there is a change recorded:



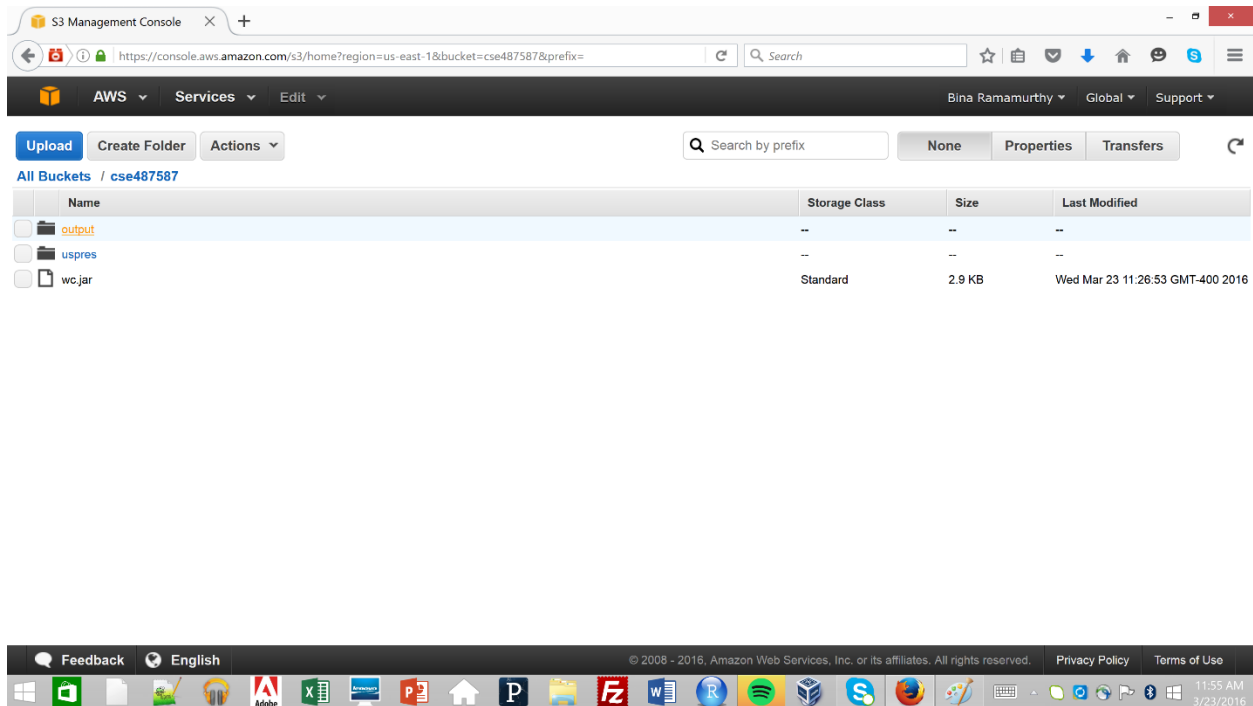
You can check the status of the active cluster by clicking on the **Cluster List**. The **Status** shows that the cluster has started running the application. Please refer to the following screenshot:



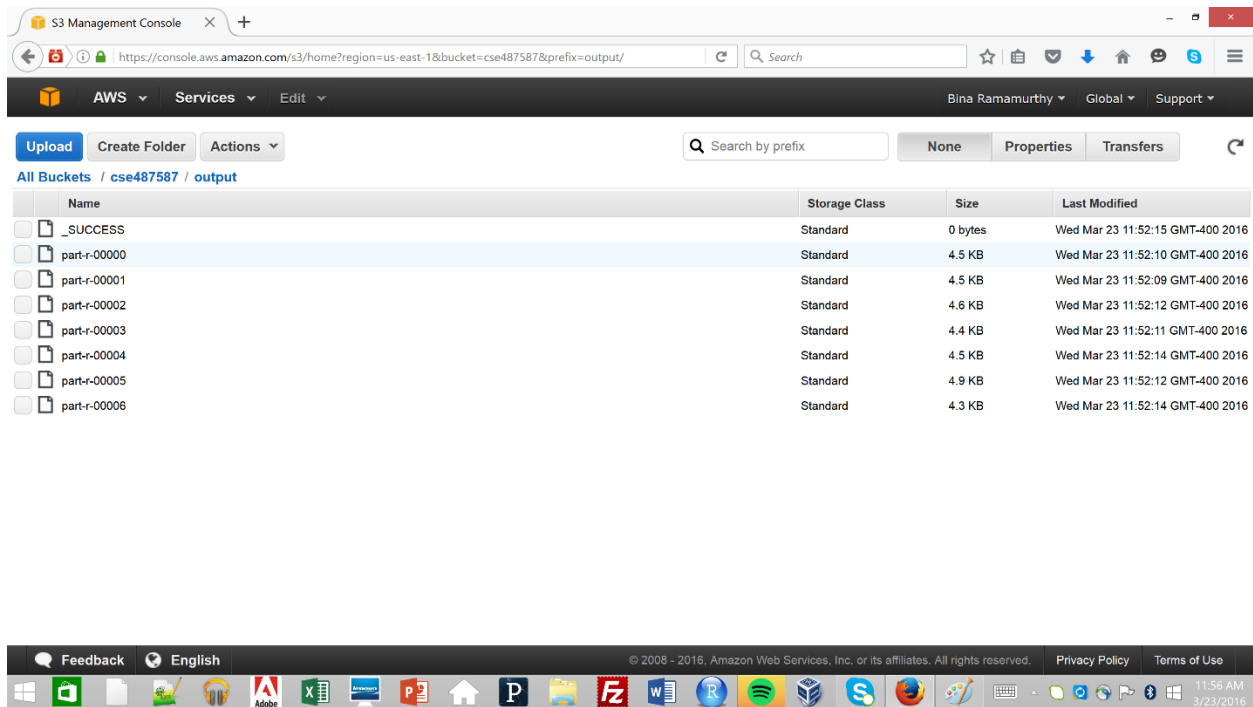
When the state changes to **Running** the color changes of the green icon changes to full green (●). This is not shown in the screenshot because the cluster terminated soon enough to notice these changes. Later, when the cluster completed the run, you will see something similar to the following screenshot:



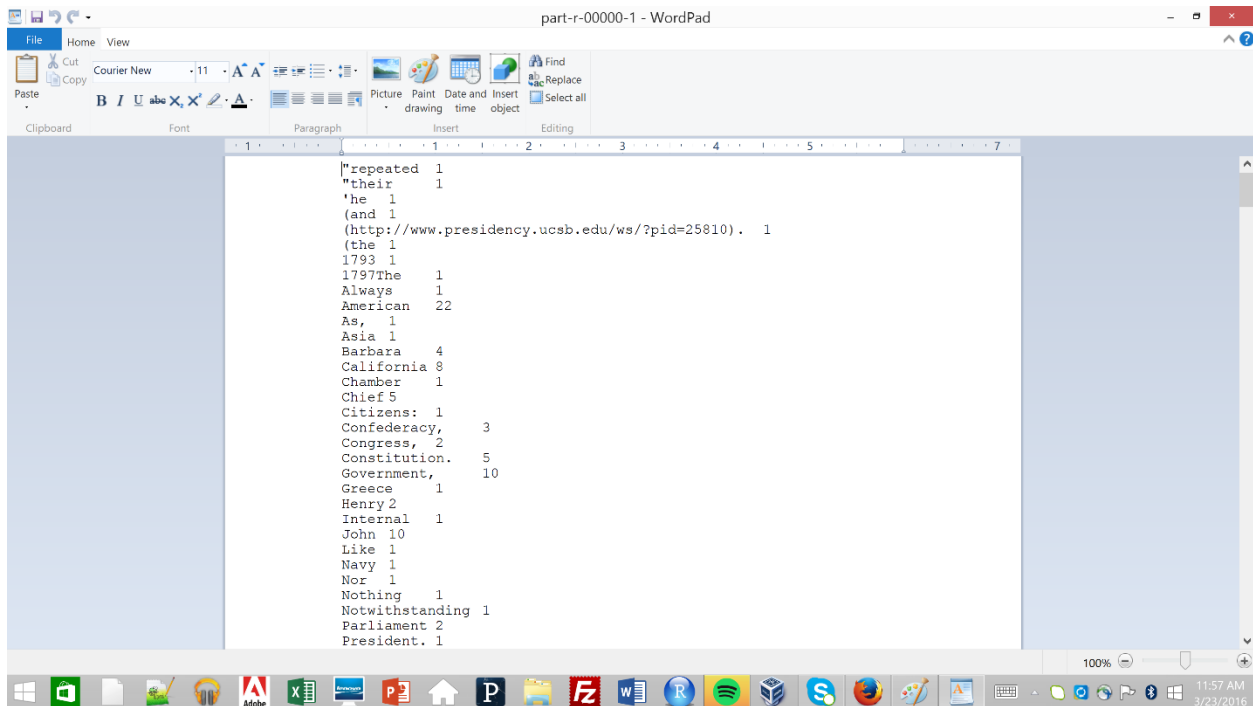
You can later check the output directories in the S3 Bucket which you created earlier. Please refer to the following screenshot:



Inside the **output** directory we can find the following list of files. Please refer to the following screenshot:



Once in this bucket you can check the details in the respective files. Here is a screenshot showing the the content in on these files:



5 NOTES AND WARNINGS

1. On creating an AWS Educate account using buffalo mail id, you would get benefits which include a free credit of \$100 which can be used to configure and run clusters. Each run will cost around \$1.
2. A bucket name must always be unique, it's a [rule](#).
3. If auto terminate is unchecked and Terminate Cluster is ignored then your cluster might run repeatedly. This might not be in your best interests as each run operation costs around \$1. It is better to continuously check the status of the active clusters, just like here below:

Name	ID	Status	Creation time (UTC-4)	Elapsed time	Normalized instance hours
My cluster	j-2EDTSWSNZKINQ	Terminated All steps completed	2016-03-22 22:14 (UTC-4)	16 minutes	24
My cluster	j-1U9338CWRUOXL	Terminated User request	2016-03-22 22:11 (UTC-4)	31 seconds	0
My cluster	j-2JN1EBWN921U	Terminated with errors Step failure	2016-03-22 21:51 (UTC-4)	16 minutes	24
My cluster	j-2QEOVA0124NIO4	Terminated with errors Step failure	2016-03-22 21:33 (UTC-4)	12 minutes	24
My cluster	j-13FEH0IHR5HRB	Terminated with errors Step failure	2016-03-22 20:59 (UTC-4)	10 minutes	24
My cluster	j-2RITG28MBMJTU	Terminated with errors Step failure	2016-03-22 10:19 (UTC-4)	14 minutes	24
My cluster	j-1QTFVH2AR2S1H	Terminated with errors Validation error	2016-03-22 10:13 (UTC-4)	33 seconds	0
My cluster	j-18A05NFGXB4NJ	Terminated User request	2016-03-11 09:12 (UTC-4)	4 days, 11 hours	2592