**Part 2 - Create a Virtual Warehouse and Run Queries [45 minutes]**

Overview: What is Cloudera Data Warehouse?

We will explore features of Cloudera Data Warehouse (CDW) by performing some data exploration and create dashboards to share our results to a wider audience
We will be taking a look at a generated data set from a mock airline company containing flights information from its fleet of aircraft.

A virtual warehouse represents virtual compute resources to access data that is stored in a database catalog. This lets you create or destroy compute resources, auto-scale, or separate resources across different workloads, all running on the same underlying data.

CDW let's you choose from a set of default resources based on your predicted workload as well as give you fine grained control over autoscaling and timeout features so you can fine tune your system to be most cost effective.

Purpose: Create a virtual warehouse and run queries, answering the questions below:
- What are the top 5 visited destinations by year from (1995-2008)?
- What are the top 10 routes (origin and dest) that have seen maximum diversions?
- Which three months have seen the most number of cancellation due to bad weather?

1) Open CDP, using the "admin" user within the Test Drive link.

Your link should look something like (remember click the link in your email not the link below)
http://login.trycdp.com/auth/realms/trycdp-trialxx/protocol/saml/clients/samlclient?tn=trialxx_admin@trycdp.com&p=X
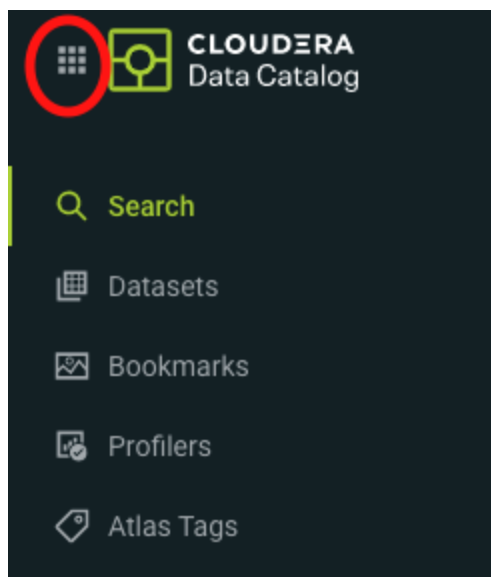*xx represents the trial user #
*X represents the password

2) Click the "Data Warehouse" within the CDP Home Screen



How do you get to the CDP Home Screen?

- From any experience such as "Data Catalog", click the 9 square at the top left and then click "Home"

CLOUDERA
Data Platform          ✕

Home

# *WE HAVE DONE THIS FOR YOU – DO NOT CREATE A NEW VIRUTAL WAREHOUSE – READ THROUGH THIS FOR BACKGROUND INFO ONLY…*

3) DO NOT Click the "+" at the top right next to "Virtual Warehouses"

4) DO NOT  Enter a name for your New Virtual Warehouse



5)   ) DO NOT  Select the Size of "xsmall - 2 Executor Nodes"
*How do I choose a size? Initial concurrent users

6) To save money you stop the instances you aren't using.  Cloudera lets you define if you spin down to zero, if you have some Kubernetes pods running all the time, and how long these live when there is no workload.   DO NOT Set the AutoSuspend Timeout (in seconds) between 4500 and 5500:

*What is AutoSuspend Timeout? Automatically spin-down unused resources after timeout occurs.

Virtual Warehouses | 1                         ⇕  Q  +

New Virtual Warehouse                                    X

Name *

testvirtualwarehouse1

Type *

HIVE      IMPALA

Database Catalog *

cdptrialuser24-dl-default

Size *

xsmall - 2 Executor Nodes

AutoSuspend Timeout (in seconds): 5000

0      1000    2000    3000    4000    5000    6000    7000

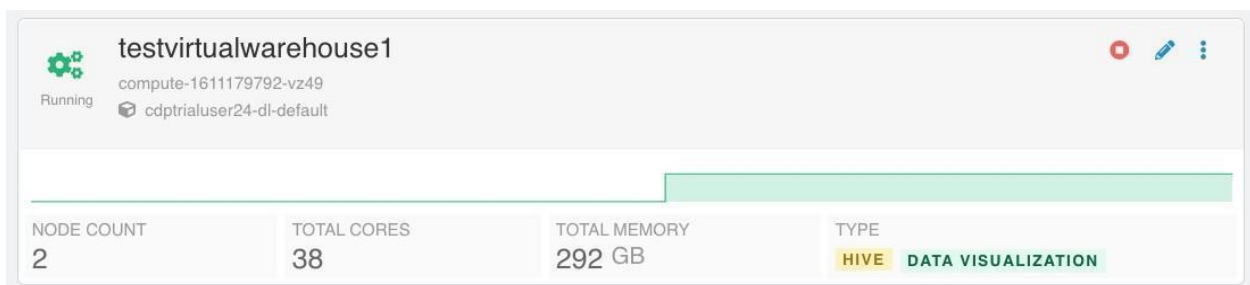7) DO NOT Choose "Install Data Visualization" to be on *Allowing for Data Visualizations in Part 3

8) DO NOT , DO NOT , DO NOT,  REALLY DO NOT: Click "Create" to create your Virtual Warehouse

*Allow for approximately 5 minutes for your Virtual Warehouse to become available for use

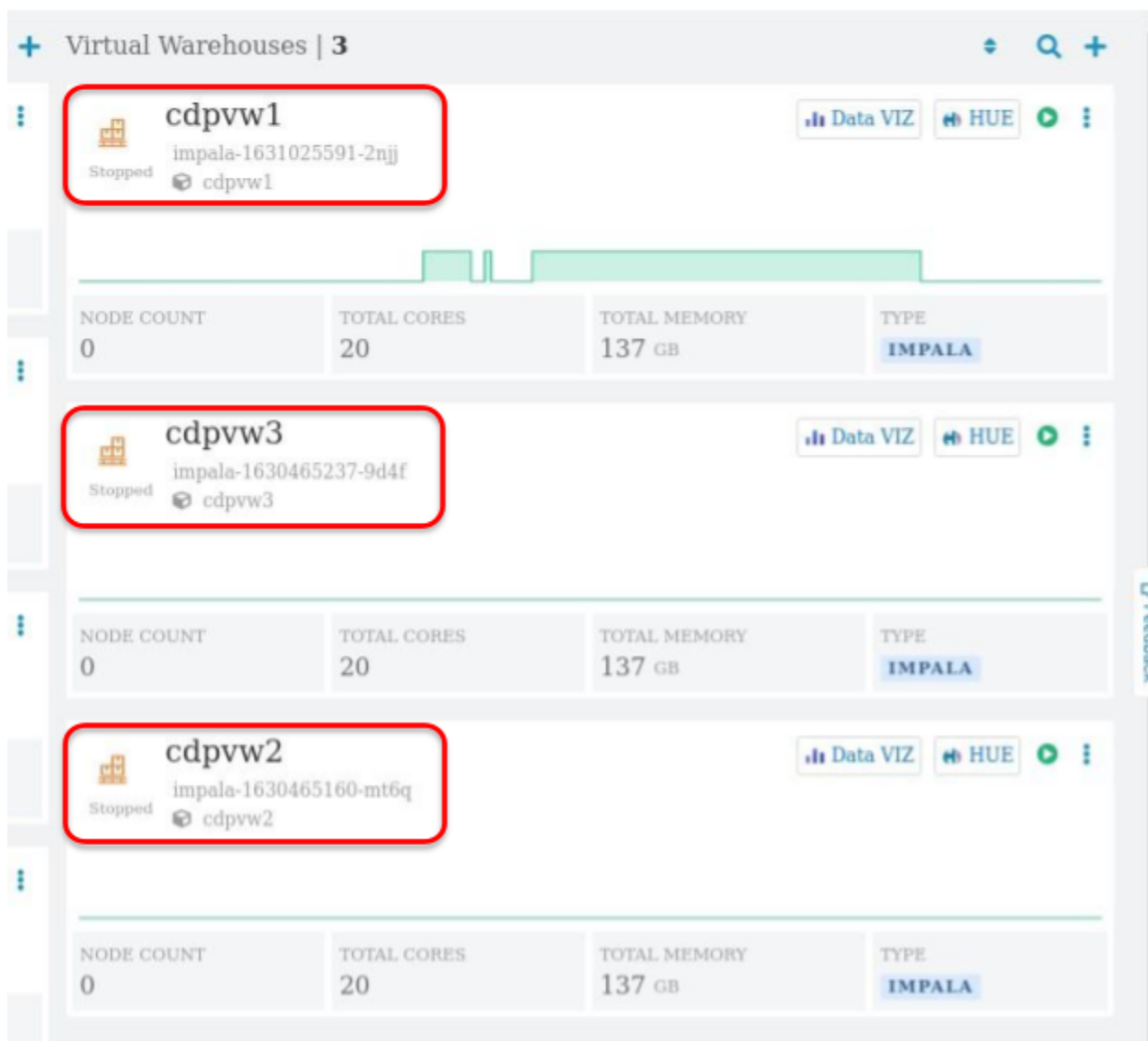When available for use, "Starting" will change to "Running" as shown below

If you can read this it is the end of the background information and it is

# TIME FOR YOU TO JUMP BACK IN AND DO THE LAB. *Please, Please, not drop any tables. Do not alter any tables. This is a shared environment. You all have*

## *the same userid and you have admin powers.*

9) Notice there multiple Virtual Warehouses (VWs).  You will be working on one of the VWs. If you are in Zoom Breakout Room 1 use cdpvw1, room2 uses cdpvw2 etc
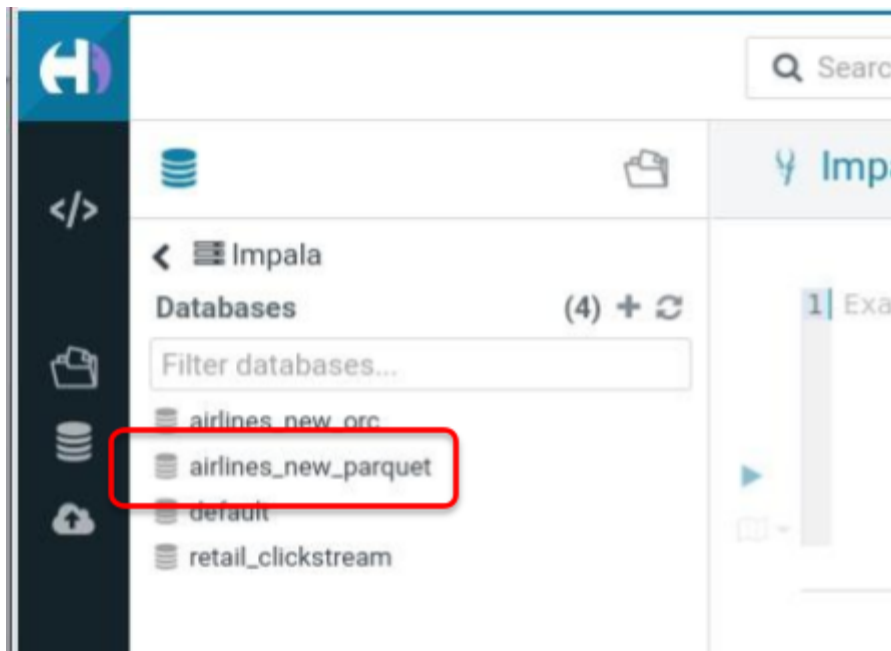
10) Click on HUE to enter the "Hadoop User Experience" in your designated room.



11) The landing page takes you to the "default" database. Click on the ◀ to the left of the default database to select a different database

12) Click on the database "airlines_new_parquet" that we saw in Part 1 "Data Catalog."  Both Impala and Hive work with both Parquet and ORC files.  As a rule of thumb if you're mostly using Impala use Parquet format or Kudu.  When working with Hive ORC is the preferred format.



If you ever need to get back to this screen layout click on the "editor"  shown here:

13) Enter the following query, answering the question "show me the top 5 visited destination by year from (1995-2008)"   Click on the blue triangle to run the query.

```
SELECT dest,year,COUNT(dest) as Times_Visited FROM flights
GROUP BY dest,year
ORDER BY Times_Visited DESC
LIMIT 5;
```

Why doesn't it run right away?  You promised it was fast ☺

**Click on blue triangle to run**

**What is "waiting for query executors to start?"**

**Impala went to sleep to save money on cloud costs. You can see it wake back up to answer your query. You learned about this in the configuration section where you did not provision a virtual warehouse.**

You may get the query result without having to wait for Impala to "wake up."  You get to decide how long an idle time you want to wait before you scale down to zero, or if you have the budget you can have Impala always available.

14) Click "EXPLAIN" to see the explain plan prior to running the query
*Not required to execute the query - this gives us a plan on exactly what the query is doing

Click on this down-arrow to expose the menu

The "Explain Plan" shows you how the query will execute.   It is asking you to update statistics for the tables in the query.  This is a good idea for performance.  After everyone reaches this point in the lab designate one person to run: **compute stats airlines_new_parquet.flights;**

We are in a shared environment.  If someone else has done "compute stats" you won't see the message.  In STEP 4, after the DataViz lab,  you will have a chance to create your own tables and work on compute stats with your unique tables.

We can discuss optimizer paths and table statistic  in more detail as part of the breakout room.

15) After you run the query you will have a link to lots of information about the query. Click on the link shown below



16) Explore the different tabs in the query pop-up.  The Visual Plan shows how the query was executed.  These get more interesting with multi table joins. The Summary show how much time was spent in each stage of the query, the memory used, and the rows produced.  The Profile includes the summary and great detail of everything that happened on each node running the query.
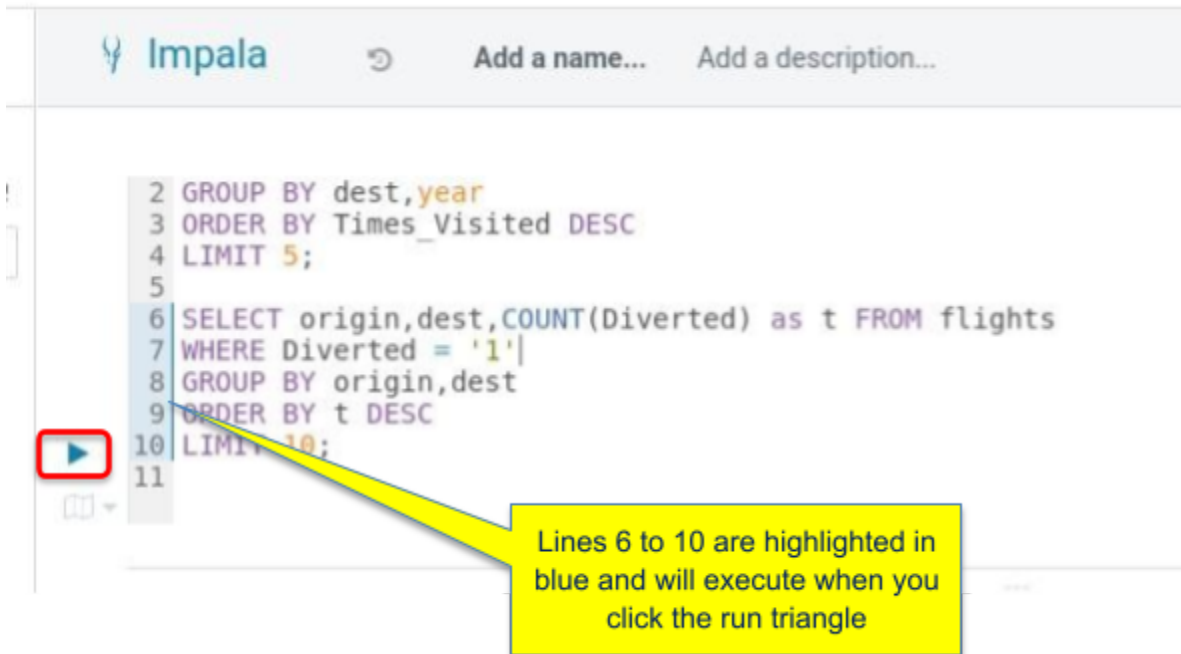
## 17) Add another query to the editor

```
SELECT origin,dest,COUNT(Diverted) as t FROM flights
WHERE Diverted = "1"
GROUP BY origin,dest
ORDER BY t DESC
LIMIT 10;
```
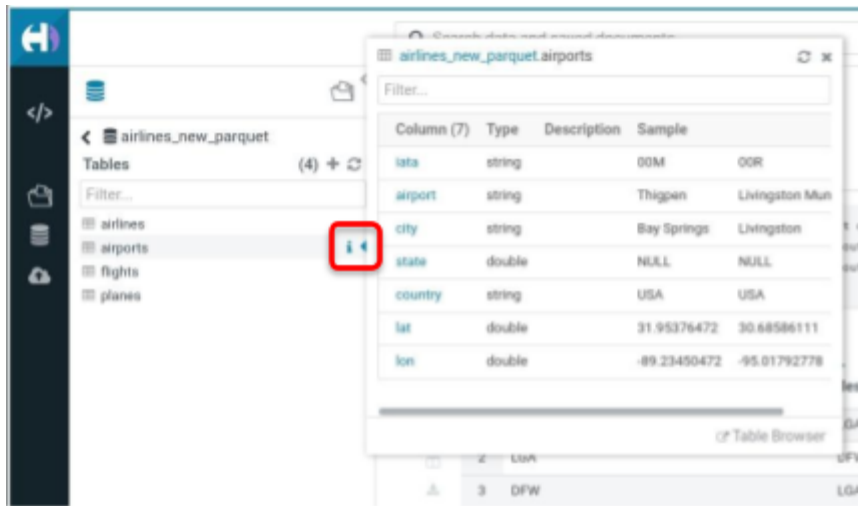
18) Notice the highlighting on the left edge. HUE is parsing based on the semi-colon and the execution arrow will run whatever is highlighted in blue, or whatever has been highlighted by the cursor.  This way you can have multiple queries in the same canvas.



19) Click the blue triangle to execute the query, answering the question "What are the top 10 routes (origin and dest) that have seen maximum diversions?"
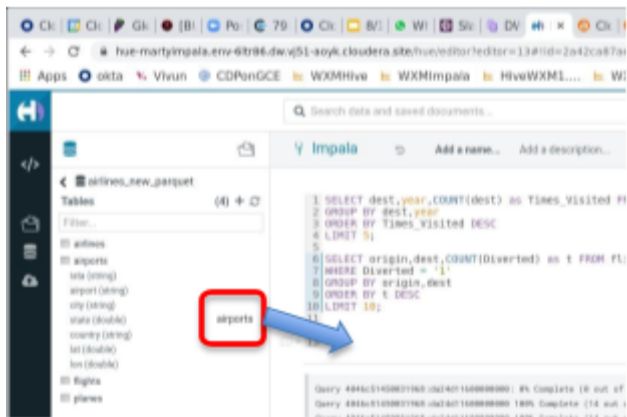
| | origin | dest | t |
|---|---|---|---|
| 1 | ORD | LGA | 845 |
| 2 | LGA | DFW | 749 |
| 3 | DFW | LGA | 653 |
| 4 | DAL | HOU | 615 |
| 5 | ATL | LGA | 567 |
| 6 | MDW | STL | 512 |
| 7 | ATL | DFW | 482 |
| 8 | ORD | DFW | 450 |

20) Hover over the disappearing  `i`  next to the airports table to see more information about the table.  We're going to use the geo location of the airports to do a marker map in HUE.
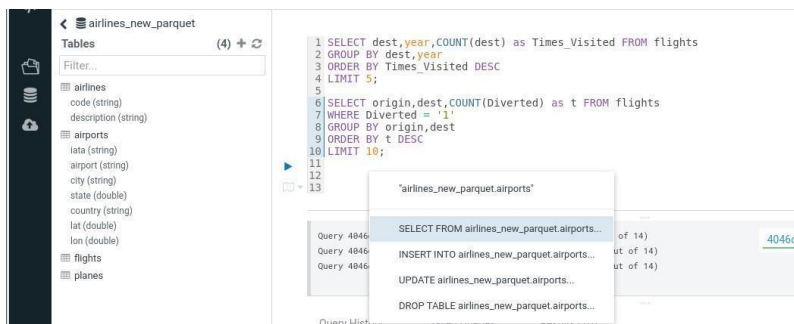
You can also click on the table name to expand all the columns.  On the far right of the browser is help tooling.

21) Try out the "drag and drop" option, drag and drop the table name "airports over to line 12



When you drop the table name you'll get to choose what SQL you want auto generated. Take the "select" option and run the query with the blue-triangle

22) Let's now build a marker map of the airports with the most cancellations.  This will correlate
with the airports that have the most flights.  Run the SQL  shown below

```
SELECT origin, lat, lon, COUNT(Cancelled) as num_of_cancellations ,
concat(origin, " ", cast(count(cancelled) as string)) as airport_label
FROM flights, airports
WHERE origin = iata and cancelled = 1 AND cancellationcode = 'B'
GROUP BY origin, lat, lon order by num_of_cancellations desc limit 50;
```
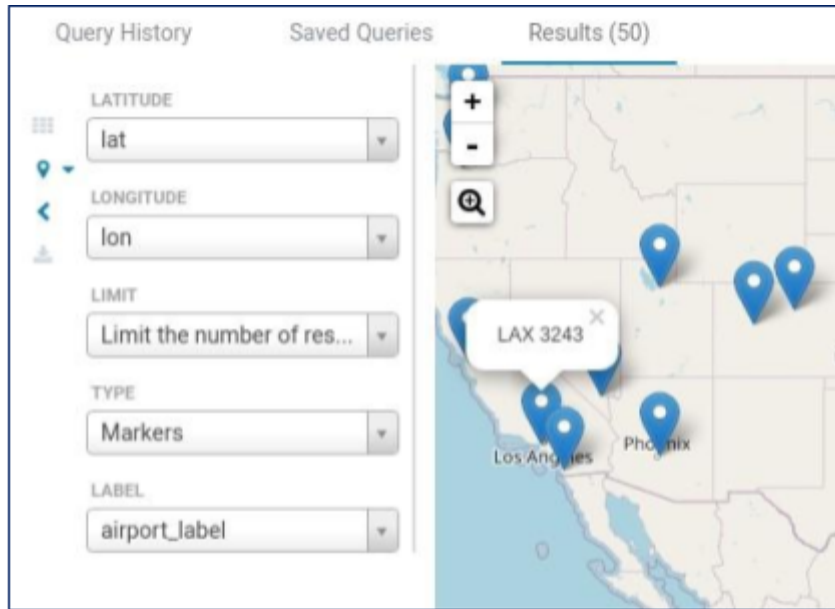
23) After you run the SQL use the down arrow to choose the type of output formatting.  You've
been using the data grid, we're now going to choose the marker-map

24) Configure the marker map per shown below. Clicking on one of the markers will pop up the value of the "ariport_label" column



25) Look at the query plan – notice we now have a join in the tree

26) The summary shows details of all the stages in the join and their metrics



27) Time to save our work.  Give your Impala SQL a name.  For the lab use yourNameXXX where XXX is where you get to be creative.  Use something unique, this is a multi user environment.

```
5
6 SELECT origin,dest,COUNT(Diverted) as t FROM flights
7 WHERE Diverted = '1'
```

Then click "Save" and then "Save" in the popup





Your saved queries will show up under the "Saved Queries"  heading.

```
20 WHERE origin = iata and cancelled = 1 AND cancellationcode = 'B'
21 GROUP BY origin, lat, lon order by num_of_cancellations desc limit 50;
22 |;
```

Query 6c419d4a75c5fabe:772554f000000000 100% Complete (15 out of 15)
Query 6c419d4a75c5fabe:772554f000000000 100% Complete (15 out of 15)

6c419d4a75c5fabe:772554f000000000

Query History        Saved Queries        Results (50)

| Name | Description | Owner | Last Modified |
|------|-------------|-------|---------------|
| martyImpalaQueries123 | | trial10_admin_0 | 08/17/2021 1:55 PM -04:00 |