



INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY



HARVARD
School of Engineering
and Applied Sciences

Guide: Hadoop Cluster on AWS

Ignacio M. Llorente

v2.0 - 2 March 2020

Abstract

This is a screenshot document of how to run EMR Hadoop cluster and run MapReduce jobs on AWS environment.

Requirements

- **First you should have followed the Guide “First Access to AWS”.** It is assumed you already have an AWS account and a key pair, and you are familiar with the AWS EC2 environment.
- We strongly recommend cluster instances with at least 4 vCPUs (**m4.xlarge**) to be able to evaluate parallel implementation within each node.
- The files needed to do the exercises are available for download from **Canvas**.

Acknowledgments

The author is grateful for constructive comments and suggestions from David Sondak, Charles Liu, Matthew Holman, Keshavamurthy Indireskumar, Kar Tong Tan, Zudi Lin, Nick Stern, Dylan Randle, Hayoun Oh, Zhiying Xu and Zijie Zhao.



1. Launch Hadoop EMR cluster

- Go to the EMR dashboard (<https://console.aws.amazon.com/elasticmapreduce/home>) and click “Create cluster”. We recommend the following configuration
 - ClusterName: MyHadoop
 - Launch mode “Cluster”
 - Release: 5.29.0
 - Applications: Core Hadoop
 - Instance type: m4.xlarge
 - Number of Instances: 3
 - Key pair: course-key (or any other key you want to use, see Guide “First Access to AWS”)

General Configuration

Cluster name

Logging ⓘ

S3 folder

Launch mode Cluster ⓘ Step execution ⓘ

Software configuration

Release ⓘ

Applications Core Hadoop: Hadoop 2.8.5 with Ganglia 3.7.2, Hive 2.3.6, Hue 4.4.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.2

HBase: HBase 1.4.10 with Ganglia 3.7.2, Hadoop 2.8.5, Hive 2.3.6, Hue 4.4.0, Phoenix 4.14.3, and ZooKeeper 3.4.14

Presto: Presto 0.227 with Hadoop 2.8.5 HDFS and Hive 2.3.6 Metastore

Spark: Spark 2.4.4 on Hadoop 2.8.5 YARN with Ganglia 3.7.2 and Zeppelin 0.8.2

Use AWS Glue Data Catalog for table metadata ⓘ

Hardware configuration

Instance type ⓘ The selected instance type adds 64 GiB of GP2 EBS storage per instance by default. [Learn more](#)

Number of instances (1 master and 2 core nodes)



CS205: Computing Foundations for Computational Science, Spring 2020

- Click on “Create Cluster”

Clone Terminate AWS CLI export

Cluster: MyHadoop **Starting** Configuring cluster software

Summary Application history Monitoring Hardware Configurations Events Steps Bootstrap actions

Connections: [Enable Web Connection](#) – Hue, Ganglia, Resource Manager ... (View All)

Master public DNS: ec2-100-25-12-63.compute-1.amazonaws.com [SSH](#)

History service: --

Tags: -- [View All / Edit](#)

Summary	Configuration details
ID: j-1AV4CMVFSM36	Release label: emr-5.29.0
Creation date: 2020-03-03 18:57 (UTC+1)	Hadoop distribution: Amazon 2.8.5
Elapsed time: 2 minutes	Applications: Ganglia 3.7.2, Hive 2.3.6, Hue 4.4.0, Mahout 0.13.0, Pig 0.17.0, Tez 0.9.2
After last step completes: Cluster waits	Log URI: s3://aws-logs-196331178428-us-east-1/elasticmapreduce/
Termination protection: Off Change	EMRFS consistent view: Disabled
	Custom AMI ID: --

Network and hardware	Security and access
Availability zone: us-east-1a	Key name: course-key
Subnet ID: subnet-38252002	EC2 instance profile: EMR_EC2_DefaultRole
Master: Bootstrapping 1 m4.xlarge	EMR role: EMR_DefaultRole
Core: Provisioning 2 m4.xlarge	Visible to all users: All Change
Task: --	Security groups for sg-f02adb8f (ElasticMapReduce-Master: master)
	Security groups for sg-ee2adb91 (ElasticMapReduce-Core & Task: (ElasticMapReduce-slave))

- Wait for the cluster to be ready. The cluster is ready when its state is “Waiting” and the Master and Core under the **Networks and hardware** section are both in “Running” state

Amazon EMR

- Clusters
- Security configurations
- Block public access
- VPC subnets
- Events
- Notebooks
- Git repositories
- Help
- What's new

Clone Terminate AWS CLI export

Cluster: MyHadoop **Waiting** Cluster ready after last step completed.

Summary Application history Monitoring Hardware Configurations Events Steps Bootstrap actions

Connections: [Enable Web Connection](#) – Hue, Ganglia, Resource Manager ... (View All)

Master public DNS: ec2-100-25-12-63.compute-1.amazonaws.com [SSH](#)

History service: --

Tags: -- [View All / Edit](#)

Summary	Configuration details
ID: j-1AV4CMVFSM36	Release label: emr-5.29.0
Creation date: 2020-03-03 18:57 (UTC+1)	Hadoop distribution: Amazon 2.8.5
Elapsed time: 10 minutes	Applications: Ganglia 3.7.2, Hive 2.3.6, Hue 4.4.0, Mahout 0.13.0, Pig 0.17.0, Tez 0.9.2
After last step completes: Cluster waits	Log URI: s3://aws-logs-196331178428-us-east-1/elasticmapreduce/
Termination protection: Off Change	EMRFS consistent view: Disabled
	Custom AMI ID: --

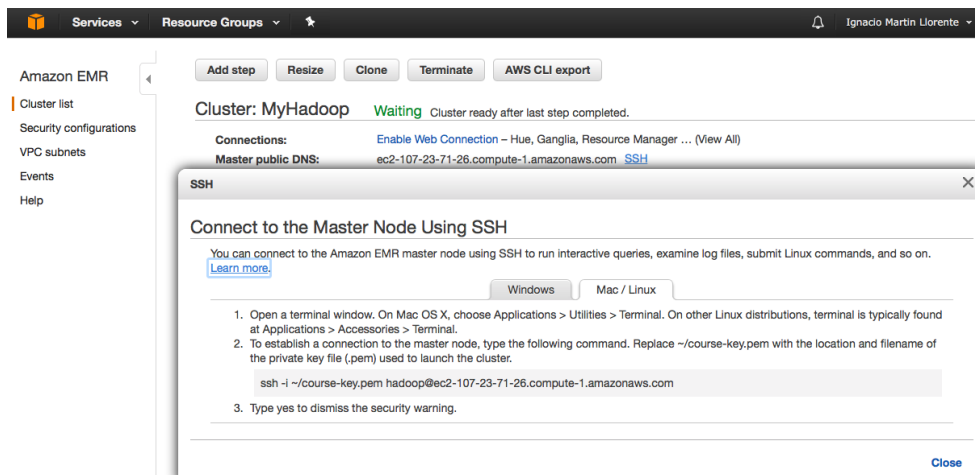
Network and hardware	Security and access
Availability zone: us-east-1a	Key name: course-key
Subnet ID: subnet-38252002	EC2 instance profile: EMR_EC2_DefaultRole
Master: Running 1 m4.xlarge	EMR role: EMR_DefaultRole
Core: Running 2 m4.xlarge	Visible to all users: All Change
Task: --	Security groups for sg-f02adb8f (ElasticMapReduce-Master: master)
	Security groups for sg-ee2adb91 (ElasticMapReduce-Core & Task: (ElasticMapReduce-slave))



2. Login to the cluster

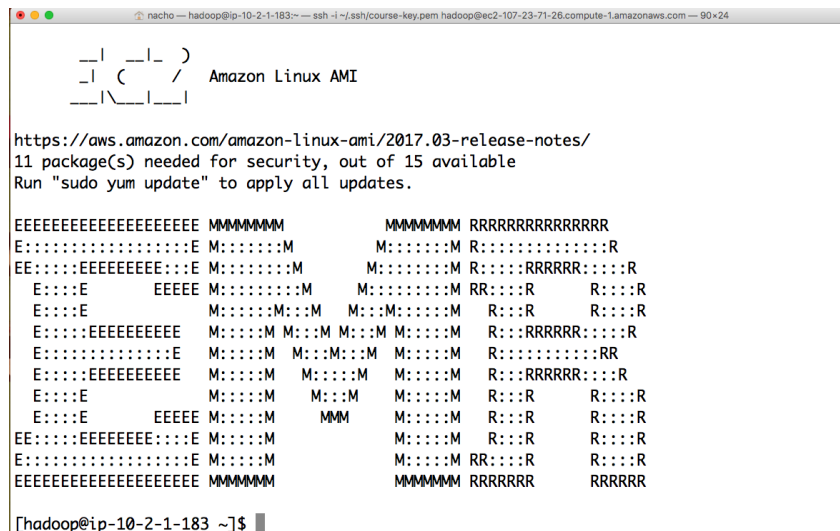
This section is for illustrative purposes to show how EMR is a Hadoop cluster automatically installed and configured on-demand on EC2 instances. You can skip this section to complete this guide because, as it is described in Section 3, you can submit basic MapReduce jobs from the AWS web interface

- Write down the “Master public DNS” and click on the SSH link next to it. The SSH link gives you the commands you might use to login to your cluster



- Most likely you will need to open port 22 to be able to login. Make sure that the security groups (firewalls) of the EMR cluster master node opens the port 22 to the outside world (see Guide “First Access to AWS”). Click the link to the security group next to **Security groups for Master**, click the Master security group and add an SSH rule with port 22 and source 0.0.0.0/0.
- SSH to the machine using the private key

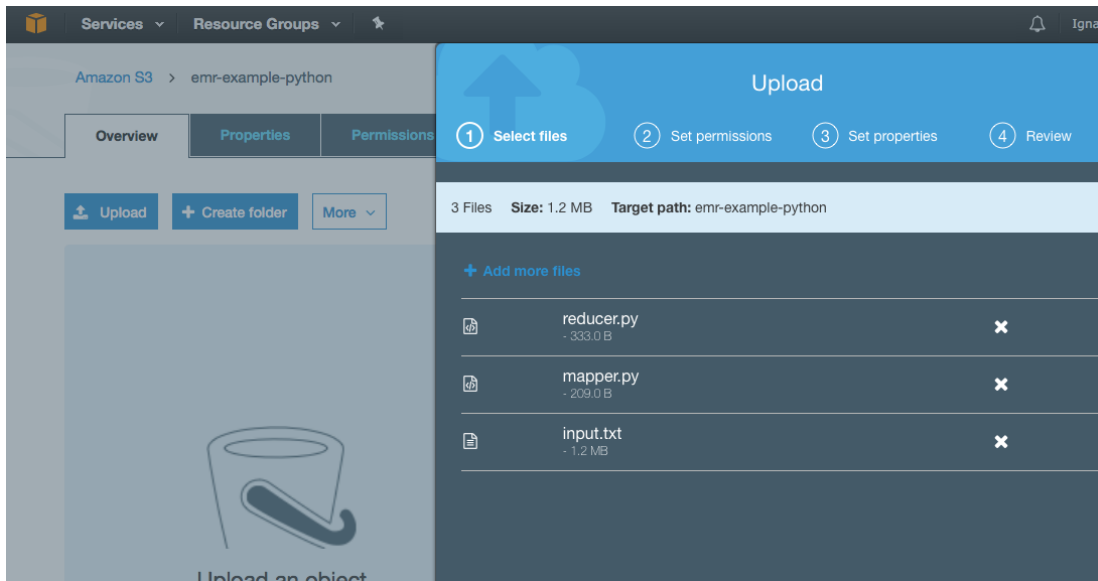
```
$ ssh -i your/course/ssh-key.pem hadoop@your-master-public-dns
```



3. Submit a MapReduce job

Hadoop Streaming is a utility that comes with Hadoop that enables you to develop MapReduce executables in languages other than Java. A Streaming application reads input from standard input and then runs a script or executable (called a mapper) against each input. The result from each of the inputs is saved locally on a Hadoop Distributed File System (HDFS) partition. After all the input is processed by the mapper, a second script or executable (called a reducer) processes the mapper results. The results from the reducer are sent to standard output.

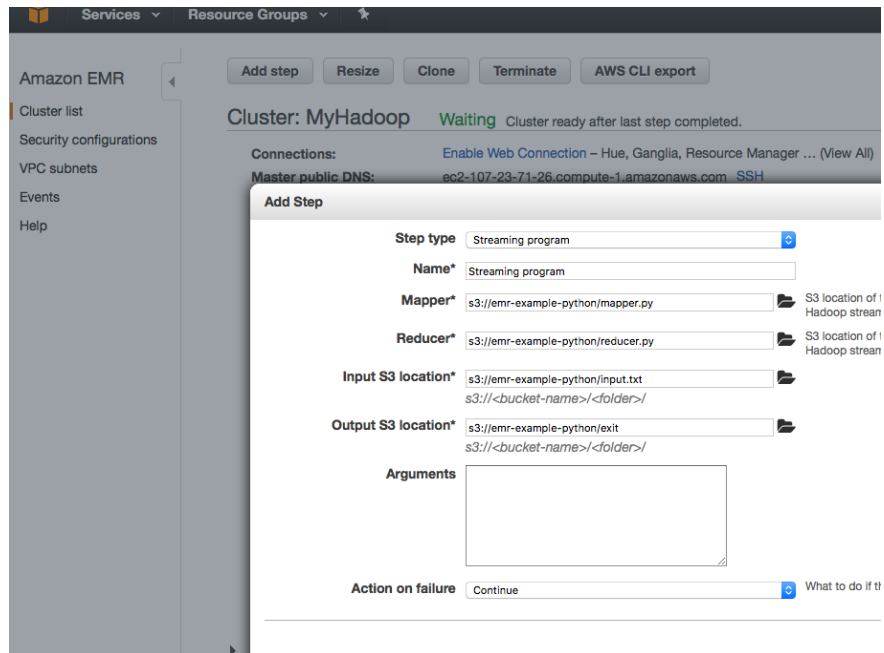
- Upload `mapper`, `reducer` and `input` files to a new S3 bucket. Create a S3 bucket, I named it `emr-example-python`. Remember this name should be unique. Moreover, because of Hadoop requirements, S3 bucket names used with Amazon EMR have the following constraints: must contain only lowercase letters, numbers, periods (`.`), and hyphens (`-`); and cannot end in numbers
 - Both mapper and reducer assume that lines are fed in through `sys.stdin`. Good sources of available text to play with are in Project Gutenberg.



- Go to the Hadoop cluster dashboard's **Steps** tab and click on "Add Step" with the following configuration
 - Step type: Streaming program
 - Name: MyHadoopJob
 - Mapper: Complete path to uploaded mapper
 - Reducer: Complete path to uploaded reducer
 - Input: Complete path to uploaded input
 - Output: Complete path to new folder to be created with the output (**it should not exist**)



CS205: Computing Foundations for Computational Science, Spring 2020



- Wait for the “step” to be “completed”
- After “completed” you can check the execution time in the `controller` log file

```
INFO total process run time: 72 seconds
```

- If the job is not successfully “completed”, you can check the logging files for further information
- Finally, check the results in the bucket, Hadoop creates one output file for each executed reducer task

				Viewing 1 to 8
<input type="checkbox"/>	Name	Last modified	Size	Storage class
<input type="checkbox"/>	_SUCCESS	Mar 3, 2020 7:18:26 PM GMT+0100	0 B	Standard
<input type="checkbox"/>	part-00000	Mar 3, 2020 7:18:16 PM GMT+0100	24.7 KB	Standard
<input type="checkbox"/>	part-00001	Mar 3, 2020 7:18:17 PM GMT+0100	25.4 KB	Standard
<input checked="" type="checkbox"/>	part-00002	Mar 3, 2020 7:18:25 PM GMT+0100	25.7 KB	Standard
<input type="checkbox"/>	part-00003	Mar 3, 2020 7:18:25 PM GMT+0100	25.0 KB	Standard
<input type="checkbox"/>	part-00004	Mar 3, 2020 7:18:24 PM GMT+0100	25.7 KB	Standard
<input type="checkbox"/>	part-00005	Mar 3, 2020 7:18:24 PM GMT+0100	25.8 KB	Standard
<input type="checkbox"/>	part-00006	Mar 3, 2020 7:18:21 PM GMT+0100	26.1 KB	Standard

Terminate the cluster when you are sure you are done for the day to avoid incurring charges