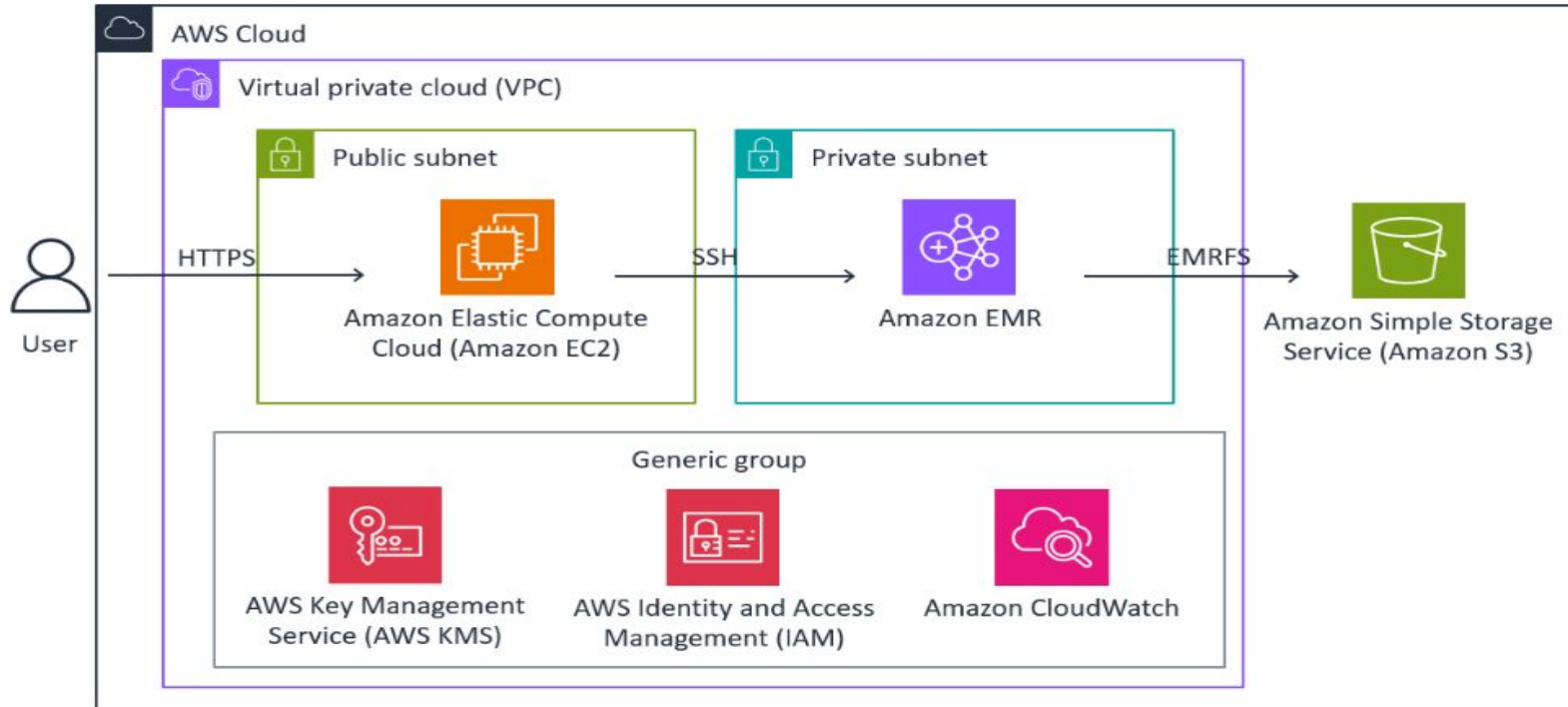


Load the sample data in Amazon Simple Storage Service (Amazon S3) and connect to the Amazon EMR cluster. Next, you create an Apache Hive table, load the data from Amazon S3, and run queries using HiveQL



Objectives

By the end of this lab, you should be able to do the following:

- Review how Amazon EMR and Apache Hive can be used together to ingest and query data
- Identify key components of an EMR cluster.
- Connect to an EMR cluster with SSH.
- Create a table using Apache Hive and load batch data from Amazon S3.
- Run queries using HiveQL

Trade_Date	Ticker	High	Low	Open	Close	Volume	Adj_Close
2020-01-02	aapl	75.1500015258789	73.79750061035156	74.05999755859375	75.0875015258789	135480400.0	74.2074661254888
2020-01-02	sq	64.05000305175781	62.95000076293945	62.9900016784668	63.83000183105469	5264700	63.83000183105469
2020-01-02	amzn	1898.010009765625	1864.1500244140625	1875.0	1898.010009765625	4029000	1898.010009765625
2020-01-02	ge	11.960000038146973	11.229999542236328	11.229999542236328	11.930000305175781	87421800.0	11.86101913452184
2020-01-02	m	17.270000457763672	16.389999389648438	17.18000030517578	16.520000457763672	26388100.0	15.8619861602782
2020-01-02	tsla	86.13999938964844	84.34200286865234	84.9000015258789	86.052001953125	47660500.0	86.052001953125
2020-01-02	msft	160.72999572753906	158.3300018310547	158.77999877929688	160.6199951171875	22622100.0	158.205764770508

Connect to the EMR leader node using Session Manager

use Session Manager, a capability of AWS Systems Manager, to connect to your EMR leader node

Get EMR cluster ID and export to the Environment.

```
export ID=$(aws emr list-clusters | jq '.Clusters[0].Id' | tr -d '"')
```

Use the ID to get the PublicDNS name of the EMR cluster

and export to the Environment.

```
export HOST=$(aws emr describe-cluster --cluster-id $ID | jq '.Cluster.MasterPublicDnsName' | tr -d '"')
```

SSH to the EMR cluster

```
ssh -i ~/EMRKey.pem hadoop@$HOST
```

Access your Amazon S3 data with Hive

In this task, you start an interactive Hive session with the leader node. Then, you create the Hive table from the CSV file on S3.

Note: For the purpose of this lab, you use the Open CSV SerDe hive driver to read the CSV files on Amazon S3 and create a table in Amazon EMR.

10. Command: To create a logging directory that is used by Hive, copy and paste the following commands into the SSH window:

```
sudo chown hadoop -R /var/log/hive
```

```
mkdir /var/log/hive/user/hadoop
```

The *hive.log* file is stored in this directory, which contains logs related to Hive.

11. Command: To connect to the Hive CLI, paste the following command into the SSH window:

```
Hive
```

You should be presented with a *hive>* prompt. It might take about 10 seconds to appear.

To create a table, copy and paste the following Hive statement in a text editor:

Warning: Replace *<dataBucket>* with the dataBucket value shown to the left of these instructions.

```
CREATE TABLE stockprice (  
  
  `Trade_Date` string,  
  
  `Ticker` string,  
  
  `High` double,  
  
  `Low` double,  
  
  `Open` double,  
  
  `Close` double,  
  
  `Volume` double,  
  
  `Adj_Close` double  
  
) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' WITH SERDEPROPERTIES (  
  
  "separatorChar" = ",",  
  
  "quoteChar"      = "\"",  
  
  "escapeChar"     = "\\\""  
  
)  
  
STORED AS TEXTFILE  
  
LOCATION 's3://<dataBucket>/data/'  
  
  TBLPROPERTIES ("skip.header.line.count"="1");
```

```
CREATE TABLE movies (  
  
  `year` int,  
  
  `title` string,  
  
  `directors_0` string,  
  
  `rating` string,  
  
  `genres_0` string,  
  
  `genres_1` string,  
  
  `rank` string,  
  
  `running_time_secs` string,  
  
  `actors_0` string,  
  
  `actors_1` string,  
  
  `actors_2` string,  
  
  `directors_1` string,  
  
  `directors_2` string  
  
)  
  
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'  
  
WITH SERDEPROPERTIES (  
  
  "separatorChar" = ",",  
  
  "quoteChar"    = "\"",  
  
  "escapeChar"   = "\\")
```