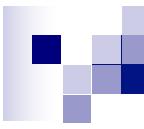


## **Module 7**

# **Understanding Sqoop**

Thanachart Numnonda, Executive Director, IMC Institute

Thanisa Numnonda, Faculty of Information Technology,  
King Mongkut's Institute of Technology Ladkrabang



# Introduction

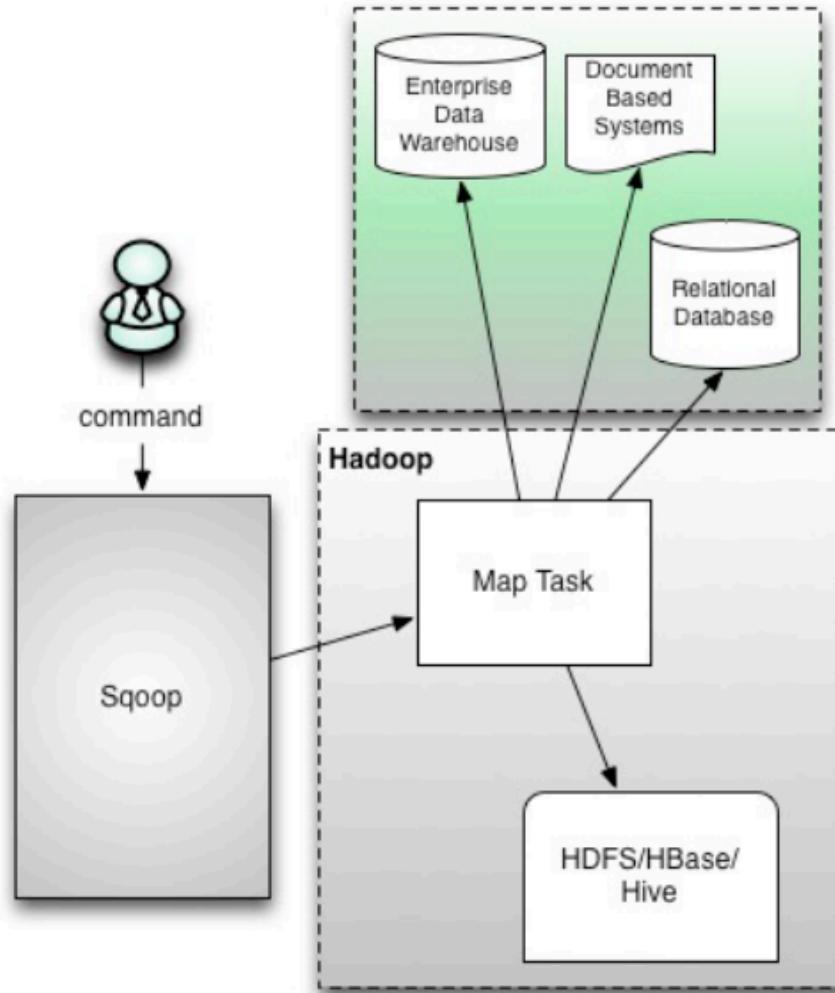


**Sqoop (“SQL-to-Hadoop”) is a straightforward command-line tool with the following capabilities:**

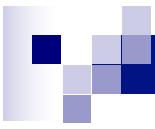
- .Imports individual tables or entire databases to files in HDFS**
- .Generates Java classes to allow you to interact with your imported data**
- .Provides the ability to import from SQL databases straight into your Hive data warehouse**

Text

# Architecture Overview



เป็น dataware  
house ได้ ไม่จำเป็น  
ต้อง relation ขอแค่  
มี .....

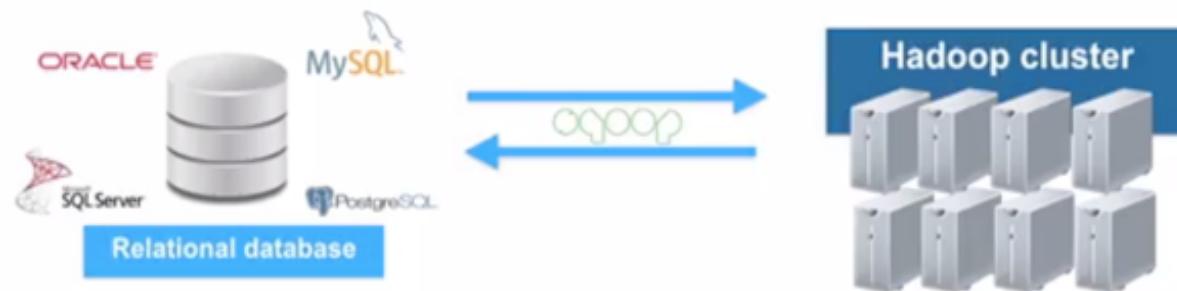


# Sqoop Benefit

- Leverages RDBMS metadata to get the column data types
- It is simple to script and uses SQL
- It can be used to handle change data capture by importing      อาจจะใช้ sqoop import ไว้เห็บใน hadoop daily transactional data to Hadoop
- It uses MapReduce for export and import that enables parallel and efficient data movement

# Sqoop Mode

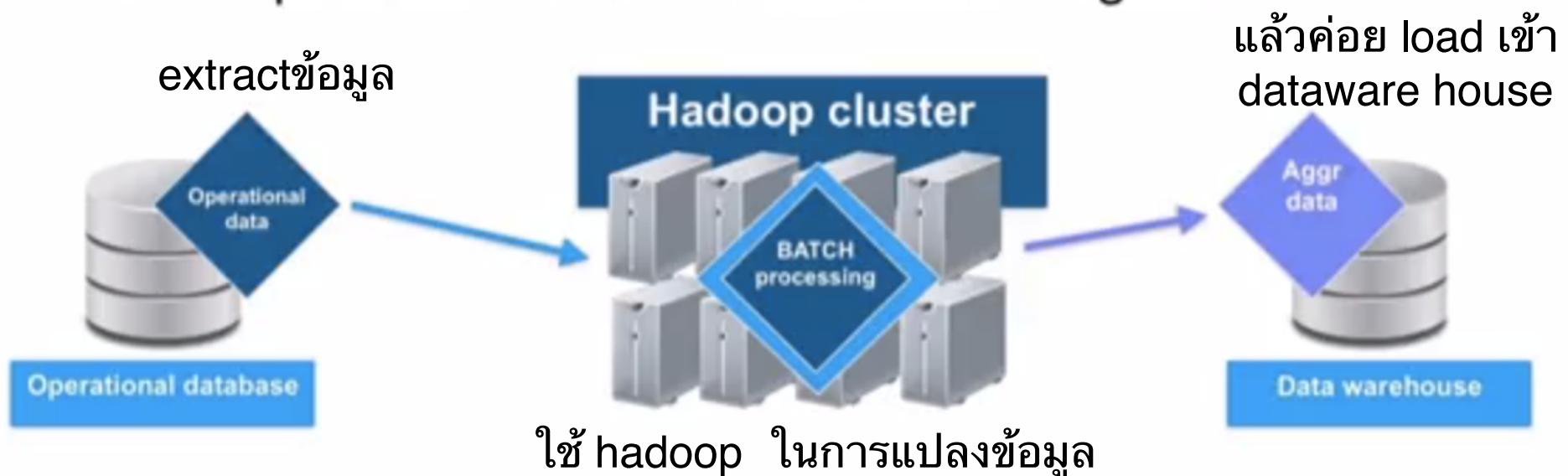
- Sqoop import: Data moves from RDBMS to Hadoop
- Sqoop export: Data moves from Hadoop to RDBMS



# Use Case #1: ETL for Data Warehouse

extract tranforms

- Transform operational data for data warehouse reports in Hadoop as the batch transformation “engine”



# Use Case #2: ELT

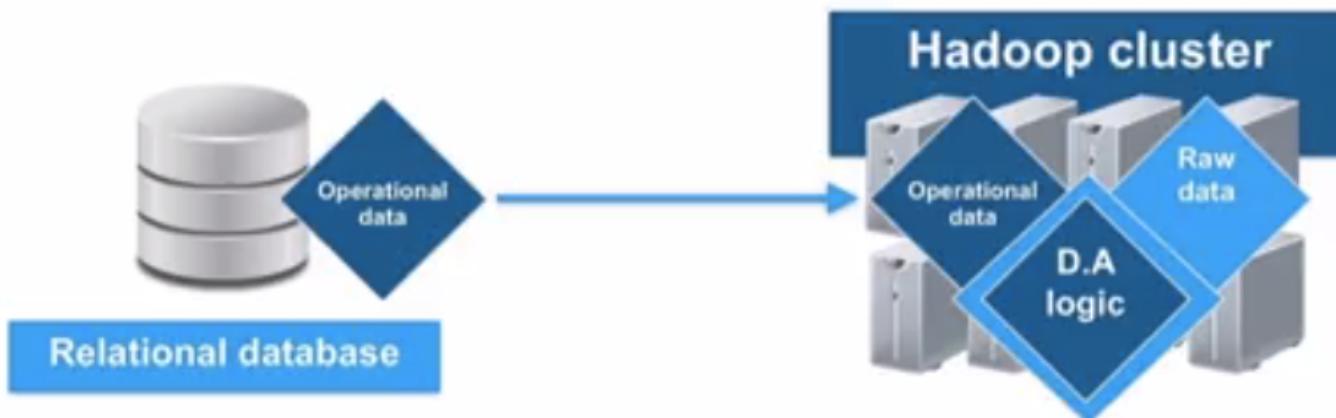
- Extract operational data from RDBMS, process in Hadoop, return **result** to RDBMS



extract เข้ามาโหลดเข้ามาเก็บก่อน ค่อย transform (ส่วนมากทำแบบนี้)

# Use Case #3: Data Analysis

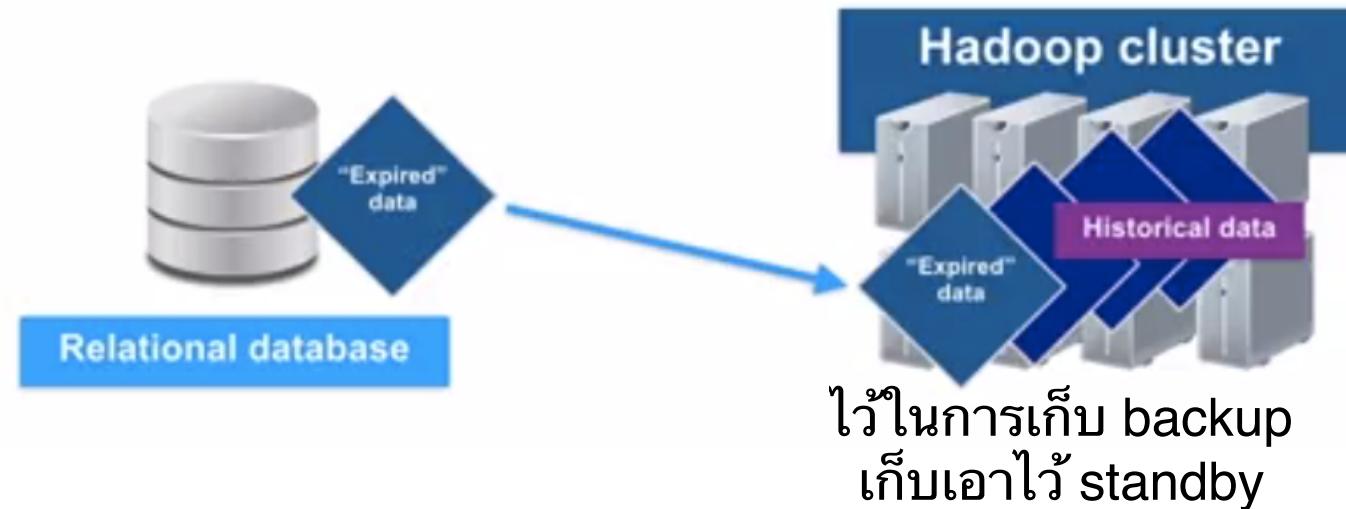
- Copy real-time data from RDBMS, combine with raw data on Hadoop using complex data analysis logic (not just SQL!)



เพิ่มความสามารถ ในการวิเคราะห์ อาจจะมีข้อมูลดิบที่เพิ่มมา ให้เรา วิเคราะห์

# Use Case #4: Data Archival

- Move data from RDBMS after it expires to Hadoop, keeping the RDBMS “clean and lean”



# Use Case #5: Data Consolidation

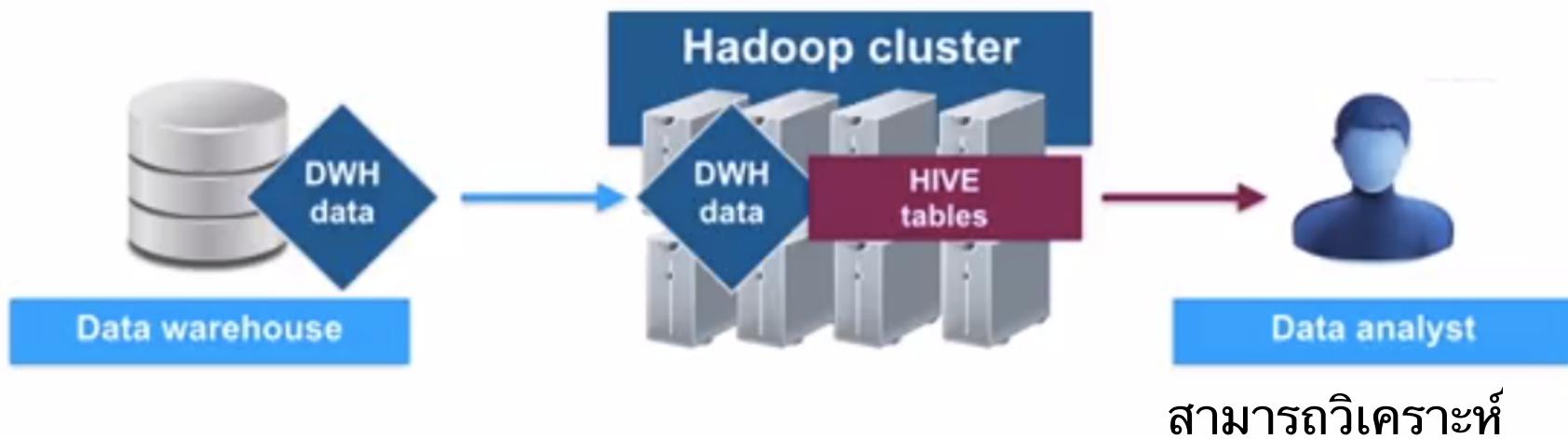
- Integrate data from various organizational “data stores” to Hadoop for various data processing requirements



เอา hadoop มาเป็น consol แบบ  
เก็บข้อมูลไว้หลายๆที่ และดึงข้อมูลให้มาอยู่ที่เดียวกัน  
เวลาประมาณข้อมูลจะง่ายขึ้น ใช้ sqoop ในการดึง

# Use Case #6: Move reports to Hadoop

- Easily allow traditional data analysis and business intelligence using Hadoop's power



# Import Commands

--connect <jdbc-uri>	Specifies the server or database to connect to. It also specifies the port. For example:  --connect jdbc:mysql://host:port/databaseName
--connection-manager <class-name>	Specifies the connection manager class name.
--driver <class-name>	Specifies the fully qualified name of the JDBC driver class.
--hadoop-home <dir>	This parameter is used to override the <code>\$HADOOP_HOME</code> environment variable.

# Import Commands

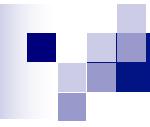
--password <password>	Sets the authentication password required to connect to the input source.
--username <username>	Sets the authentication username.
--connection-param-file <properties-file>	Specifies the connection parameter's file.
--help	This option will provide the usage instructions.
--verbose	Prints more information during a query execution.

If a user doesn't want to specify the database password along with the command, the -P option can be used to read the password from the console.

# Incremental import

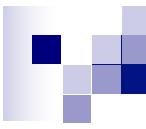
Parameter/argument	Description
--check-column <i>&lt;column-name&gt;</i>	The value of this column is used to determine the rows to be imported during the import process.
--incremental <i>&lt;incremental-type&gt;</i>	Specifies the type of incremental mode. Possible values are <code>append</code> and <code>lastmodified</code> .
--last-value <i>&lt;value&gt;</i>	Specifies the last value or the maximum value of the <code>check</code> column from the previous import. All the records whose <code>check</code> column value is greater than the value of the <code>--last-value</code> argument will be imported to HDFS.

```
bin/sqoop import -connect jdbc:mysql://localhost:3306/db1 -username root -password password -table student -target-dir /user/abc/student -columns "student_id,address,name" --incremental lastmodified -last-value "2012-11-06 19:01:35"--check-column col4
```



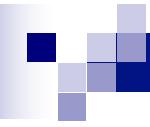
# Export Commands

Parameters	Description
<code>--direct</code>	Use the direct mode to perform the export quickly. Note that it is only supported for MySQL.
<code>--export-dir&lt;dir&gt;</code>	The location of input files in HDFS.
<code>--table &lt;table-name&gt;</code>	Name of the output table (the RDBMS table).
<code>-m, --num-mappers &lt;n&gt;</code>	Refers to the number of map tasks.



# **Hands-On: Loading Existing Data from RDBMS to Hive**

---



# Connect to MySQL in Cloudera Quickstart

**mysql -uroot -pcloudera**

```
[root@quickstart /]# mysql -uroot -pcloudera
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 5069
Server version: 5.1.73 Source distribution

Copyright (c) 2000, 2013, Oracle and/or its affiliates. All
rights reserved.

Oracle is a registered trademark of Oracle Corporation and/
or its
affiliates. Other names may be trademarks of their respecti
ve
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the curre
nt input statement.

mysql> █
```

# Show all MySQL DBs in Cloudera by Default

```
show databases;
```

```
mysql> show databases;
+-----+
| Database      |
+-----+
| information_schema |
| cm           |
| firehose      |
| hue          |
| metastore     |
| mysql         |
| nav          |
| navms         |
| oozie         |
| retail_db     |
| rman         |
| sentry        |
+-----+
12 rows in set (0.14 sec)

mysql> █
```

# Change Database and Show all Tables

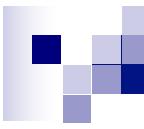
```
use retail_db;
```

```
show tables;
```

```
mysql> use retail_db;
Reading table information for completion names
You can turn off this feature to get .
-A

Database changed
mysql> show tables;
+-----+
| Tables_in_retail_db |
+-----+
| categories          |
| customers           |
| departments         |
| order_items         |
| orders               |
| products             |
+-----+
6 rows in set (0.00 sec)

mysql> ■
```



## Importing the “orders” table to Hive

import export เป็นการทำ  
map reduce

map reduce ตรง -m 1  
เปลี่ยนเป็น -m2,3,4

```
$sqoop import --connect "jdbc:mysql://quickstart.cloudera:3306/retail_db"  
--username root --password cloudera --table orders --hive-import  
--hive-overwrite --create-hive-table -m 1
```

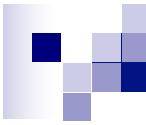
Data sample for orders

[View more...](#)

	orders.order_id	orders.order_date	orders.order_customer_id	orders.order_status
1	1	2013-07-25 00:00:00.0	11599	CLOSED
2	2	2013-07-25 00:00:00.0	256	PENDING_PAYMENT
3	3	2013-07-25 00:00:00.0	12111	COMPLETE
4	4	2013-07-25 00:00:00.0	8827	CLOSED
5	5	2013-07-25 00:00:00.0	11318	COMPLETE
6	6	2013-07-25 00:00:00.0	7130	COMPLETE
7	7	2013-07-25 00:00:00.0	4530	COMPLETE
8	8	2013-07-25 00:00:00.0	2911	PROCESSING
9	9	2013-07-25 00:00:00.0	5657	PENDING_PAYMENT
10	10	2013-07-25 00:00:00.0	5648	PENDING_PAYMENT
11	11	2013-07-25 00:00:00.0	918	PAYMENT REVIEW
12	12	2013-07-25 00:00:00.0	1837	CLOSED
13	13	2013-07-25 00:00:00.0	9149	PENDING_PAYMENT

# **Exercise** : Use Impala to show all product names that have price = 50

```
+-----+  
| product_name |  
+-----+  
| Nike Adult Vapor Jet 3.0 Receiver Gloves |  
| Reebok Women's Performance Regular Motion Fit |  
| Nike Men's Dri-FIT Victory Golf Polo |  
| adidas Men's Germany Home Replica White Top |  
| adidas Men's Germany Home Replica White Top |  
| Stitches Men's Pittsburgh Pirates Black Poly |  
| adidas Men's Germany Away Replica Black Top |  
| Majestic Youth Replica New York Yankees Derek |  
| Majestic Youth Replica Pittsburgh Pirates And |  
| Majestic Youth Replica Boston Red Sox David O |  
| Majestic Youth Replica Los Angeles Dodgers Ya |  
| Majestic Youth Replica New York Yankees Jacob |  
+-----+  
Fetched 12 row(s) in 0.41s  
[quickstart.cloudera:21000] > 
```



# **Hands-On: Loading Data from RDBMS to Hadoop**

---

# Running MySQL Docker

```
root@hadoop-docker:~# docker pull mysql:5.6.40
root@hadoop-docker:~# docker run --name mysqlServer
-e MYSQL_ROOT_PASSWORD=itkm1t -p 3306:3306
-d mysql:5.6.40
root@hadoop-docker:~# docker exec -it mysqlServer bash
```

```
root@262789f9c117:/# mysql -uroot -p"itkm1t"
```

Oracle is a registered trademark of Oracle Corporation and/or its affiliates. Other names may be trademarks of their respective owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

```
mysql> 
```

# Prepare a test database table

```
mysql> CREATE DATABASE itkmrl_db;
mysql> USE itkmrl_db;
mysql> CREATE TABLE country_tbl(id INT NOT NULL, country
VARCHAR(50), PRIMARY KEY (id));
```

```
mysql> INSERT INTO country_tbl VALUES(1, 'USA');
mysql> INSERT INTO country_tbl VALUES(2, 'CANADA');
mysql> INSERT INTO country_tbl VALUES(3, 'Mexico');
mysql> INSERT INTO country_tbl VALUES(4, 'Brazil');
mysql> INSERT INTO country_tbl VALUES(61, 'Japan');
mysql> INSERT INTO country_tbl VALUES(65, 'Singapore');
mysql> INSERT INTO country_tbl VALUES(66, 'Thailand');
```

## View data in the table

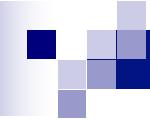
```
mysql> SELECT * FROM country_tbl;
```

id	country
1	USA
2	CANADA
3	Mexico
4	Brazil
61	Japan
65	Singapore
66	Thailand

```
7 rows in set (0.00 sec)
```

```
mysql> exit;
```

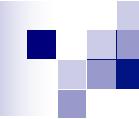
Then exit from the container by press Ctrl-P & Ctrl-Q



## Inspect imcMysql IP Address

```
root@hadoop-docker:~# docker inspect mysqlServer | grep IPAddress
```

```
"IPAddress": "172.17.0.7",
```



## Restart the Cloudera docker with linking to the MySQL Docker

---

```
root@hadoop-docker:~# docker run --hostname=quickstart.cloudera  
--privileged=true --link mysqlServer:mysqlDb -t -i -p  
8888:8888 cloudera/quickstart /usr/bin/docker-quickstart
```

---

# Case I: Importing data from MySQL to HDFS

```
[root@quickstart /]# sqoop import --connect  
jdbc:mysql://172.17.0.7/itkmil_db --username root --password  
itkmil --table country_tbl --target-dir  
/user/cloudera/testtable -m 1
```

The screenshot shows the Apache Hue interface with the 'File Browser' tab selected. The top navigation bar includes links for 'HUE', 'Home', 'Query Editors', 'Data Browsers', 'Workflows', 'Search', and 'Security'. Below the navigation is a toolbar with icons for file operations. The main area is titled 'File Browser' and shows a list of files under the path '/user/cloudera/testtable/part-m-00000'. A red oval highlights this path. On the left, a sidebar titled 'ACTIONS' lists options: 'View as binary', 'Edit file', 'Download', 'View file location', and 'Refresh'. The right side displays the contents of the file, which is a text file containing the following data:

Country ID	Country Name
1	USA
2	CANADA
3	Mexico
4	Brazil
61	Japan
65	Singapore
66	Thailand

## Case II: Importing data from MySQL to HBase

```
[root@quickstart /]# sqoop import --connect  
jdbc:mysql://172.17.0.7/itkmil_db --username root --password  
itkmil --table country_tbl --hbase-table country --column-  
family hbase_country_cf --hbase-row-key id --hbase-create-  
table -m 1
```

### Start HBase

```
root@quickstart /]# hbase shell  
hbase(main):001:0> list
```

```
TABLE  
country  
1 row(s) in 0.2570 seconds  
  
=> ["country"]
```

# Viewing HBase data

```
hbase(main):003:0> scan 'country'
ROW                                COLUMN+CELL
 1        column=hbase_country_cf:country, timestamp=1468081466623, value=USA
 2        column=hbase_country_cf:country, timestamp=1468081466623, value=CANADA
 3        column=hbase_country_cf:country, timestamp=1468081466623, value=Mexico
 4        column=hbase_country_cf:country, timestamp=1468081466623, value=Brazil
 61       column=hbase_country_cf:country, timestamp=1468081466623, value=Japan
 65       column=hbase_country_cf:country, timestamp=1468081466623, value=Singapore
 66       column=hbase_country_cf:country, timestamp=1468081466623, value=Thailand
7 row(s) in 0.1670 seconds
```

# Viewing data from HBase browser

The screenshot shows the Hue HBase Browser interface. The top navigation bar includes links for Home, Query Editors, Data Browsers, Workflows, Search, and Security, along with various icons for file operations and cluster management.

The main title is "HBase Browser". Below it, the path "Home - Cluster / country" is displayed, with a "Switch Cluster" dropdown menu.

The search bar contains the query: "row\_key, row\_prefix\* +scan\_len [col1, family:col2, fam3:, col\_prefix\*]" and a search icon. To the right of the search bar are buttons for "hbase\_country\_cf", "Filter Columns/Families", "All" (checked), and "Sort By ASC".

The data table displays three rows of data:

- Row 1: hbase\_country\_cf: country, USA
- Row 2: hbase\_country\_cf: country, CANADA
- Row 3: hbase\_country\_cf: country, Mexico

## Case III: Importing data from MySQL to Hive Table

```
root@quickstart /]# sqoop import --connect  
jdbc:mysql://172.17.0.7/itkmctl_db --username root --password  
itkmctl --table country_tbl --hive-import --hive-table  
country -m 1
```

The screenshot shows the Hue interface with the 'File Browser' tab selected. The URL in the address bar is `/user/hive/warehouse/country/part-m-00000`. A red oval highlights this path. The page displays a list of country names:

Country
USA
CANADA
Mexico
Brazil
Japan
Singapore
Thailand

# Reviewing data from Hive Table

```
[root@quickstart /]# hive
```

```
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties
```

```
WARNING: Hive CLI is deprecated and migration to Beeline is recommended.
```

```
hive> show tables;  
hive> select * from country;
```

```
1      USA  
2      CANADA  
3      Mexico  
4      Brazil  
61     Japan  
65     Singapore  
66     Thailand
```

```
Time taken: 0.587 seconds, Fetched: 7 row(s)
```

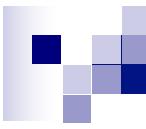
# Running from Hue: Beewax

The screenshot shows the Hue interface with the 'Query Editor' tab selected. In the top navigation bar, the 'Hive Editor' tab is also visible. On the left, there's a sidebar with 'Assist' and 'Settings' tabs, and a 'Tables' section showing a single table named 'country'. The main area contains a code editor with the following SQL query:

```
1 | SELECT * FROM country;
```

Below the code editor are several action buttons: 'Execute' (highlighted in blue), 'Save as...', 'Explain', 'Format', and 'or create a New query'. The results section is titled 'Results' and displays the data from the 'country' table:

	country.id	country.country
1	1	USA
2	2	CANADA
3	3	Mexico
4	4	Brazil
5	61	Japan
6	65	Singapore
7	66	Thailand



# **Hands-On: Exporting Data from Hadoop to RDBMS**

---

# Exporting data from HDFS to MySQL

```
[root@quickstart /]# wget  
https://s3.amazonaws.com/imcbucket/data/country.txt  
  
[root@quickstart /]# hdfs dfs -put country.txt  
/user/cloudera/input/  
  
[root@quickstart /]# sqoop export --connect  
jdbc:mysql://172.17.0.7/itkmmitl_db --username root --password  
itkmmitl --table country_tbl --export-dir  
/user/cloudera/input/country.txt
```

## Viewing Result in a MySQL table on a Ubuntu server

```
root@hadoop-docker:~# docker exec -it mysqlServer bash  
root@f1922a70e09c:/# mysql -uroot -p"itkmmitl"  
mysql> USE itkmmitl_db;  
mysql> SELECT * FROM country_tbl;
```