

Big Data Infrastructure and Technologies

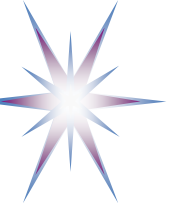
Practice 02 – Getting started with Hadoop

Creating and accessing AWS EMR Cluster



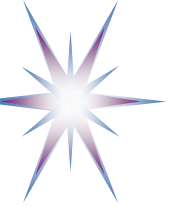
Installing and running AWS EMR cluster

- Prerequisite
- EMR cluster deployment steps
- Accessing Hadoop cluster with web browser
 - Configuring tunnel and Hadoop services
- <https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-gs.html>
- <https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-plan-instances-guidelines.html>



Prerequisites

- Create an Amazon S3 Bucket
 - To be referred as *s3://mybucket/MyHiveQueryResults*
 - Can be also created automatically during EMR creation
- Create EC2 keypair
 - Create and store during VM instance creation
 - Store created keypair with remembered name
 - Example: keypair2020bdit4da.pem => for Windows convert to .ppk with PuTTYgen
 - It will be prompted during EMR cluster creation
- For more information -
<https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-gs.html>



What to consider before starting

- Estimate Pricing for Amazon EMR and Amazon EC2 (On-Demand) - <https://aws.amazon.com/emr/pricing/>
 - The Amazon EMR price is in addition to the Amazon EC2 price (the price for the underlying servers) and Amazon EBS price (if attaching Amazon EBS volumes - billed per-second, with a one-minute minimum).
- **Save money with Reserved and Spot Instances**
- Estimate your monthly bill using the AWS Pricing Calculator <https://calculator.aws/>
 - Helps only with EC2/S3/RDS services, some Management services
 - Old pricing calculator <https://calculator.s3.amazonaws.com/index.html?s=EMR> (available until June 2020)



EMR Price Comparison table (2020-04)

Instance type	Amazon EC2 Price	Amazon EMR Price
m5.xlarge	\$0.192 per Hour	\$0.048 per Hour
m5.2xlarge	\$0.384 per Hour	\$0.096 per Hour
m5a.2xlarge	\$0.344 per Hour	\$0.086 per Hour
m4.large	\$0.10 per Hour	\$0.03 per Hour
m4.xlarge	\$0.20 per Hour	\$0.06 per Hour
m4.2xlarge	\$0.40 per Hour	\$0.12 per Hour
c5.xlarge	\$0.17 per Hour	\$0.043 per Hour
c4.large	\$0.10 per Hour	\$0.026 per Hour
c4.xlarge	\$0.199 per Hour	\$0.052 per Hour

The screenshot displays the Amazon EMR console interface. At the top, the navigation bar shows 'Services', 'Resource Groups', and a user profile. The left sidebar contains links to 'Amazon EMR', 'Clusters', 'Security configurations', 'Block public access', 'VPC subnets', 'Events', 'Notebooks', 'Git repositories', 'Help', and 'What's new'. The main content area shows details for cluster 'cluster02', which is in a 'Waiting' state. The 'Summary' tab is selected, displaying information such as the cluster ID (j-34U69132OC1XC), creation date (2020-04-15 19:45 UTC+2), elapsed time (4 hours, 27 minutes), and the master public DNS (ec2-54-235-53-87.compute-1.amazonaws.com). The 'Configuration details' section shows the release label (emr-5.29.0), Hadoop distribution (Amazon 2.8.5), and 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). The 'Network and hardware' section shows the availability zone (us-east-1b) and subnet (subnet-449d906a). The 'Security and access' section shows the key name (keypair2020bdt4da) and EC2 instance profile (EMR_EC2_DefaultRole). Below the cluster details, the 'Launch mode' is set to 'Cluster'. The 'Software configuration' section shows the release label (emr-5.23.0) and a list of applications: Core Hadoop (selected), HBase, Presto, and Spark.

Cluster: cluster02 Waiting Cluster ready after last step completed.

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

Master public DNS: [ec2-54-235-53-87.compute-1.amazonaws.com](#) [SSH](#)

History service: --

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

Summary

ID: j-34U69132OC1XC
 Creation date: 2020-04-15 19:45 (UTC+2)
 Elapsed time: 4 hours, 27 minutes
 After last step completes: Cluster waits
 Termination protection: [Off](#) [Change](#)

Configuration details

Release label: emr-5.29.0
 Hadoop Amazon 2.8.5 distribution:
 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
 Log URI: [s3://aws-logs-405547798672-us-east-1/elasticmapreduce/](#)

EMRFS consistent view: Disabled
 Custom AMI ID: --

Network and hardware

Availability zone: us-east-1b
 Subnet ID: [subnet-449d906a](#)
 Master: Running 1 m5.xlarge

Security and access

Key name: [keypair2020bdt4da](#)
 EC2 instance profile: [EMR_EC2_DefaultRole](#)

Launch mode ☒ Cluster ☐ Step execution

Software configuration

Release: [emr-5.23.0](#)

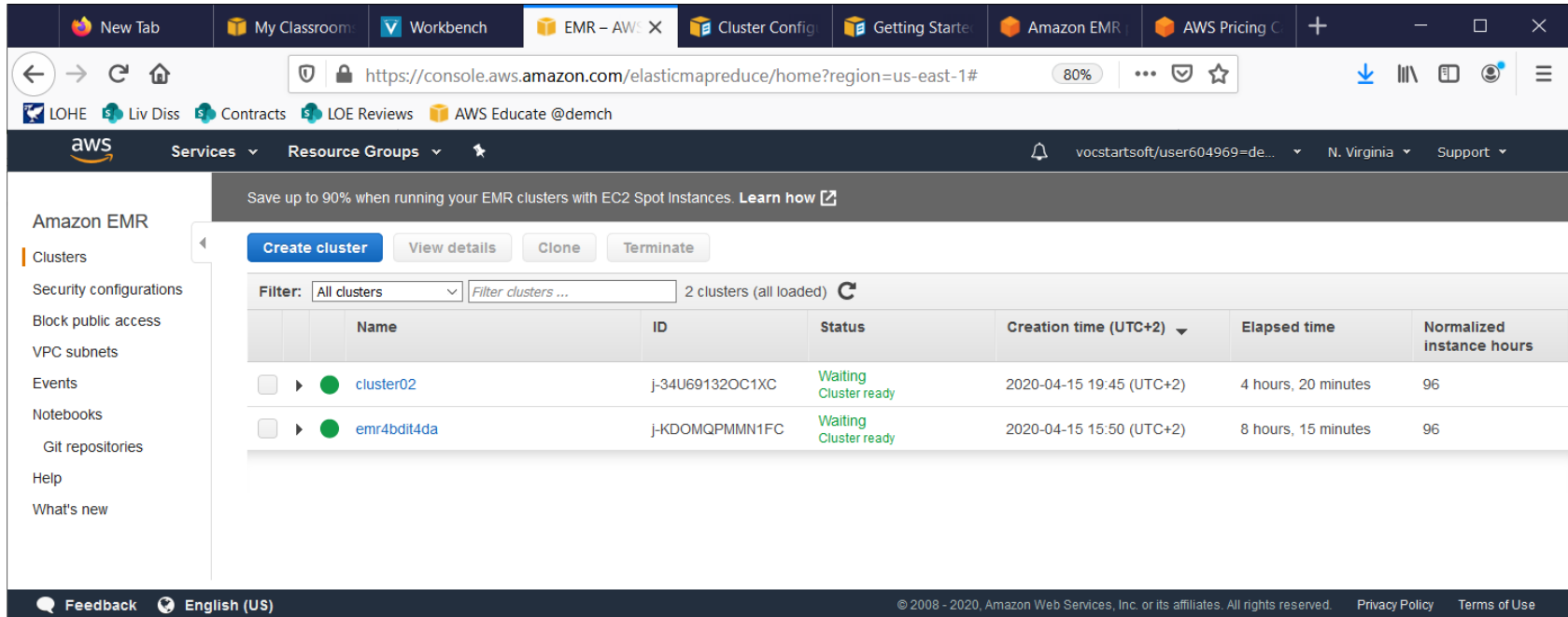
Applications

- ☒ Core Hadoop: Hadoop 2.8.5 with Ganglia 3.7.2, Hive 2.3.4, Hue 4.3.0, Mahout 0.13.0, Pig 0.17.0, and Tez 0.9.1
- ☐ HBase: HBase 1.4.9 with Ganglia 3.7.2, Hadoop 2.8.5, Hive 2.3.4, Hue 4.3.0, Phoenix 4.14.1, and ZooKeeper 3.4.13
- ☐ Presto: Presto 0.214 with Hadoop 2.8.5 HDFS and Hive 2.3.4 Metastore
- ☐ Spark: Spark 2.4.0 on Hadoop 2.8.5 YARN with

EMR cluster deployment steps

- Follow steps “Create Cluster” on Amazon EMR console
- Creation takes time of 10+ min
- Depends on the size of VM instance
- Check View Details for access details
- **Terminate cluster after finishing practice tasks**
- Use Clone to restore previously used configurations
- Note: All your tunnels and ports configuration will remain

Actual Cluster compute cost by total cores



Save up to 90% when running your EMR clusters with EC2 Spot Instances. [Learn how](#)

[Create cluster](#) [View details](#) [Clone](#) [Terminate](#)

Filter: All clusters 2 clusters (all loaded)

	Name	ID	Status	Creation time (UTC+2)	Elapsed time	Normalized instance hours
<input type="checkbox"/>	cluster02	j-34U69132OC1XC	Waiting Cluster ready	2020-04-15 19:45 (UTC+2)	4 hours, 20 minutes	96
<input type="checkbox"/>	emr4bdit4da	j-KDOMQPMMN1FC	Waiting Cluster ready	2020-04-15 15:50 (UTC+2)	8 hours, 15 minutes	96

- Note: Your charge time is calculated based on total number of cores used – See Normalised instance hours

Cluster configuration and Connection URL

The screenshot displays the AWS Management Console interface for an Amazon EMR cluster. The browser address bar shows the URL <https://console.aws.amazon.com/elasticmapreduce/home>. The left-hand navigation pane lists various services, with 'Security configurations' and 'Block public access' highlighted by a red circle. The main content area shows the details for a cluster named 'cluster02', which is currently in a 'Waiting' state. The cluster's status is described as 'Cluster ready after last step completed.' Below this, there are tabs for 'Summary', 'Application history', 'Monitoring', 'Hardware', 'Configurations', 'Events', 'Steps', and 'Bootstrap actions'. The 'Summary' tab is selected, displaying the cluster's ID (j-34U69132OC1XC), creation date (2020-04-15 19:45 (UTC+2)), elapsed time (4 hours, 27 minutes), and other details. The 'Configuration details' section shows the release label (emr-5.29.0), Hadoop distribution (Amazon 2.8.5), and 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). The 'Security and access' section at the bottom right shows the key name (keypair2020bdt4da) and the EC2 instance profile (EMR_EC2_DefaultRole), both of which are circled in red. The bottom of the console features a footer with 'Feedback', 'English (US)', and copyright information.

Amazon EMR

Clusters

Security configurations

Block public access

VPC subnets

Events

Notebooks

Git repositories

Help

What's new

Clone Terminate AWS CLI export

Cluster: cluster02 **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: ec2-54-235-53-87.compute-1.amazonaws.com [SSH](#)

History service: --

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

Summary

ID: j-34U69132OC1XC

Creation date: 2020-04-15 19:45 (UTC+2)

Elapsed time: 4 hours, 27 minutes

After last step Cluster waits completes:

Termination Off [Change](#) protection:

Configuration details

Release label: emr-5.29.0

Hadoop Amazon 2.8.5 distribution:

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

Log URI: s3://aws-logs-405547798672-us-east-1/elasticmapreduce/

EMRFS consistent Disabled view:

Custom AMI ID: --

Network and hardware

Availability zone: us-east-1b

Subnet ID: [subnet-449d906a](#)

Master: **Running** 1 m5.xlarge

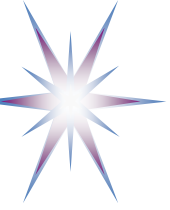
Security and access

Key name: keypair2020bdt4da

EC2 instance profile: EMR_EC2_DefaultRole

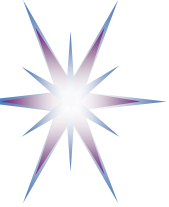
Feedback English (US)

© 2008 - 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use



View web interfaces to EMR cluster

- <https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-web-interfaces.html>
- Option 1: Set Up an SSH Tunnel to the Master Node Using Local Port Forwarding
 - Use tunneled port <http://localhost:8157/>
- Option 2, Part 1: Set Up an SSH Tunnel to the Master node using Static and Dynamic Port Forwarding
- Option 2, Part 2: Configure Proxy Settings to view web interface hosted on the Master node
- Access the Web Interfaces on the Master Node using web browser
- **NOTE: Check Master Node Security Group has open 22 port via Security group > Modify Rules**



Important EMR ports and services

Name of interface	URI
Ganglia	http://master-public-dns-name/ganglia/
Hadoop HDFS NameNode	https://master-public-dns-name:50470/
Hadoop HDFS DataNode	https://coretask-public-dns-name:50475/
HBase	http://master-public-dns-name:16010/
Hue	http://master-public-dns-name:8888/
JupyterHub	https://master-public-dns-name:9443/
Livy	http://master-public-dns-name:8998/
Spark HistoryServer	http://master-public-dns-name:18080/
Tez	http://master-public-dns-name:8080/tez-ui
YARN NodeManager	http://coretask-public-dns-name:8042/
YARN ResourceManager	http://master-public-dns-name:8088/
Zeppelin	http://master-public-dns-name:8890/

Connect to the Master Node Using SSH

The screenshot shows the AWS Management Console interface. At the top, there's a navigation bar with the AWS logo, 'Services', 'Resource Groups', and a user profile. Below this, the 'Amazon EMR' section is active, showing a cluster named 'emr4bdi4da' in a 'Waiting' state. An 'SSH' modal window is open in the foreground, titled 'Connect to the Master Node Using SSH'. The modal provides instructions on how to connect to the master node using SSH, including downloading PuTTY, starting it, and configuring the host name and private key file. The instructions are numbered 1 through 8. The modal also has tabs for 'Windows' and 'Mac / Linux'. At the bottom of the modal, there's a 'Close' button. Below the modal, the 'Network and hardware' and 'Security and access' sections are visible, showing the 'Availability zone: us-east-1a' and 'Key name: keypair2020bdi4da' respectively. The footer of the console contains a 'Feedback' button, 'English (US)' language selector, and copyright information.

Amazon EMR

Clusters

Cluster: `emr4bdi4da` Waiting Cluster ready after last step completed

SSH

Connect to the Master Node Using SSH

You can connect to the Amazon EMR master node using SSH to run interactive queries, examine log files, submit Linux commands, and so on. [Learn more](#).

Windows **Mac / Linux**

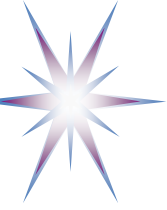
1. Download PuTTY.exe to your computer from:
<http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html>
2. Start PuTTY.
3. In the Category list, click Session.
4. In the Host Name field, type `hadoop@ec2-3-228-220-23.compute-1.amazonaws.com`
5. In the Category list, expand Connection > SSH, and then click Auth.
6. For Private key file for authentication, click Browse and select the private key file (`keypair2020bdi4da.ppk`) used to launch the cluster.
7. Click Open.
8. Click Yes to dismiss the security alert.

Close

Network and hardware **Security and access**

Availability zone: `us-east-1a` Key name: `keypair2020bdi4da`

Feedback English (US) © 2008 - 2020, Amazon Web Services, Inc. or its affiliates. All rights reserved. Privacy Policy Terms of Use



Connect to the Master Node Using SSH

(Instruction pops up when starting ER cluster)

1. Connect to the Amazon EMR master node using SSH to run interactive queries, examine log files, submit Linux commands, and so on.
2. Download PuTTY.exe to your computer from:
<http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html>
3. Start PuTTY.
4. In the Category list, click Session.
5. In the Host Name field, type **hadoop@ec2-3-228-220-23.compute-1.amazonaws.com**
6. In the Category list, expand Connection > SSH, and then click Auth.
7. For Private key file for authentication, click Browse and select the private key file (**keypair2020bdit4da.ppk**) used to launch the cluster.
8. Click Open.
9. Click Yes to dismiss the security alert.

NOTE:

Check Master Node Security Group has open 22 port via Security group
If not, add port 22 forwarding (0.0.0.0/0) via Modify Rules

Enabling Web Connection (via Tunnel configuration)

Enable Web Connection

Setup Web Connection

Hadoop, Ganglia, and other applications publish user interfaces as web sites hosted on the master node. For security reasons, these web sites are only available on the master node's local web server.

To reach the web interfaces, you must establish an SSH tunnel with the master node using either dynamic or local port forwarding. If you establish an SSH tunnel using dynamic port forwarding, you must also configure a proxy server to view the web interfaces.

Step 1: Open an SSH Tunnel to the Amazon EMR Master Node - [Learn more](#)

Windows | **Mac / Linux**

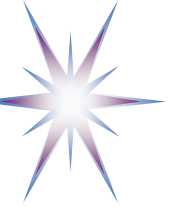
1. Download PuTTY.exe to your computer from:
<http://www.chiark.greenend.org.uk/~sgtatham/putty/download.html>
2. Start PuTTY.
3. In the Category list, click Session
4. In the Host Name field, type **hadoop@ec2-54-235-53-87.compute-1.amazonaws.com**
5. In the Category list, expand Connection > SSH > Auth
6. For Private key file for authentication, click Browse and select the private key file (**keypair2020bdit4da.ppk**) used to launch the cluster.
7. In the Category list, expand Connection > SSH, and then click Tunnels.
8. In the Source port field, type **8157** (a randomly chosen, unused local port).
9. Select the Dynamic and Auto options.
10. Leave the Destination field empty and click Add.
11. Click Open.
12. Click Yes to dismiss the security alert.

Step 2: Configure a proxy management tool - [Learn more](#)

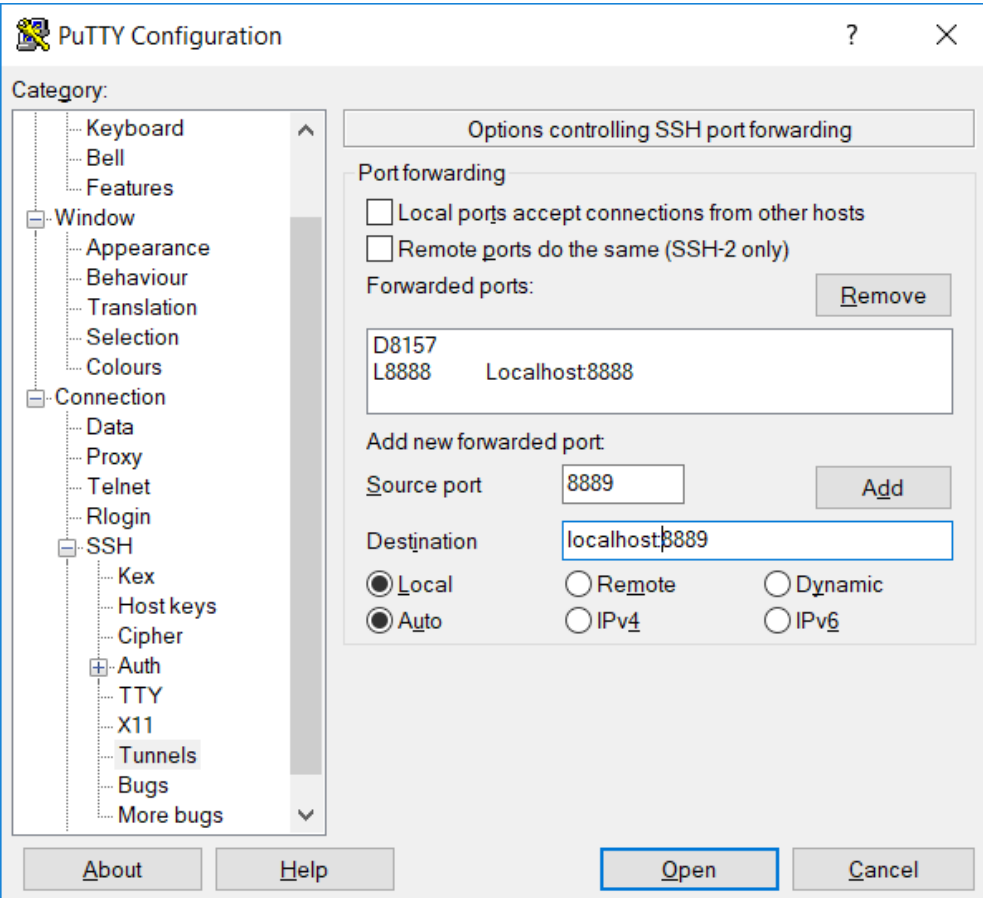
Chrome | **Firefox**

NOTE: Check Master Node Security Group has open 22 port via Security group > Modify Rules

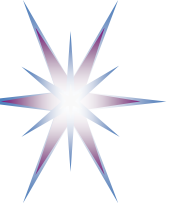
Close | Terms of Use



Configuring Tunnels



- Using PuTTY SSH client
- First, configure dynamic tunnel to EMR cluster
 - Port 8157 (or any random)
 - Using Source port and Dynamic radio button
- Configure other local ports
 - See EMR ports table
- This will create a secure tunnel by forwarding a port (the “destination port”) on the remote server to a port (the “source port”) on the local host (127.0.0.1 or localhost).



Connecting to Hue web interface

<http://localhost:8888/>

- Use <http://localhost:8888/> - or other port you assigned to browser access
 - Note: Because of security reasons using <http://master-public-dns-name:8888/> is not allowed
- At first login assign username and password
 - It is recommended to use username “hadoop” (the same as SSH login) to have access to HDFS space, otherwise you will get message of type:
 - Cannot access: /user/cluster02/. Note: you are a Hue admin but not a HDFS superuser, "hdfs" or part of HDFS supergroup, "hadoop".
- Hue web interface invocation takes 2-3 min
 - Started with a simple guide
- Use Query dropdown menu to select Query>Editor Pig, Hive, Java, Scala, etc
- Use top-left dropdown menu to select:
 - Apps: Editor, Scheduler
 - Browser: Documents, Files, Tables, Indexes, Jobs
- Use Browser > Files to upload your datasets or script (works for S3 directories)



The screenshot shows the Hue web interface. At the top, there's a navigation bar with the Hue logo and a search bar. Below it, a banner indicates that the user is accessing a non-optimized version of Hue and provides a link to switch to the optimized version. The main interface is divided into three panels: a left sidebar with navigation icons, a central query editor, and a right sidebar with 'Jobs' and 'Functions' tabs. The 'Tools' menu is open, displaying various options. The query editor shows a Hive query and its results, followed by a Pig query and its results. The right sidebar shows a table named 'default.students2'.

Tools Menu Options:

- Editor
- Dashboard
- Scheduler
- Pig
- Java
- Spark
- MapReduce
- Shell
- Sqoop 1
- Distcp
- Solr SQL
- Add more...

Query Editor Content:

Hive Query:

```
hive> select * from students2
nd(queryId=hive_20181017033638_20b811c3-409
3); Time taken: 0.001 seconds
```

Pig Query:

```
select * from students2

INSERT INTO TABLE students2
VALUES ('sara johns', 222222,
'bdit', '20170901', 'street2,
house2', 'master BA', 'master',
'NL'); ('Jaun Jansen', 3333333,
'bdit', '20160901', 'street3,
house3', 'master NatSc', 'master',
'UK')
```

Pig Results:

```
VALUES ('john smith', 111111,
'bdit', '20170901', 'street,
house', 'master MSc', 'master',
'NL')
```

Pig Query:

```
create1 CREATE TABLE students2
(last_name STRING, student_id
BIGINT, course STRING,
start_date DATE, address STRING
COMMENT 'Permanent home address
of the student',
highest_qualification STRING,
degree_type STRING, country
STRING) STORED AS SEQUENCEFILE
```

Right Sidebar:

Jobs

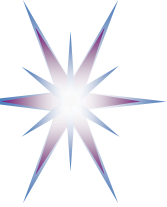
Functions

Tables

Filter...

default.students2

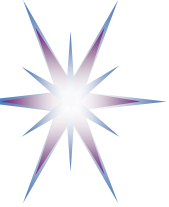
- Normally, you can upload file/data via File Browser in Hue
- However the main way of upload data is via S3 bucket and File Browser in Hue



Supported features in EMR

https://docs.amazonaws.cn/en_us/emr/latest/ReleaseGuide/emr-hue-supported-features.html

- Amazon S3 and Hadoop File System (HDFS) Browser
 - With the appropriate permissions, you can browse and move data between the ephemeral HDFS storage and S3 buckets belonging to your account.
 - By default, superusers in Hue can access all files that Amazon EMR IAM roles are allowed to access. Newly created users do not automatically have permissions to access the Amazon S3 filebrowser and must have the `filebrowser.s3_access` permissions enabled for their group.
- Hive - Run interactive queries on your data. This is also a useful way to prototype programmatic or batched querying.
- Pig - Run scripts on your data or issue interactive commands.
- Oozie - Create and monitor Oozie workflows.
- Metastore Manager - View and manipulate the contents of the Hive metastore (import/create, drop, and so on).
- Job browser - See the status of your submitted Hadoop jobs.
- User management - Manage Hue user accounts and integrate LDAP users with Hue.
- AWS Samples - There are several "ready-to-run" examples that process sample data from various AWS services using applications in Hue. When you log in to Hue, you are taken to the Hue Home application where the samples are pre-installed.
- Livy Server is supported only in Amazon EMR version 5.9.0 and later.
- To use the Hue Notebook for Spark, you must install Hue with Livy and Spark.
- **The Hue Dashboard is not supported.**
- **Files upload from browser is recently not supported. You need to do this via S3 bucket.**

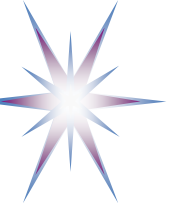


How to Get Data Into Amazon EMR

Amazon EMR provides several ways to get data onto a cluster.

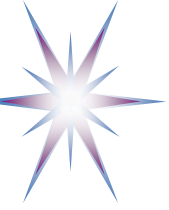
<https://docs.aws.amazon.com/emr/latest/ManagementGuide/emr-plan-get-data-in.html>

- The most common way is to upload the data to Amazon S3 and use the built-in features of Amazon EMR to load the data onto your cluster.
https://docs.cloudera.com/documentation/enterprise/6/6.3/topics/hue_use_s3_source_sink.html
- Use the Distributed Cache feature of Hadoop to transfer files from a distributed file system to the local file system.
- The implementation of Hive provided by Amazon EMR (Hive version 0.7.1.1 and later) includes functionality that you can use to import and export data between DynamoDB and an Amazon EMR cluster.



Error messages and problems

- HUE File Browser interface:
Cannot access: /user/cluster02/. Note: you are a Hue admin but not a HDFS superuser, "hdfs" or part of HDFS supergroup, "hadoop".
 - Hue superusers are different from hdfs superusers. With this error message, it complains that the admin and hdfs users you created through hue are not part to hdfs superuser.
 - You should be able to fix the issue by fixing the permission of the directories to the user/group on the cluster.
 - Simple solution assign/use username “hadoop” when accessing Hue
- **SSH connection: SSH connection timed out**
 - Check if cluster allows external access with port 22 (optionally HTTPS) and Master Node security group has open the same ports 22 and optionally 80 or 453
- Cannot connect to Hue via browser, site is not reachable
 - This is disable due to security reasons. Use SSH tunnel for this
- Hue interface is not responsive, doesn't browse files, tables, etc.
 - Logout and login to refresh Hue state
-



Additional Materials
