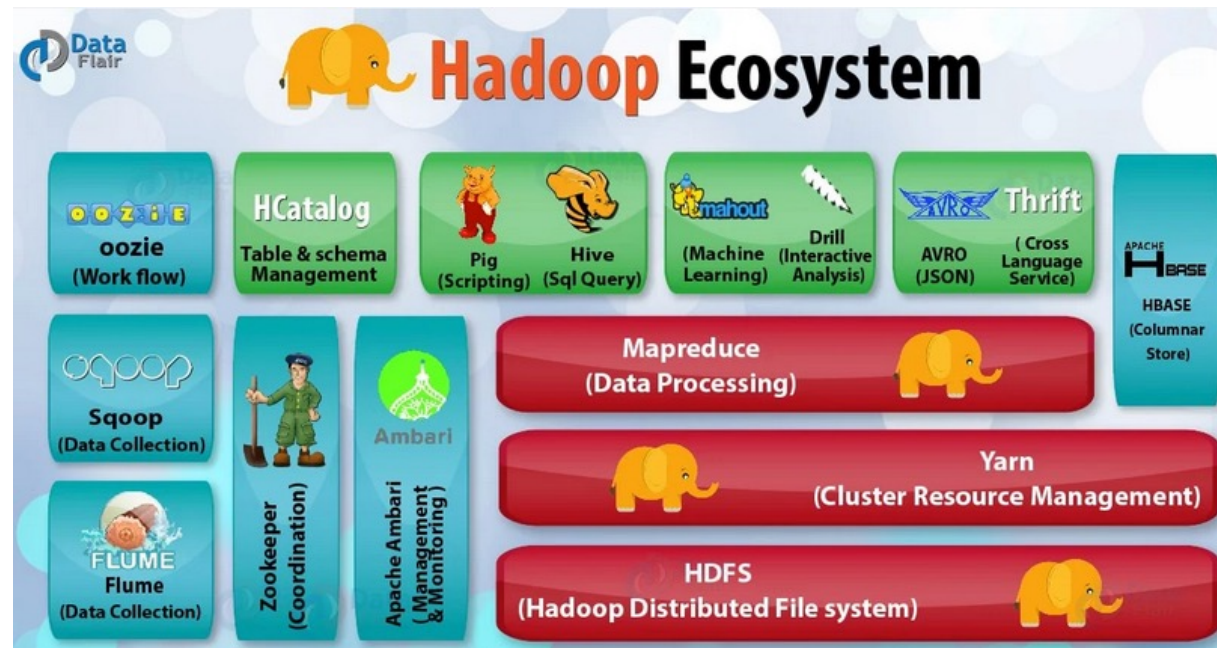# Querying the Data interactively

Bernard Lee Kok Bang

**allow us to execute SQL quesries :)**

# Core Hadoop Ecosystem

# External Data Storage

AC

cassandra

MySQL

CP

mongoDB

# Query Engines

analyzing the data

# Apache Drill – SQL for NoSQL

- A SQL query engine for a variety of non-relational databases and data files

  *- Hive, MongoDB, HBase*

  <span style="color:red">wh</span>

  *- Flat JSON or Parquet files on HDFS, S3, Azure, Google Cloud, local file system*

- Based on Google's Dremel

- **<span style="color:red">let us execute actual SQL</span>**

# Parquet – columnar format

- stores the data in row group



*https://shorturl.at/efno0*

# It's real SQL

- Not SQL-like
- And it has ODBC (Microsoft) / JDBC (Sun Microsystems) driver so other tools can connect to it just like any relational database
    - *allow connection to external tools such as Tableau*

# Drill - Fast and easy to set up

- Still non-relational databases under the hood

- Allow SQL analysis of disparate data sources without having to transform and load

  - *Internally data is represented as JSON and so has no fixed schema*

- Doing JOINS across different database technologies or with JSON files that are sitting around in HDFS

- **having a SQL without a schema**

# Let's Drill

- Import data into Hive (movie ratings) and MongoDB (movie users) – MovieLens datasets

- Set up Drill on top of both

- And do some queries

# Setting up Drill

1. Import *u.user* into MongoDB

2. Upload *u.data* into Hive
   - *Login to Ambari as admin*
   - *Start MongoDB*

- Navigate to Hive View to create a new database
  - *CREATE DATABASE movielens;*

- upload u.data into movielens database

# Upload data into Hive

# Import data into MongoDB

- Log into puTTY as root:

  *su root*

- run the previous *MongoSpark.py* script to import *u.user* data into mongoDB database [make sure the ~~u.ser~~ data is still in the ml-100k folder in HDFS]

  ***u.user***

- submit the script using the following command

  *spark-submit --packages org.mongodb.spark:mongo-spark-connector_2.11:2.3.2 MongoSpark.py*

# Installing Drill

- Not a built-in part of Ambari on Hortonworks platform
- Need to get specific version of *Drill: version 1.12*

```
[root@sandbox-hdp maria_dev]# wget http://archive.apache.org/dist/drill/drill-1.
12.0/apache-drill-1.12.0.tar.gz
```

- Decompress the downloaded file
  - *tar –xvf apache-drill-1.12.0.tar.gz*

- To start running Drill, run the following script [open the port 8765]
  - *cd apache-drill-1.12.0*
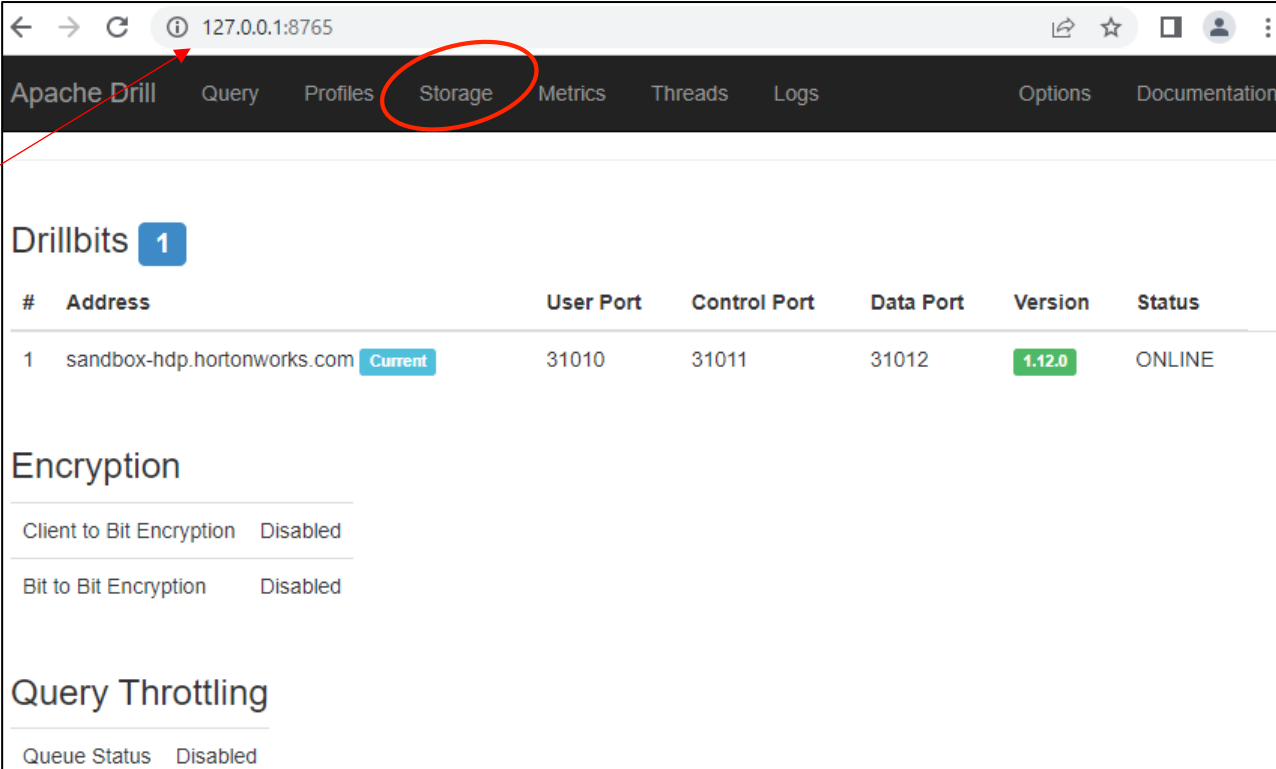  - *bin/drillbit.sh start –Ddrill.exec.http.port=8765*  **this is to start Drill**

# Connect Drill through browser

- Log into browser address specific for Drill
  - *127.0.0.1:8765*

**Key in the Drill address**

# Connect Drill to Hive and MongoDB

- Click *Storage* tab
- *Enable* both *mongo and hive*
- Update both hive and mongo

# Configure Hive

- Update *hive.metastore.uris*
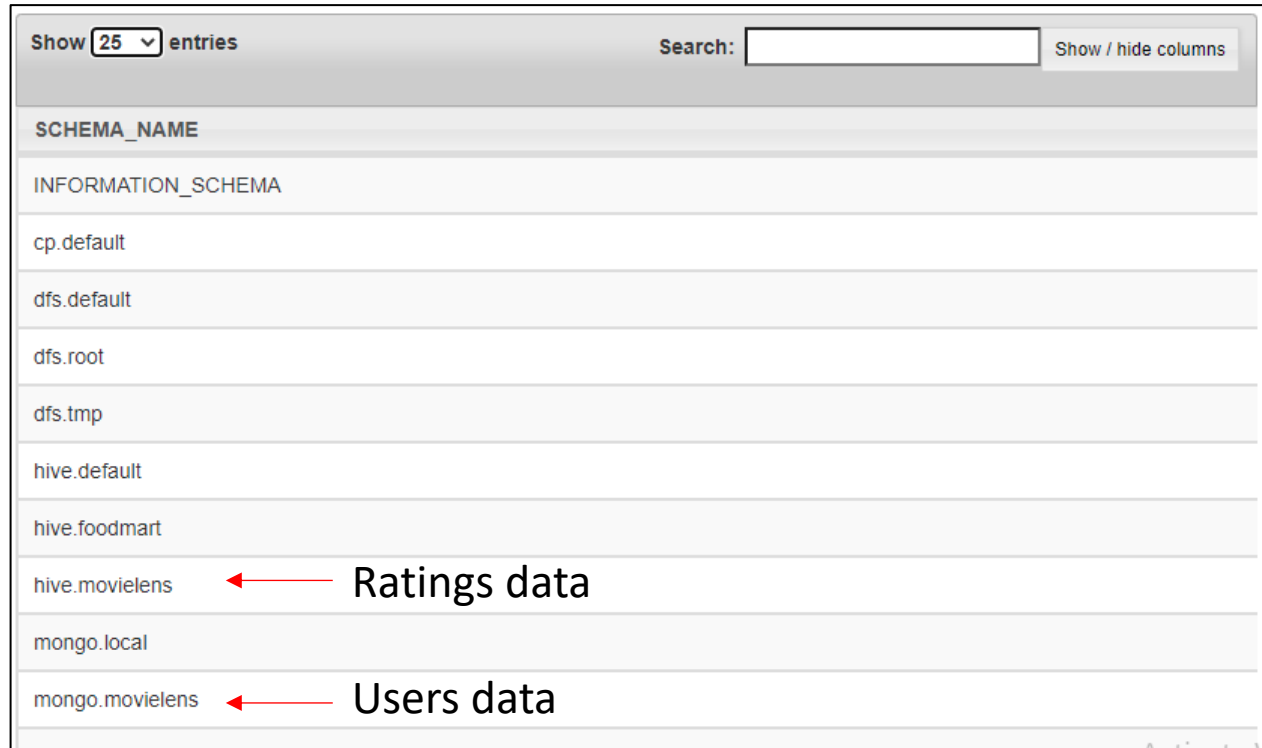


*Update hive metastore accordingly: thrift://localhost:9083*

# Start querying across multiple databases using Drill

- *SHOW DATABASES;*

| SCHEMA_NAME | |
|---|---|
| INFORMATION_SCHEMA | |
| cp.default | |
| dfs.default | |
| dfs.root | |
| dfs.tmp | |
| hive.default | |
| hive.foodmart | |
| hive.movielens | ← Ratings data |
| mongo.local | |
| mongo.movielens | ← Users data |

# Start querying (cont...)

- *SELECT * FROM hive.movielens.ratings LIMIT 10;*

| user_id | movie_id | rating | epoch_seconds |
|---------|----------|--------|---------------|
| 196 | 242 | 3 | 881250949 |
| 186 | 302 | 3 | 891717742 |
| 22 | 377 | 1 | 878887116 |
| 244 | 51 | 2 | 880606923 |
| 166 | 346 | 1 | 886397596 |
| 298 | 474 | 4 | 884182806 |
| 115 | 265 | 2 | 881171488 |
| 253 | 465 | 5 | 891628467 |
| 305 | 451 | 3 | 886324817 |
| 6 | 86 | 3 | 883603013 |

Show 10 entries    Search: ____    Show / hide columns

# Start querying (cont…)

- *SELECT * FROM mongo.movielens.users LIMIT 10;*

| _id | age | gender | occupation | user_id | zip |
|---|---|---|---|---|---|
| [B@3157f22e | 24 | M | technician | 1 | 85711 |
| [B@1fe2c00 | 53 | F | other | 2 | 94043 |
| [B@1d439b97 | 23 | M | writer | 3 | 32067 |
| [B@5494494c | 24 | M | technician | 4 | 43537 |
| [B@5a5c1772 | 33 | F | other | 5 | 15213 |
| [B@2a32be15 | 42 | M | executive | 6 | 98101 |
| [B@2b2c07e7 | 57 | M | administrator | 7 | 91344 |
| [B@5d84e5e7 | 36 | M | administrator | 8 | 05201 |
| [B@67599ea2 | 29 | M | student | 9 | 01002 |
| [B@5824564a | 53 | M | lawyer | 10 | 90703 |

Show 10 entries — Search: — Show / hide columns

# Start querying (cont…)

- *SELECT u.occupation, COUNT(*) FROM hive.movielens.ratings r JOIN mongo.movielens.users u ON r.user_id = u.user_id GROUP BY u.occupation;*

**How many ratings were provided by each occupation?**

| occupation | EXPR$1 |
|---|---|
| administrator | 7479 |
| artist | 2308 |
| doctor | 540 |
| educator | 9442 |
| engineer | 8175 |
| entertainment | 2095 |
| executive | 3403 |
| healthcare | 2804 |
| homemaker | 299 |
| lawyer | 1345 |
| librarian | 5273 |
| marketing | 1950 |

Show 25 entries    Search: ___    Show / hide columns

# Remember to clean up the mess!!!

# Stop drillbit

- bin/drillbit.sh stop

- Stop mongoDB in Ambari

- Exit from puTTY and VirtualBox

**this is to stop Drill**