



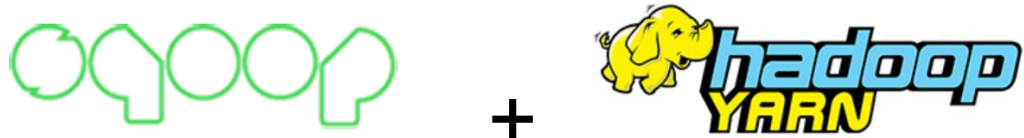
## Tutorial 3: Apache Sqoop with Cloudera

### CN7022 - Big Data Analytics

Dr Amin Karami ([a.karami@uel.ac.uk](mailto:a.karami@uel.ac.uk))

**LEARNING OUTCOMES:** After completing this tutorial, you should:

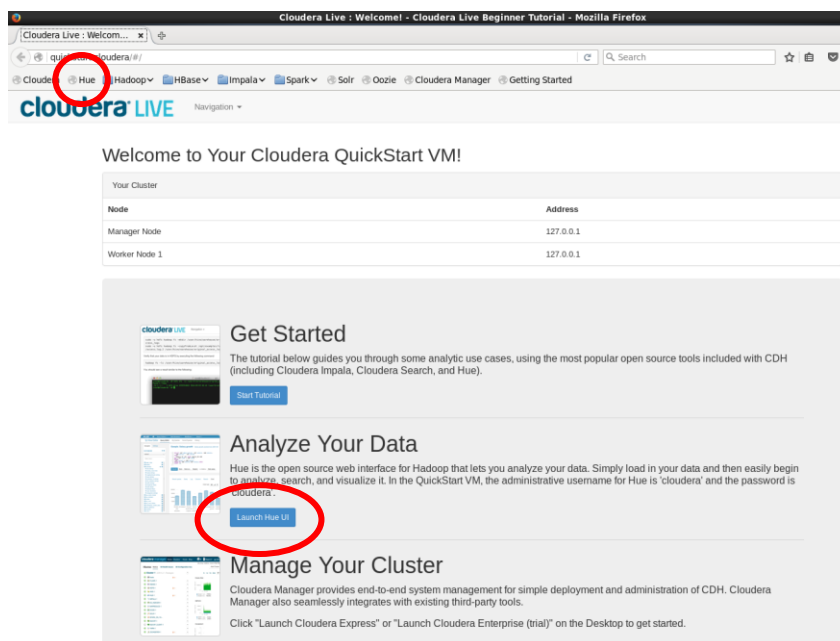
- Have gotten a hands-on experience in deploying Apache Sqoop with Cloudera
- Learn to transfer data from a relational database into Hadoop using Sqoop
- Learn hands-on practices on Sqoop commands



### Phase 1: Getting Ready for Hue Cloudera

Hue is an open-source SQL Cloud Editor, licensed under the Apache v2 license, to dynamically interact and visualize big sized data.

Launch the “cloudera-quickstart-vm” from VMWare workstation. It takes a couple of minutes to load up. Then, find and click on **Hue** icon in the browser from one of the two ways as shown in the following screenshot.



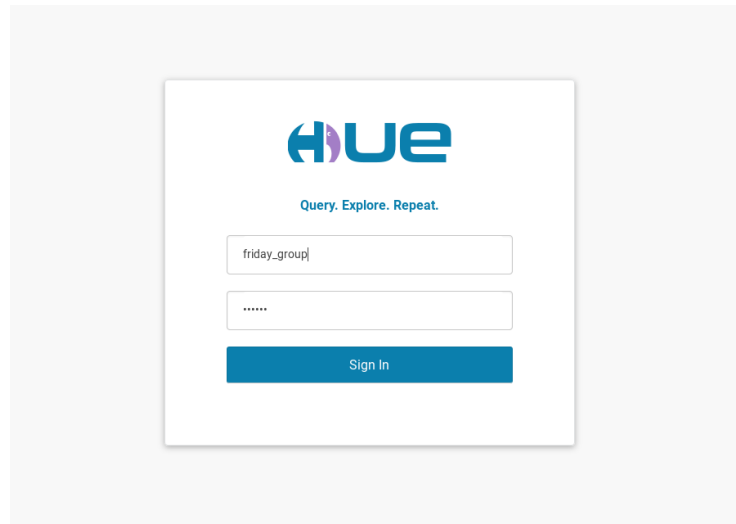


Then, enter your username/password:

Username: **friday\_group** Password: **friday** [default]

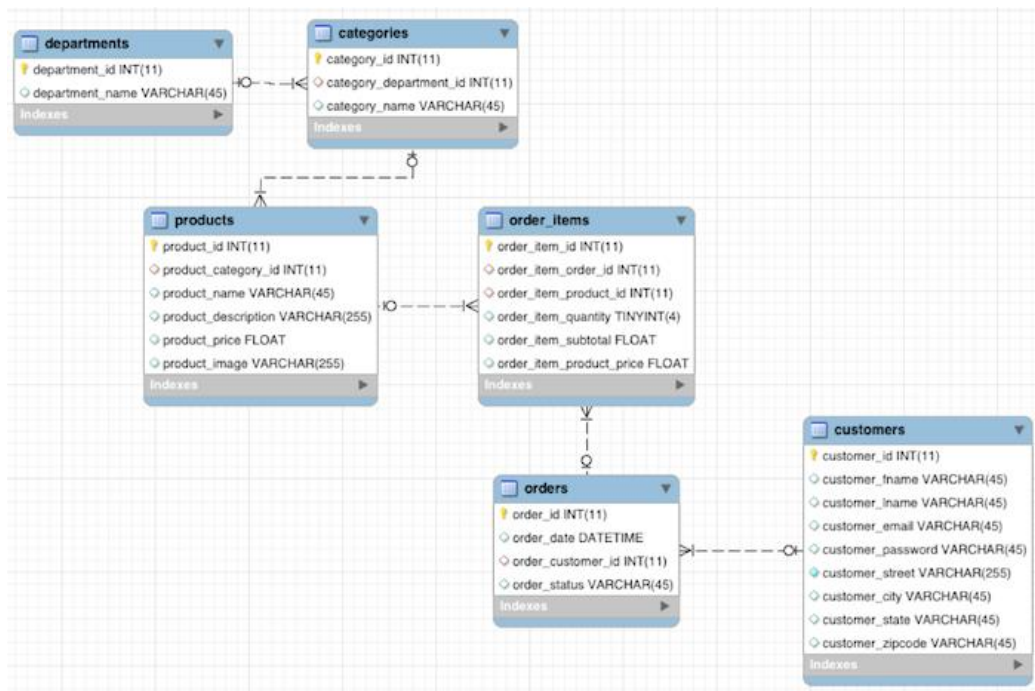
Username: **tuesday\_group** Password: **tuesday**

Username: **wednesday\_group** Password: **wednesday**



## Phase 2: Working with Sqoop

The flowchart of `retail_db` is in the following picture. We are going to import it into the Hadoop with Sqoop. This database already exists in the Cludera. Let us check the database. Source: <https://www.cludera.com/developers/get-started-with-hadoop-tutorial/exercise-1.html>





a) Open cmd. Login to MySQL Database Server (password: cloudera).

```
mysql -u root -p
```

b) Show all the databases and the tables.

```
show databases;  
use retail_db;  
show tables;  
quit
```

```
+-----+  
| Tables_in_retail_db |  
+-----+  
| categories  
| customers  
| departments  
| order_items  
| orders  
| products  
+-----+  
6 rows in set (0.00 sec)
```

c) Open a new terminal and run the following Sqoop command line by line. **It would take a while to be completed, because Sqoop launches MapReduce job.**

```
[cloudera@quickstart ~]$ sqoop import-all-tables \  
-m 2 \  
--connect jdbc:mysql://localhost:3306/retail_db \  
--username=root \  
--password=cloudera \  
--compression-codec=snappy \  
--as-parquetfile \  
--warehouse-dir=/user/hive/warehouse \  
--hive-overwrite \  
--hive-import
```

After about 10 minutes, you have to see the following screenshot as the confirmation:



```

cloudera@quickstart:~
File Edit View Search Terminal Help
HDFS: Number of write operations=20
Job Counters
  Launched map tasks=2
  Other local map tasks=2
  Total time spent by all maps in occupied slots (ms)=27937
  Total time spent by all reduces in occupied slots (ms)=0
  Total time spent by all map tasks (ms)=27937
  Total vcore-milliseconds taken by all map tasks=27937
  Total megabyte-milliseconds taken by all map tasks=28607488
Map-Reduce Framework
  Map input records=1345
  Map output records=1345
  Input split bytes=236
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=532
  CPU time spent (ms)=9040
  Physical memory (bytes) snapshot=857354240
  Virtual memory (bytes) snapshot=3157467136
  Total committed heap usage (bytes)=711983104
File Input Format Counters
  Bytes Read=0
File Output Format Counters
  Bytes Written=0
19/10/20 07:34:12 INFO mapreduce.ImportJobBase: Transferred 55.624 KB in 112.2176 seconds (507.5763 bytes/sec)
19/10/20 07:34:12 INFO mapreduce.ImportJobBase: Retrieved 1345 records.
[cloudera@quickstart ~]$

```

d) **Verification:** When this command is complete, confirm that your data files exist in HDFS. These commands will show the directories and the files inside them that make up your tables:

- `hdfs dfs -ls /user/hive/warehouse/`
- `hdfs dfs -ls /user/hive/warehouse/categories/`

```

[cloudera@quickstart ~]$ hdfs dfs -ls /user/hive/warehouse/
Found 6 items
drwxrwxrwx - cloudera supergroup      0 2019-10-20 14:57 /user/hive/warehouse/categories
drwxrwxrwx - cloudera supergroup      0 2019-10-20 14:59 /user/hive/warehouse/customers
drwxrwxrwx - cloudera supergroup      0 2019-10-20 15:02 /user/hive/warehouse/departments
drwxrwxrwx - cloudera supergroup      0 2019-10-20 15:04 /user/hive/warehouse/order_items
drwxrwxrwx - cloudera supergroup      0 2019-10-20 15:06 /user/hive/warehouse/orders
drwxrwxrwx - cloudera supergroup      0 2019-10-20 15:08 /user/hive/warehouse/products
[cloudera@quickstart ~]$ hdfs dfs -ls /user/hive/warehouse/categories
Found 4 items
drwxr-xr-x - cloudera supergroup      0 2019-10-20 14:56 /user/hive/warehouse/categories/.metadata
drwxr-xr-x - cloudera supergroup      0 2019-10-20 14:57 /user/hive/warehouse/categories/.signals
-rw-r--r-- 1 cloudera supergroup    1491 2019-10-20 14:57 /user/hive/warehouse/categories/0e01f6e7-4f7a-4b66-9225-34c2bbe9bfd6.parquet
-rw-r--r-- 1 cloudera supergroup    1520 2019-10-20 14:57 /user/hive/warehouse/categories/c141345a-2a78-48cc-8801-058e7e53d102.parquet

```

**Note:** The number of `.parquet` files shown will be equal to what was passed to Sqoop with the `-m` parameter. This is the number of 'mappers' that Sqoop will use in its MapReduce jobs. It could also be thought of as the number of simultaneous connections to your database, or the number of disks/Data Nodes you want to spread the data across. So on a single-node you will just see one, but larger clusters will have a greater number of files.

OR, complete the verification step visually (localhost:50070):



/user/hive/warehouse								Go!
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
drwxrwxrwx	cloudera	supergroup	0 B	Sun Oct 20 14:57:39 -0700 2019	0	0 B	<a href="#">categories</a>	
drwxrwxrwx	cloudera	supergroup	0 B	Sun Oct 20 14:59:25 -0700 2019	0	0 B	<a href="#">customers</a>	
drwxrwxrwx	cloudera	supergroup	0 B	Sun Oct 20 15:02:03 -0700 2019	0	0 B	<a href="#">departments</a>	
drwxrwxrwx	cloudera	supergroup	0 B	Sun Oct 20 15:04:20 -0700 2019	0	0 B	<a href="#">order_items</a>	
drwxrwxrwx	cloudera	supergroup	0 B	Sun Oct 20 15:06:39 -0700 2019	0	0 B	<a href="#">orders</a>	
drwxrwxrwx	cloudera	supergroup	0 B	Sun Oct 20 15:08:46 -0700 2019	0	0 B	<a href="#">products</a>	

## Phase 3: Working with Impala

We are going to use Hue's Impala to query our tables. *Impala is a query engine that runs on Hadoop.* Once you are inside the Hue editor, click on **Query** → **Editor** → **Impala** to launch the Impala Editor.

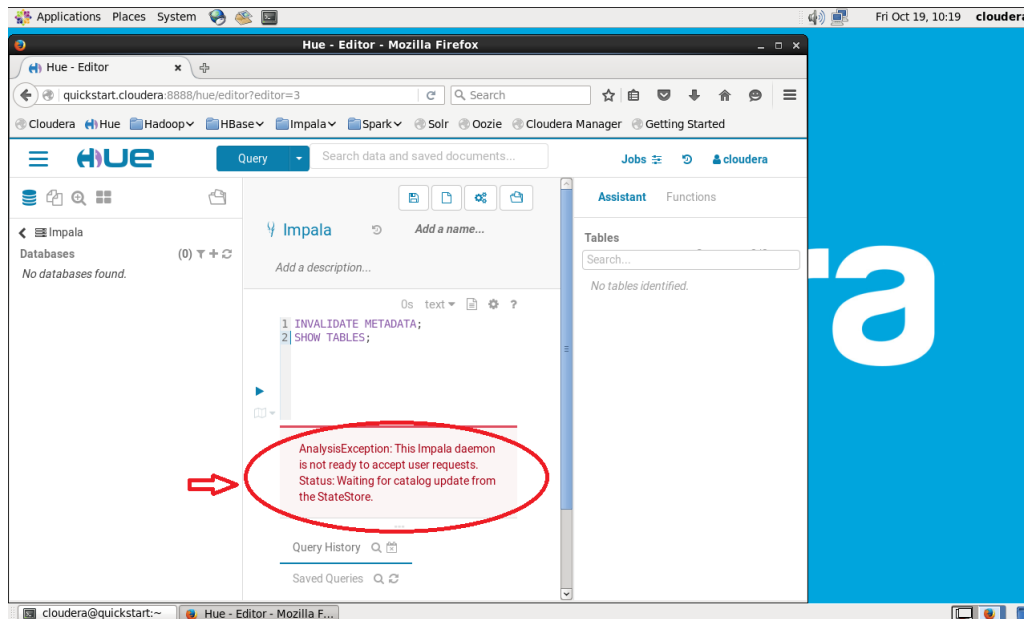
To save time during queries, Impala does not poll constantly for metadata changes. So the first thing we must do is tell Impala that its metadata is out of date. Then we should see our tables show up, ready to be queried:

```
invalidate metadata;  
show tables;
```

The screenshot shows the Hue Impala interface. At the top, there's a header with the Impala logo and options to 'Add a name...' and 'Add a description...'. Below this is a query editor with a single line of code: `show tables`. To the left of the editor is a sidebar with icons for query history, saved queries, and results. Below the editor, there's a section for 'Results (6)' which displays a table with the following data:

	name
1	categories
2	customers
3	departments
4	order_items
5	orders
6	products

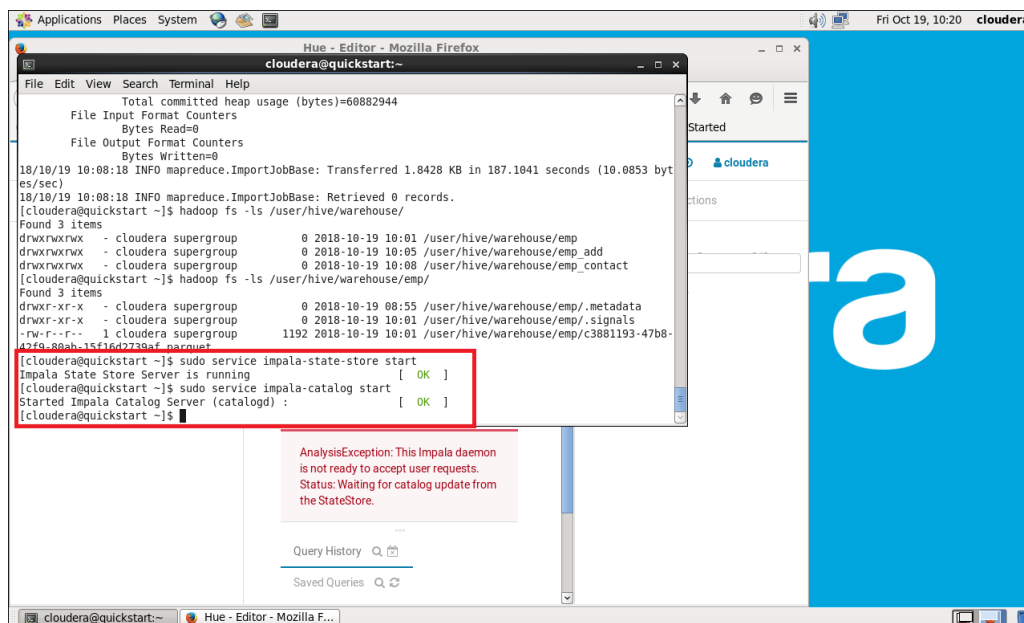
**[Exceptional Error:]** You may get the following error when you run the commands. This is because both services `impala-state-store` and `impala-catalog` are down and need to be restarted.



To fix the error, **back to the terminal and type the following commands:**

```
sudo service impala-state-store start
```

```
sudo service impala-catalog start
```



Now click on Hue icon again to refresh the page and then run the query. You can now see the new tables in the result menu.

- e) Now, your transaction data is readily available for structured queries in CDH, *it's time to address a few questions. You can make a few basic/advanced SQL queries to get familiar with the retail\_db database. Also, you can see the tables' properties visually in the Hue platform, such as:*



The screenshot shows a database management interface. On the left, a sidebar lists tables: categories, customers, departments, order\_items, orders, and products. The main area displays a query 'show tables' with results showing a list of table names: categories, customers, departments, order\_items, orders, and products.

name
1 categories
2 customers
3 departments
4 order_items
5 orders
6 products

f) SQL query 1: calculate the total revenue per product and showing the top 10 revenue generating products:

```
select c.category_name, count(order_item_quantity) as count
from order_items oi
inner join products p on oi.order_item_product_id = p.product_id
inner join categories c on c.category_id = p.product_category_id
group by c.category_name
order by count desc
limit 10;
```

The screenshot shows the Impala query results page. The query is the same as the one in the previous block. The results table shows the top 10 revenue-generating products, ordered by count in descending order.

category_name	count
1 Cleats	24551
2 Meris Footwear	22246
3 Women's Apparel	21035
4 Indoor/Outdoor Games	19298
5 Fishing	17325
6 Water Sports	15540
7 Camping & Hiking	13729
8 Cardio Equipment	12487
9 Shop By Sport	10984
10 Electronics	3156



g) SQL query 2: get the top 10 revenue generating products:

```
select p.product_id, p.product_name, r.revenue
from products p inner join
(select oi.order_item_product_id,
sum(cast(oi.order_item_subtotal as float)) as revenue
from order_items oi inner join orders o
on oi.order_item_order_id = o.order_id
where o.order_status <> 'CANCELED'
and o.order_status <> 'SUSPECTED_FRAUD'
group by order_item_product_id) r
on p.product_id = r.order_item_product_id
order by r.revenue desc
limit 10;
```

The screenshot shows the Impala SQL interface. The query is executed, and the results are displayed in a table with 10 rows. The table has three columns: product\_id, product\_name, and revenue. The results are sorted by revenue in descending order.

	product_id	product_name	revenue
1	1004	Field & Stream Sportsman 16 Gun Fire Safe	6637668.2823181152
2	365	Perfect Fitness Perfect Rip Deck	4233794.3682899475
3	957	Diamondback Women's Serene Classic Comfort BI	3946837.0045471191
4	191	Nike Men's Free 5.0+ Running Shoe	3507549.2067337036
5	502	Nike Men's Dri-FIT Victory Golf Polo	3011600
6	1073	Pelican Sunstream 100 Kayak	2967851.6815185547
7	1014	O'Brien Men's Neoprene Life Vest	2765543.314743042
8	403	Nike Men's CJ Elite 2 TD Football Cleat	2763977.4868011475
9	627	Under Armour Girls' Toddler Spine Surge Runni	1214896.220287323
10	565	adidas Youth Germany Black/Red Away Match Soc	63490