

# Creating an EMR Cluster

## 1. Create an Amazon EC2 Key Pair and PPK File

Amazon EMR uses an Amazon EC2 key pair to ensure that you alone have access to the instances that you launch. Amazon EMR requires a PPK file. The PPK file associated with this key pair is required to connect (via ssh) directly to the leader (master) node of the cluster.

To create an Amazon EC2 key pair:

1. Go to the Amazon EC2 console
2. In the Navigation pane (left side panel), click Key Pairs
3. On the Key Pairs page, click Create Key Pair
4. In the Create Key Pair dialog box, enter a name for your key pair, such as, EMRkeypair
5. If you are on Windows → Select the file format: PPK
6. If you are on Mac → Select the file format: PEM
7. Click Create key pair
8. Save the resulting file in a safe location. We will use this file to connect to the leader node

## 2. Launch an Amazon EMR Cluster

1. Services → EMR
2. Click **Create cluster**
3. In the **Name and applications** section, configure the following:
  - Cluster Name: give your cluster a name, such as, EMRCluster
  - Application bundle: select **Core Hadoop**.
4. In the **Cluster configuration** section, configure:
  - Under Uniform instance groups, select: Primary: **m4.large**, Core: **m4.large**, Task -1: **m4.large**
5. In the **Security configuration and EC2 key pair** section, configure:
  - Amazon EC2 key pair for SSH to the cluster: select the name of the key pair file you created
6. In the **Identity and Access Management (IAM) roles** section, configure:
  - Amazon EMR service role: under Service role, select EMR\_DefaultRole
  - EC2 instance profile for Amazon EMR: under Instance profile, select EMR\_EC2\_DefaultRole
7. Click **Create cluster** to launch your EMR Cluster.

The cluster will take approximately **fifteen minutes to launch**. The cluster will go through Starting, Bootstrapping, and Running states before the status changes to WAITING. Your cluster will be ready once the status changes to **WAITING**.

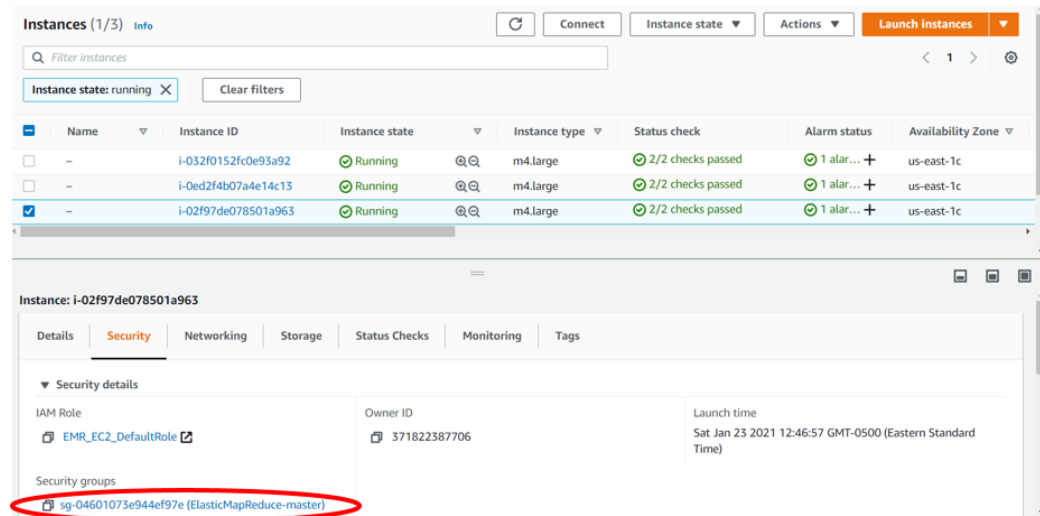
## 3. Upload files to S3

4. Services → S3
5. Create bucket
6. Give the bucket a name and select Create
7. Click in the name of your bucket and select Upload

8. Upload the files needed for this lab

#### 4. Allow your computer to connect to your EMR Cluster

1. Services → EC2
2. Click on **Instances (running)**
3. Find your Master:
  1. Click in any of the checkboxes of the instances you see running. You should only check ONE checkbox at the time.
  2. Once you have clicked in the checkbox, in the bottom portion, click on the Security tab.
  3. Under the Security Groups, you will see a description that says either [\(ElasticMapReduce-slave\)](#) or [\(ElasticMapReduce-master\)](#). If the instance you are checking says [\(ElasticMapReduce-slave\)](#), continue the process of checking instance by instance until you find the one that says: [\(ElasticMapReduce-master\)](#). Here is an example of what you will see once you find the master instance.



4. Once you have found the master instance, click where it says [ElasticMapReduce-master](#)
5. Click on Edit inbound rules
6. Scroll down to the bottom of the screen, and click on Add rule.
7. Select: **SSH for Type** and **Anywhere-IPv4 for Source**. You need to modify **BOTH** values. Leave the rest of the options in the default.
8. Select Save Rules

#### 5. Connect to your Leader (Master)

1. Services → EMR. You will see your cluster, similarly to this:

