# APACHE IMPALA - LAB

## Group Members

Taha Rushain – 05811
Muhammad Siraj -  10255
Maria Zafar – 19842
M Sarmad Ali - 19828

# Table of Contents

# Apache Impala – Lab

## Overview

To overcome the slowness of Hive Queries, Cloudera offers a separate tool and that tool is what we call Impala. Impala queries run very faster than Hive Queries even after they are more or less same as Hive Queries. Impala is not built on MapReduce, Impala has its own execution engine.

## Selling points

- High-performance, low-latency SQL queries
- To share databases and tables between both Impala and hive it integrates very well with the Hive Metastore.
- It is Compatible with HiveQL Syntax
- We can easily integrate with Hbase database system and Amazon Simple Storage System.

We can perform interactive, ad-hoc and batch queries together in the Hadoop system, by using Impala's MPP (M-P-P) style execution along with other Hadoop processing MapReduce frameworks.

For querying analytic data it offers new possibilities. In addition, to query this type of data we can use exploratory data analysis and data discovery techniques.

## Components of a Impala Architecture

# Demo

## Installation

**Step 1:** Install Cloudera QuickStart Container using the following command

docker pull cloudera/quickstart:latest

```
(base) node@machine:~$ docker pull cloudera/quickstart:latest
latest: Pulling from cloudera/quickstart
Image docker.io/cloudera/quickstart:latest uses outdated schema1 manifest format.
1d00652ce734: Already exists
Digest: sha256:f91bee4cdfa2c92ea3652929a22f729d4d13fc838b00f120e630f91c941acb63
Status: Image is up to date for cloudera/quickstart:latest
docker.io/cloudera/quickstart:latest
```

**Step 2:** Download the dataset from Kaggle

We are using UK house pricing dataset that contains information on all registered property sales in England and Wales that are sold for full market value which can be downloaded from the following link

https://www.kaggle.com/hm-land-registry/uk-housing-prices-paid?select=price_paid_records.csv

Unzip the dataset and place the zip file along with the code files in a folder that would be used inside the container, in my case it is /home/node/Desktop/BDA

**Step 3:** Start the Docker image by executing the following command (Change the '/home/node/Desktop/BDA' to your folder location)

docker run --hostname=quickstart.cloudera --privileged=true -t -v /home/node/Desktop/BDA:/user/cloudera/shared -i -p 8888:8888 -p 7180:7180 cloudera/quickstart /usr/bin/docker-quickstart

**Here**

- /home/node/Desktop/BDA:/user/cloudera/shared: a mapping from the dataset folder on my Desktop to a shared folder within the cloudera docker (can now use the dataset within docker)
- 8888:8888: port 8888 of container is mapped to 8888 of my hosted OS (can now access hue on my Ubuntu through localhost:8888)
- 7180:7180: mapping for Cloudera manager, required to start the docker
- /usr/bin/docker-quickstart: starts the Cloudera docker
- --privileged=true: This is required for HBase, MySQL-backed Hive metastore, Hue, Oozie, Sentry, and Cloudera Manager
- -t : provide a bash shell after starting
- -i : we want to use the terminal

After the image is completely loaded, you will have a prompt inside the docker instance like this

```
Using CLASSPATH:         /usr/lib/bigtop-tomcat/bin/bootstrap.jar
Using CATALINA_PID:      /var/run/oozie/oozie.pid
Starting Solr server daemon:                         [  OK  ]
Using CATALINA_BASE:     /var/lib/solr/tomcat-deployment
Using CATALINA_HOME:     /usr/lib/solr/../bigtop-tomcat
Using CATALINA_TMPDIR:   /var/lib/solr/
Using JRE_HOME:          /usr/java/jdk1.7.0_67-cloudera
Using CLASSPATH:         /usr/lib/solr/../bigtop-tomcat/bin/bootstrap.jar
Using CATALINA_PID:      /var/run/solr/solr.pid
Started Impala Catalog Server (catalogd) :           [  OK  ]
Started Impala Server (impalad):                     [  OK  ]
[root@quickstart /]#
```

Step 4: Execute the following command to start cloudera manager

/home/cloudera/cloudera-manager –express

```
Success! You can now log into Cloudera Manager from the QuickStart VM's browser:

    http://quickstart.cloudera:7180

    Username: cloudera
    Password: cloudera

[root@quickstart /]#
```

Step 5: Execute the following commands (one by one) to sync your Docker time with NTP

service ntpd restart
ntpdate -u 0.centos.pool.ntp.org
hwclock –systohc
service ntpd restart

```
[root@quickstart /]# service ntpd restart
Shutting down ntpd:                                  [  OK  ]
Starting ntpd:                                       [  OK  ]
[root@quickstart /]# ntpdate -u 0.centos.pool.ntp.org
22 May 23:26:11 ntpdate[13380]: adjust time server 162.159.200.1 offset 0.052902 sec
[root@quickstart /]# hwclock --systohc
[root@quickstart /]# service ntpd restart
Shutting down ntpd:                                  [  OK  ]
Starting ntpd:                                       [  OK  ]
[root@quickstart /]#
```

**Step 6:** Open cloudera manager at the following url with cloudera as username and password http://localhost:7180/

**Step 7:** Press the down arrow button next to 'Cloudera Quickstart' title, press Start to start all cloudera service and select start on the prompt as well, as in the given screenshot

**Step 8:** After that, we need to install CDH5 parcel, go the following Hosts > Parcels in the top menu



**Step 9:** Download the CDH 5 parcel by clicking the Download button in the second row.



# Projec Setup

**Step 10:** We would use Hue to upload to interface with HDFS and upload our dataset. Open the following url and use username/password 'cloudera'
http://localhost:8888/



**Step 11:** Open File Browser, create a folder 'impala-input'

## Create Directory     ✕

Directory Name    impala-input

Cancel    **Create**

**Step 12:** Now open the 'impala-input' folder in Hue and upload out dataset.



## Execution

**Step 13:** Open Impala editor in Hue



**Step 14:** Execute the following query to create DB and refresh

**Step 15:** Execute the following query to create the table for our dataset



CREATE TABLE uk_houses.price_paid (

 `Transaction_unique_identifier` STRING,

 `Price` STRING,

 `Date_of_Transfer` STRING,

 `Property_Type` STRING,

 `Old_New` STRING,

 `Duration` STRING,

 `Town_City` TINYINT,

 `District` STRING,

 `County` STRING,

 `PPDCategory_Type` STRING,

 `Record_Status_monthly_file_only` STRING

 )

ROW FORMAT DELIMITED FIELDS TERMINATED BY ","

LOCATION "hdfs:///user/cloudera/impala-input/";

**Step 16:** Check the table sample



**Step 17:** Execute the following query to find the number of properties by Property Type



```
select property_type, count(1) '# of Properties'
from uk_houses.price_paid
where property_type <> 'Property Type'
group by property_type;
```

| | property_type | # of properties |
|---|---|---|
| 1 | D | 5170327 |
| 2 | F | 4083424 |
| 3 | O | 100568 |
| 4 | S | 6216218 |
| 5 | T | 6918811 |