

Up and Running with Hadoop

By :
Brahmbhatt Spandan H

Ways to install :

1. Brew (MacOS only).
2. Manual Install.
3. Using Cloudera CM (GUI Based)
4. Using Apache Ambari (GUI Based)
5. Using Amazon Elastic mapreduce.

1. Using brew (PseudoDistributed mode)

Run the following command :

1. Install Brew :

```
ruby -e "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/master/install)"
```

2. Install Hadoop :

```
brew install hadoop
```

3. Hadoop is install in

```
/usr/local/Cellar/hadoop
```

4. Edit hadoop-env.sh

Replace :

```
export HADOOP_OPTS="$HADOOP_OPTS -Djava.net.preferIPv4Stack=true"
```

with

```
export HADOOP_OPTS="$HADOOP_OPTS -Djava.net.preferIPv4Stack=true -Djava.security.krb5.realm= -Djava.security.krb5.kdc="
```

5. Edit Core-site.xml

```
<configuration>
```

```
<property>
```

```
<name>hadoop.tmp.dir</name>
```

```
<value>/usr/local/Cellar/hadoop/hdfs/tmp</value>
```

```
<description>A base for other temporary directories.</description>
```

```
</property>
```

```
<property>
```

```
<name>fs.default.name</name>
```

```
<value>hdfs://localhost:9000</value>
```

```
</property>
```

```
</configuration>
```

6. Edit mapred-site.xml

Rename `mapred-site.xml.template` to `mapred-site.xml` and then add

```
<configuration>
```

```
  <property>
```

```
    <name>mapred.job.tracker</name>
```

```
    <value>localhost:9010</value>
```

```
  </property>
```

```
</configuration>
```

7. Edit hdfs-site.xml

```
<configuration>
```

```
  <property>
```

```
    <name>dfs.replication</name>
```

```
    <value>1</value>
```

```
  </property>
```

```
</configuration>
```

8. Edit bash_profile file

~\$ cd ~ //Run this command to go to home.

~\$ vi .bash_profile

Add following lines after pressing “i” (i.e go in insert mode)

alias hstart="/usr/local/Cellar/hadoop/2.6.0/sbin/start-dfs.sh;/usr/local/Cellar/hadoop/2.6.0/sbin/start-yarn.sh"

alias hstop="/usr/local/Cellar/hadoop/2.6.0/sbin/stop-yarn.sh;/usr/local/Cellar/hadoop/2.6.0/sbin/stop-dfs.sh"

Now run following command :

~\$ source ~/.bash_profile

~\$ hadoop namenode -format

9. Before you run hadoop :

~\$ ssh-keygen -t rsa //Just press enter when asked to specify location to save file. Do not type anything.

// if it ask you to overwrite the file then do type no and skip step 9.

Enable Remote Login

“System Preferences” -> “Sharing”. Check “Remote Login”

~\$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys

10. Run hadoop :

~\$ hstart

~ \$ jps // use this command to see hadoop services are properly running.

~ \$ // Use this command to stop Hadoop.

Your Hadoop installation is at :

/usr/local/Cellar/hadoop

Reference : <http://amodernstory.com/2014/09/23/installing-hadoop-on-mac-osx-yosemite/>

2. Manual Install

Steps to install Hadoop 2.5.0 on AWS :

1. Run your instances from AWS EC2 console.
 - a. Log into EC2 console. Click on EC2 → Launch Instance
 - b. Select **Ubuntu Server 14.04 LTS (HVM), SSD Volume Type - ami-9a562df2** (Free tier) (General Purpose)
 - c. Change number of instance to 4 (1 to be Namenode.1-SNN and 2 will be Datanode).
 - d. Create a new key pair and save it.
2. Connect to each instance using ssh (in MacOS) or use WinSCP and Putty (in Windows)

```
~ $ ssh -i keypair.pem ubuntu@publicDNS
```
3. Change hostname to be public DNS with command - `~ $ sudo hostname publicDNS`.
4. Change host file with following command :

```
~ $sudo vi /etc/hosts
```

Make entry in this format: `ipaddress publicDNS`
5. Repeat step 2 and 3 for all instances.
6. Repeat following commands on each instance :


```
-----
~ $sudo apt-get update
~ $sudo add-apt-repository ppa:webupd8team/java
~ $sudo apt-get update && sudo apt-get install oracle-jdk7-installer
~ $java -version
~ $wget http://apache.petsads.us/hadoop/common/hadoop-2.5.0/hadoop-2.5.0.tar.gz
~ $tar -xzvf hadoop-2.5.0.tar.gz
~ $mv hadoop-2.5.0 hadoop
~ $rm hadoop-2.5.0.tar.gz
~ $ vi .bashrc
```

Add following lines use i key to go in insert mode and esc and :wq to quit:

```
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
```

Run the following commands now

```
~ $source ~/.bashrc
~ $echo $HADOOP_PREFIX
~ $chmod 644 .ssh/authorized_keys
~ $chmod 400 keypair.pem
~ $eval `ssh-agent`
~ $ssh-add keypair.pem
```

Check connection to each using :

```
~ $ssh ubuntu@publicDNS
~ $exit
```

Now on master node run following command :

```
~ $cd /etc/hadoop
~ $vi /etc/hadoop/hadoop-env.sh
    export JAVA_HOME=/usr/lib/jvm/java-7-oracle           //Add following line

~ $cd ~
~ $mkdir /tmp/hdfstmp
~ $cd /etc/hadoop
```

7. Run on Master instance

Edit following files :

vi core-site.xml

```
-----  
<property>  
<name>fs.default.name</name>  
<value>hdfs://ec2-54-173-155-242.compute-1.amazonaws.com:8020</value>  
</property>  
<property>  
<name>hadoop.tmp.dir</name>  
<value>/home/ubuntu/hdfstmp</value>  
</property>  
-----
```

vi hdfs.site.xml

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

```
~ $ mv mapred-site.xml.template mapred.site.xml
vi mapred-site.xml
```

```
-----
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>
-----
```

```
vi yarn-site.xml
```

```
-----
<property>
<name>yarn.nodemanager.aux-services</name>
<value>mapreduce_shuffle</value>
</property>
<property>
<name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
<value>org.apache.hadoop.mapred.ShuffleHandler</value>
</property>
<property>
<name>yarn.resourcemanager.resource-tracker.address</name>
<value>ec2-54-173-155-242.compute-1.amazonaws.com:8025</value>
</property>
<property>
<name>yarn.resourcemanager.scheduler.address</name>
<value>ec2-54-173-155-242.compute-1.amazonaws.com:8030</value>
</property>
<property>
<name>yarn.resourcemanager.address</name>
<value>ec2-54-173-155-242.compute-1.amazonaws.com:8040</value>
</property>
-----
```

8. Run command for all slaves to transfer the config files.

```
~ $scp hadoop-env.sh core-site.xml hdfs-site.xml yarn-site.xml mapred-site.xml  
ubuntu@publicDNS:/home/ubuntu/hadoop/etc/hadoop
```

9. On master update file to have slaves publicDNS:

```
~ $ vi slaves  
-----  
publicDNS  
-----
```

10. Start your cluster using :

```
~ $ bin/hadoop namenode -format  
~ $sbin/hadoop-daemon.sh start namenode  
~ $ sbin/hadoop-daemons.sh start datanode  
~ $sbin/yarn-daemon.sh start resourcemanager  
~ $sbin/yarn-daemons.sh start nodemanager
```

11. To check if HDFS is working:

```
vi sample.txt  
hdfs dfs -mkdir -p /usr/spandan  
hfs pdfs -copyFromLocal sample.txt /usr/spandan
```

12. Hadoop Web UI can be accessed from you system browser using

<http://masterpublicDNS:50070>

3. Using Cloudera CM5 (Not free)

Using Cloudera CM5, Hadoop cluster can be easily created. Although using CM5 is free, if you run it on Amazon AWS the system requirements for instances do not come in free tier and hence will incur a small cost. However if you run it on your own cluster of machines, it will be completely free.

Steps :

Create 4 instances from Amazon EC2 console. Select

Ubuntu Server 14.04 LTS (HVM), SSD Volume Type - ami-9a562df2 (Free tier) -- > Select m3.large which has 2 vCPU and 7.5 Gb Memory.

Save the keypair.

ssh into master node (anyone that you design as your master):

```
$ ssh -i your-key.pem ubuntu@ec2-xx-xx-xx-xx.compute-1.amazonaws.com
```

Run following commands :

```
$ wget http://archive.cloudera.com/cm4/installer/latest/cloudera-manager-installer.bin
```

```
$ chmod +x cloudera-manager-installer.bin
```

```
$ sudo ./cloudera-manager-installer.bin
```

You will be presented with a GUI in shell. Accept everything and let it install.

After you are complete with it. Open your browser and go to following link

<http://ec2-xx-xx-xx-xx.compute-1.amazonaws.com:7180>

Username and password is: admin

Welcome to Cloudera. Which edition do you want to deploy?

Upgrading to **Cloudera Enterprise** provides important features that help you manage and monitor your Hadoop clusters in mission-critical environments.

	Cloudera Standard	Cloudera Enterprise Trial	Cloudera Enterprise
License	Free	60 Days Post trial period, the product will continue to function as Cloudera Standard . Your cluster and your data will remain unaffected .	Annual Subscription(s) Upload License
Node Limit	Unlimited	Unlimited	Unlimited
CDH	✓	✓	✓
Core Cloudera Manager Features	✓	✓	✓
Advanced Cloudera Manager Features (click link below for details)		✓	✓
Backup & Disaster Recovery †		✓	✓
Cloudera Navigator †		✓	✓
Cloudera Support			✓

† Purchased as separate products.

For full list of features available in **Cloudera Standard** and **Cloudera Enterprise**, [click here](#).



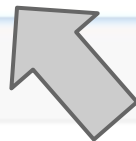
Continue

Specify hosts for your CDH cluster installation.

Cloudera recommends including Cloudera Manager server's host because it is often used for the Cloudera Management Service, and because this will enable health monitoring for that host.

Hint: Search for hostnames and/or IP addresses using [patterns](#).

SSH Port:



Enter your dns here. One per line.

Specify hosts for your CDH cluster installation.


Cloudera recommends including Cloudera Manager server's host because it is often used for the Cloudera Management Service, and because this will enable health monitoring for that host.

Hint: Search for hostnames and/or IP addresses using [patterns](#).

Every machine will come here.

11 hosts scanned, 11 running SSH.

New Search



Expanded Query	Hostname (FQDN)	IP Address	Currently Managed	Result
10.224.10.204	ip-10-224-10-204.us-west-2.compute.internal	10.224.10.204	No	✓ Host ready: 1,376 ms response time.
10.224.21.158	ip-10-224-21-158.us-west-2.compute.internal	10.224.21.158	No	✓ Host ready: 916 ms response time.
10.224.35.42	ip-10-224-35-42.us-west-2.compute.internal	10.224.35.42	No	✓ Host ready: 753 ms response time.
10.224.50.174	ip-10-224-50-174.us-west-2.compute.internal	10.224.50.174	No	✓ Host ready: 1,495 ms response time.
10.224.50.211	ip-10-224-50-211.us-west-2.compute.internal	10.224.50.211	No	✓ Host ready: 1,136 ms response time.
10.225.142.36	ip-10-225-142-36.us-west-2.compute.internal	10.225.142.36	No	✓ Host ready: 0 ms response time.
10.226.158.189	ip-10-226-158-189.us-west-2.compute.internal	10.226.158.189	No	✓ Host ready: 1,256 ms response time.
10.249.37.238	ip-10-249-37-238.us-west-2.compute.internal	10.249.37.238	No	✓ Host ready: 2,478 ms response time.
10.249.78.137	ip-10-249-78-137.us-west-2.compute.internal	10.249.78.137	No	✓ Host ready: 1,615 ms response time.
10.250.11.75	ip-10-250-11-75.us-west-2.compute.internal	10.250.11.75	No	✓ Host ready: 1,900 ms response time.

Cluster Installation

Select Repository

Cloudera Manager Parcels are the easiest way for Cloudera Manager to manage the software on your cluster, by automating the deployment and upgrade of service binaries. Electing not to use parcels will require you to manually upgrade packages on all hosts in your cluster when software updates are available, and will prevent you from using Cloudera Manager's rolling upgrade capabilities.

Choose Method:

- ☐ Use Packages
- ☒ Use Parcels (Recommended)
- > More Options



Use parcels

- SOLR ☒ SOLR-1.0.0-1.cdh4.3.0.p0.4
- ☐ None

Note: Solr is supported only on CDH 4.3 or later deployments.

- IMPALA ☒ IMPALA-1.1.1-1.p0.17
- ☐ None

Note: Impala is supported only on CDH 4.1 or later deployments.

- CDH ☒ CDH-4.4.0-1.cdh4.4.0.p0.39

Cluster Installation

Provide SSH login credentials.

Root access to your hosts is required to install the Cloudera packages. This installer will connect to your hosts via SSH and log in either directly as root or as another user with password-less sudo/pbrun privileges to become root.

Login to all hosts as:

☐ root

☒ Another User:

ubuntu

Use ubuntu as Login host

You may connect via password or public-key authentication for the user selected above.

Authentication Method: ☐ All hosts accept same password

☒ All hosts accept same private key

Private Key File:

Choose File

hadoop.pem

Use your keypair file.

Enter Passphrase:

Confirm Passphrase:

SSH Port:

22

Number of simultaneous

20

(Running a large number of installations at once can consume large amounts of network bandwidth and other system resources)

10 of 11 host(s) completed successfully.

✖ Installation failed on 1 host(s). [Uninstall Failed Hosts](#) [Retry Failed Hosts](#)

Hostname	IP Address	Progress	Status	
ip-10-224-10-204.us-west-2.compute.internal	10.224.10.204	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-224-21-158.us-west-2.compute.internal	10.224.21.158	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-224-35-42.us-west-2.compute.internal	10.224.35.42	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-224-50-174.us-west-2.compute.internal	10.224.50.174	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-224-50-211.us-west-2.compute.internal	10.224.50.211	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-225-142-36.us-west-2.compute.internal	10.225.142.36	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-226-158-189.us-west-2.compute.internal	10.226.158.189	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-249-37-238.us-west-2.compute.internal	10.249.37.238	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-249-78-137.us-west-2.compute.internal	10.249.78.137	<div></div>	✓ Installation completed successfully.	Details ⓘ
ip-10-250-11-75.us-west-2.compute.internal	10.250.11.75	<div></div>	✓ Installation completed successfully.	Details ⓘ

Choose the CDH4 services that you want to install on your cluster.

Choose a combination of services to install.

☒ **Core Hadoop**

HDFS, MapReduce, ZooKeeper, Oozie, Hive, and Hue

☐ **Core with Real-Time Delivery**

HDFS, MapReduce, ZooKeeper, HBase, Oozie, Hive, and Hue

☐ **Core with Real-Time Query**

HDFS, MapReduce, ZooKeeper, **Impala**, Oozie, Hive, and Hue

☐ **All Services**

HDFS, MapReduce, ZooKeeper, HBase, Impala, Oozie, Hive, Hue and Sqoop.

☐ **Custom Services**

Choose your own services. Services required by chosen services must also be selected. Note that Flume, Solr and Keystore Indexer services can be added after your initial cluster has been set up.

This wizard will also install the **Cloudera Management Services**. These are a set of components that enable monitoring, reporting, events, and alerts; these components require databases to store information, which will be configured on the next page.

☐ Include Cloudera Navigator

Database Setup

On this page you configure and test database connections. If using custom databases, create the databases first according to the **Installing and Configuring an External Database** section of the [Installation Guide](#).

When using the Embedded Database, passwords are auto generated. Please copy them down.

- ☒ Use Embedded Database
- ☐ Use Custom Databases

Hive

Database Host Name:	Database Type:	Database Name :	Username:	Password:
ip-10-225-142-36.us-west-2.compute.intern	PostgreSQL	hive	hive	IITSSVpTqx

Service Monitor

Currently assigned to run on ip-10-225-142-36.us-west-2.compute.internal.

Database Host Name:	Database Type:	Database Name :	Username:	Password:
ip-10-225-142-36.us-west-2.compute.intern	PostgreSQL	smon	smon	BdögmjDw45

Activity Monitor

Test Connection

Completed 3 of 17 steps.



Waiting for ZooKeeper Service to initialize
Finished waiting



Starting ZooKeeper Service
Service started successfully.



Checking if the name directories of the NameNode are empty. Formatting HDFS only if empty.
Successfully formatted NameNode.



Starting HDFS Service

Creating HDFS /tmp directory

Starting MapReduce Service

Creating Hive Metastore Database

Creating Hive Metastore Database Tables

Creating Hive user directory

Creating Hive warehouse directory

Starting Hive Service

Creating Oozie database

Installing Oozie ShareLib in HDFS

Congratulations!

The Hadoop services are installed, configured, and running on your cluster.

cloudera manager



Search by Service, Q

Support

admin

Home Services Hosts Activities Diagnose Audits Charts Reports Administration

Status All Health Issues 1 All Configuration Issues 2 All Recent Commands

September 26 2013, 4:45:29 PM UTC

Status



Cluster 1 - CDH4

Hosts		
hdfs1	1	1
hive1		
hue1		
mapreduce1		
oozie1		
zookeeper1	1	

Cloudera Management Services

mgmt1	
-------	--

Charts

30m 1h 2h 6h 12h 1d



Cluster 1 - CDH4

Cluster CPU

percent



Cluster Disk IO

bytes / second



Cluster Network IO

bytes / second



HDFS IO

bytes / second



Voila...!!!! You are all set.

5. Using Amazon Elastic Mapreduce

Amazon EMR is a paid service with very easy setup for MapReduce.

Create a JAR file with your mapper and reducer and upload it to S3 bucket.

Amazon Web Services

Compute



EC2

Virtual Servers in the Cloud



Lambda PREVIEW

Run Code in Response to Events

Storage & Content Delivery



S3

Scalable Storage in the Cloud



Storage Gateway

Integrates On-Premises IT Environments with Cloud Storage



Glacier

Archive Storage in the Cloud



CloudFront

Global Content Delivery Network

Database



RDS

MySQL, Postgres, Oracle, SQL Server, and Amazon Aurora



DynamoDB

Predictable and Scalable NoSQL Data Store



ElastiCache

In-Memory Cache



Redshift

Managed Petabyte-Scale Data Warehouse

Administration & Security



Directory Service

Managed Directories in the Cloud



Identity & Access Management

Access Control and Key Management



Trusted Advisor

AWS Cloud Optimization Expert



CloudTrail

User Activity and Change Tracking



Config PREVIEW

Resource Configurations and Inventory



CloudWatch

Resource and Application Monitoring

Deployment & Management



Elastic Beanstalk

AWS Application Container



OpsWorks

DevOps Application Management Service



CloudFormation

Templated AWS Resource Creation



CodeDeploy

Automated Deployments

Analytics



EMR

Managed Hadoop Framework

Application Services



SQS

Message Queue Service



SWF

Workflow Service for Coordinating Application Components



AppStream

Low Latency Application Streaming



Elastic Transcoder

Easy-to-use Scalable Media Transcoding



SES

Email Sending Service



CloudSearch

Managed Search Service

Mobile Services



Cognito

User Identity and App Data Synchronization



Mobile Analytics

Understand App Usage Data at Scale



SNS

Push Notification Service

Enterprise Applications



WorkSpaces

Desktops in the Cloud



WorkDocs

[AWS](#) ▾[Services](#) ▾[Edit](#) ▾[Spandan Brahmhatt](#) ▾[N. Virginia](#) ▾[Support](#) ▾[Elastic MapReduce](#) ▾[Cluster List](#)[EMR Help](#)[Create cluster](#)[Clone](#)[Terminate](#)

Filter:

[All clusters](#)[Filter clusters ...](#)

5 clusters (all loaded)



Name

ID

Status

Creation time (UTC-5) ▾

Elapsed time

Normalized
instance hours

Cluster Configuration

Cluster name

Termination protection ☒ Yes
☐ No

Logging ☒ Enabled

Log folder S3 location

s3://<bucket-name>/<folder>/

Debugging ☒ Enabled

Any Name you want

Prevents accidental termination of the cluster. To shut down the cluster, you must turn off termination protection. [Learn more](#)

Copy the cluster's log files automatically to Amazon S3. [Learn more](#)

Index logs to enable console debugging for your cluster. (requires logging). [Learn more](#)

Tags

i Optional: Add up to 10 tags to your EMR cluster. A tag consists of a case-sensitive key-value pair. Tags on EMR clusters are applied to the underlying EC2 instances. [Learn more](#) about tagging your Amazon EMR clusters.

Key

Value (optional)

Software Configuration

Hadoop distribution ☒ Amazon

Use Amazon's Hadoop distribution. [Learn more](#)

AMI version

3.3.1



Determines the base configuration of the instances in your cluster, including the Hadoop version. [Learn more](#)

☐ MapR

Use MapR's Hadoop distribution. [Learn more](#)

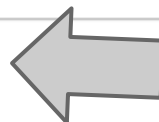
Applications to be installed	Version			
Hive	0.13.1			
Pig	0.12.0			
Hue	3.6.0			

Additional applications

Select an application



Configure and add



You can add additional Application from here

File System Configuration

i The [EMR File System \(EMRFS\)](#) and the Hadoop Distributed File System (HDFS) are both installed on your EMR cluster. HDFS stores data on an EMR cluster, while EMRFS allows EMR clusters to store data on S3. You can enable [server-side encryption](#) and [consistent view](#) for EMRFS below, or use a bootstrap action to configure additional settings for EMRFS.

Hardware Configuration

i Specify the [networking](#) and [hardware](#) configuration for your cluster. If you need more than 20 EC2 instances, [complete this form](#). [Request Spot instances](#) (unused EC2 capacity) to save money.

Network

Use a Virtual Private Cloud (VPC) to process sensitive data or connect to a private network. [Create a VPC](#)

EC2 Subnet

[Create a Subnet](#)

Type	Name	EC2 instance type	Count	Request spot	Bid price		
Master	<input type="text" value="Master instance group - 1"/>	<input type="text" value="m3.xlarge"/>	<input type="text" value="1"/>	<input type="checkbox"/>			
Core	<input type="text" value="Core instance group - 2"/>	<input type="text" value="m3.xlarge"/>	<input type="text" value="2"/>	<input type="checkbox"/>			
Task	<input type="text" value="Task instance group - 3"/>	<input type="text" value="m3.xlarge"/>	<input type="text" value="0"/>	<input type="checkbox"/>			

[Add task instance group](#)

Specify Number of Instances

Security and Access

EC2 key pair

Use an existing EC2 key pair to SSH into the master node of the Amazon EMR cluster. [Learn more](#)

IAM user access ☒ All other IAM users

☐ No other IAM users

Control the visibility of this cluster to other IAM users. [Learn more](#)

Bootstrap Actions

i Bootstrap actions are scripts that are executed during setup before Hadoop starts on every cluster node. You can use them to install additional software and customize your applications. [Learn more](#)

Bootstrap action type	Name	S3 location	Optional arguments		
-----------------------	------	-------------	--------------------	--	--

Add bootstrap action

Select a bootstrap action



Configure and add

Steps

i A step is a unit of work you submit to the cluster. A step might contain one or more Hadoop jobs, or contain instructions to install or configure an application. You can submit up to 256 steps to a cluster. [Learn more](#)

Name	Action on failure	JAR location	Arguments		
------	-------------------	--------------	-----------	--	--

Add step

Select a step

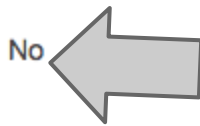


Configure and add

Auto-terminate

☐ Yes

☒ No



Select a step to add.

Automatically terminate cluster after the last step is completed.

Keep cluster running until you terminate it.

Steps

i A step is a unit of work you submit to the cluster. A step might contain one or more Hadoop jobs, or contain instructions to install or configure an application. You can submit up to 256 steps to a cluster. [Learn more](#)

Name	Action on failure	JAR location	Arguments		
------	-------------------	--------------	-----------	--	--

Add step

- ✓ Select a step
- Streaming program
- Hive program
- Pig program
- Impala program
- Custom JAR**

Auto-terminate

☒ No

Automatically terminate cluster after the last step is completed.

Keep cluster running until you terminate it.

i 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

Create cluster

Add Step

Step type

Custom JAR

Name*

Custom JAR

JAR location*

Arguments

Action on failure

Continue

What to do if the step fails.

Cancel

Add

JAR location maybe a path into S3 or a fully qualified java class in the classpath.

These are passed to the main function in the JAR. If the JAR does not specify a main class in its manifest file you can specify another class name as the first argument.

Click on Create Cluster and sit back and enjoy your Coffee!!!