

Module 7

Understanding Sqoop

Thanachart Numnonda, Executive Director, IMC Institute

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

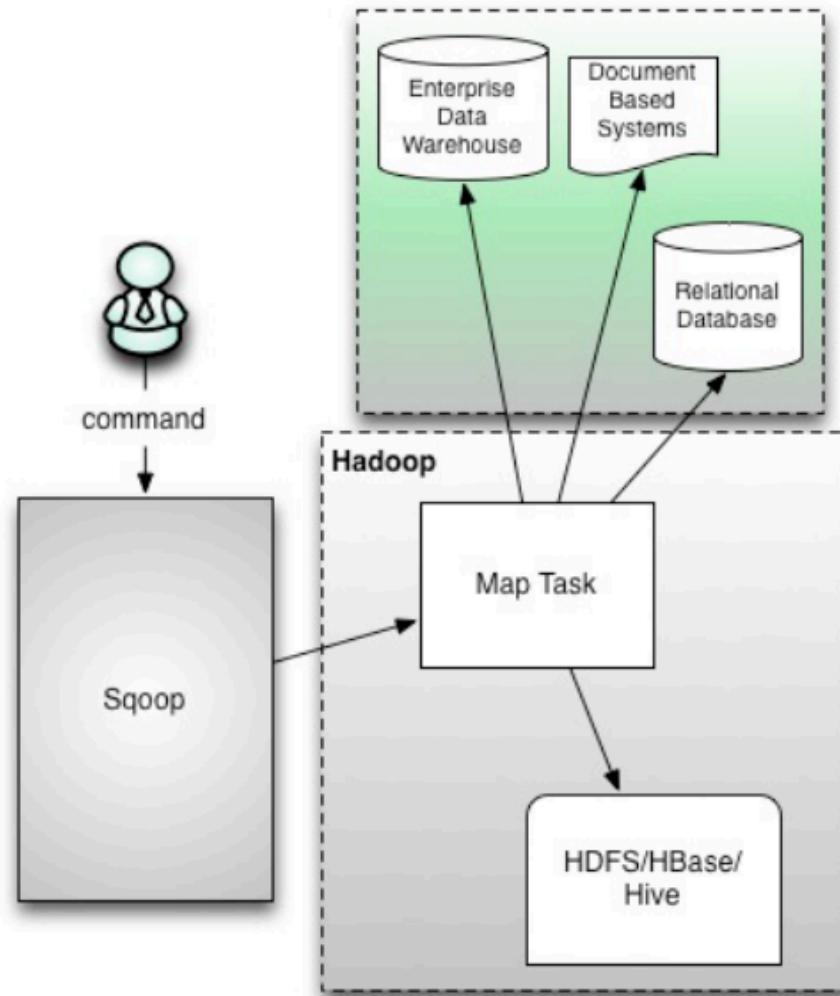
Introduction

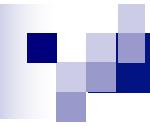


Sqoop (“SQL-to-Hadoop”) is a straightforward command-line tool with the following capabilities:

- .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**

Architecture Overview



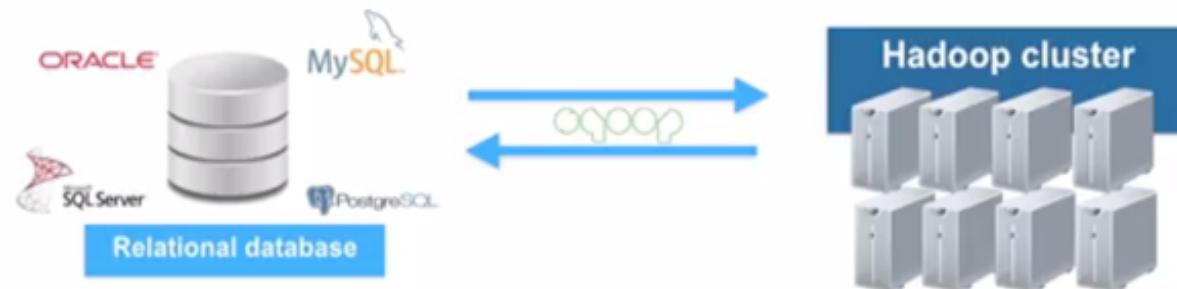


Sqoop Benefit

- 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 daily transactional data to Hadoop
- It uses MapReduce for export and import that enables parallel and efficient data movement

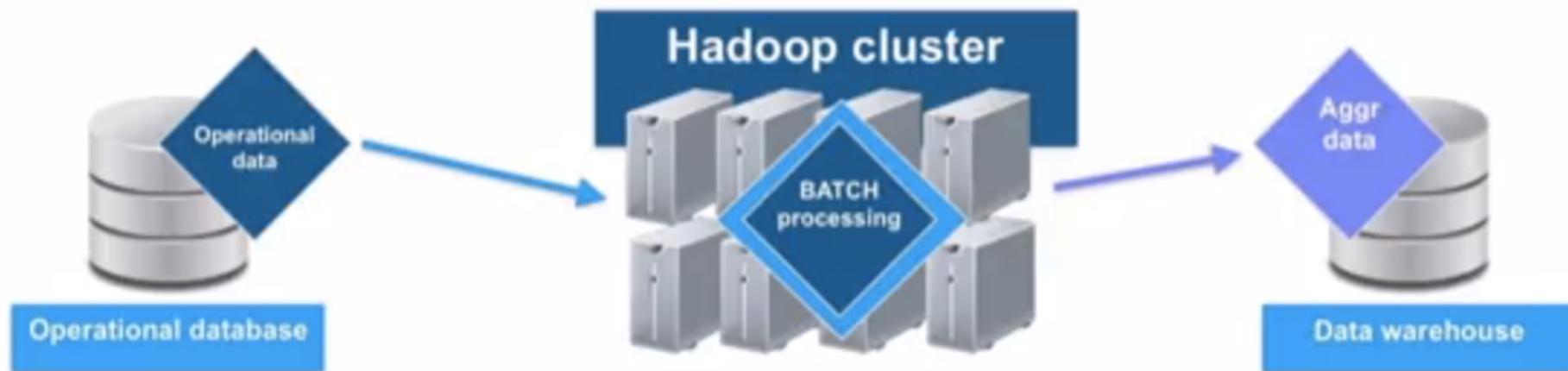
Sqoop Mode

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



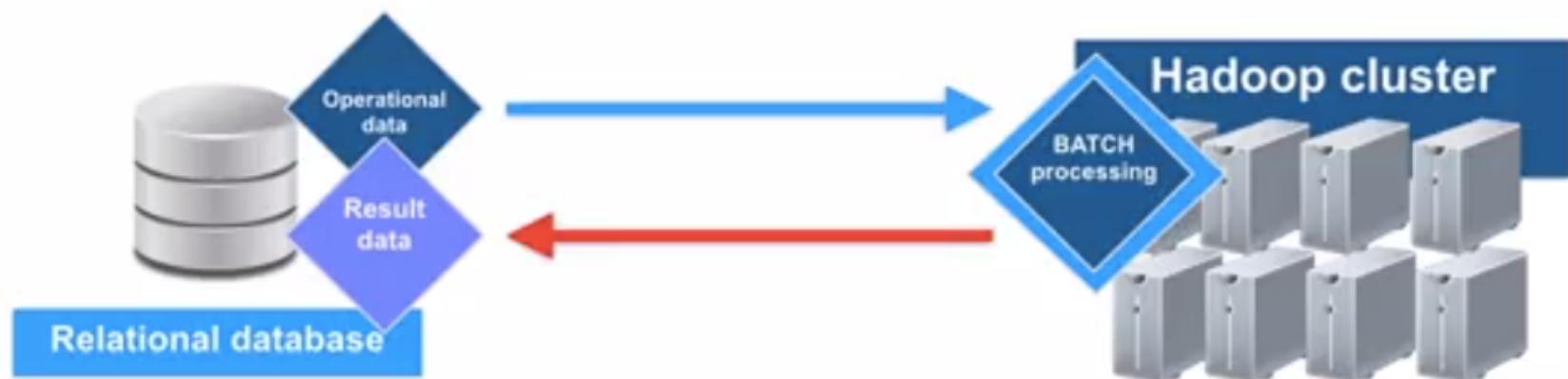
Use Case #1: ETL for Data Warehouse

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



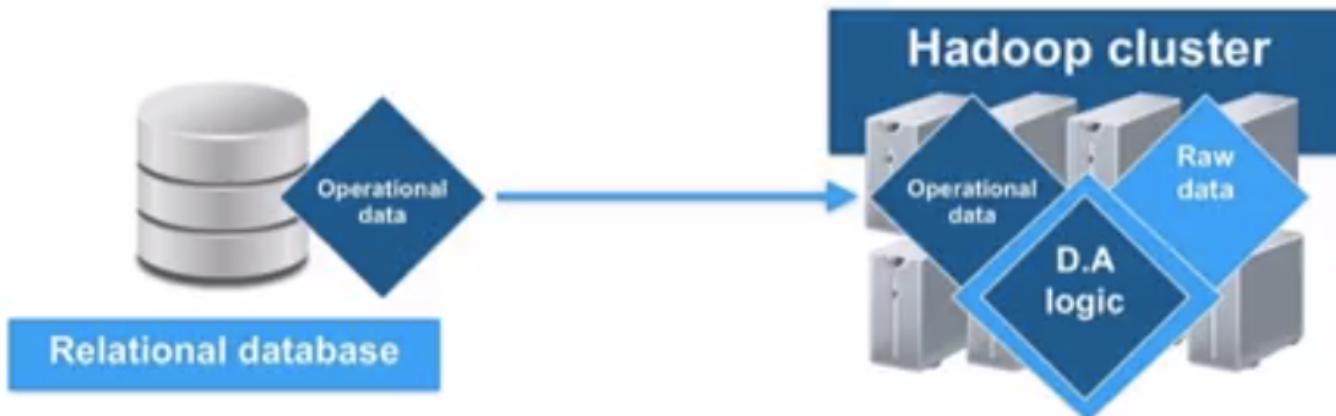
Use Case #2: ELT

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



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!)



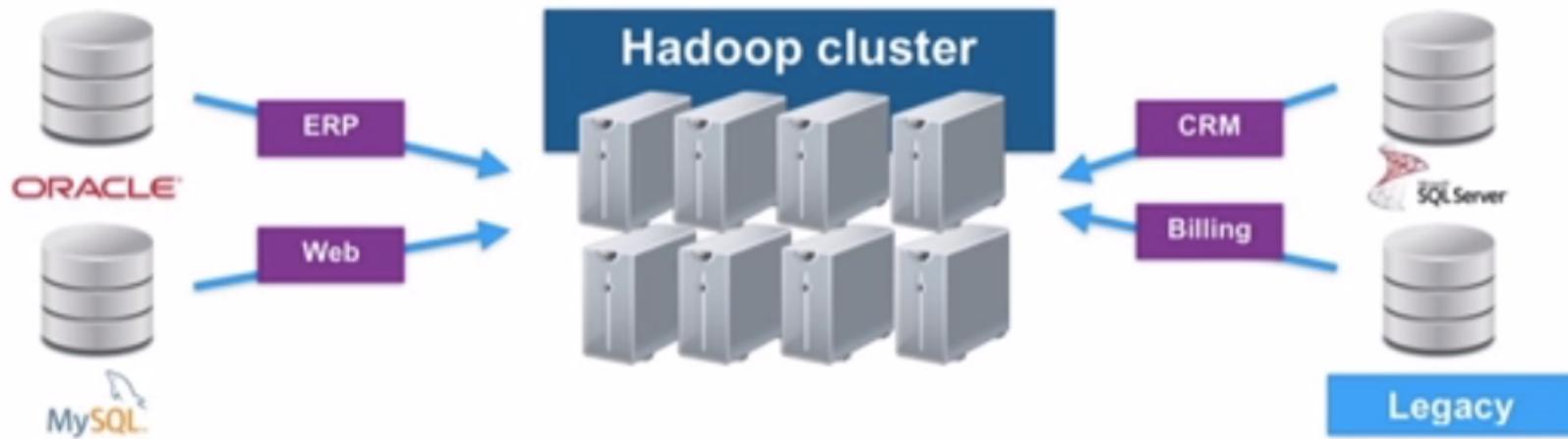
Use Case #4: Data Archival

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



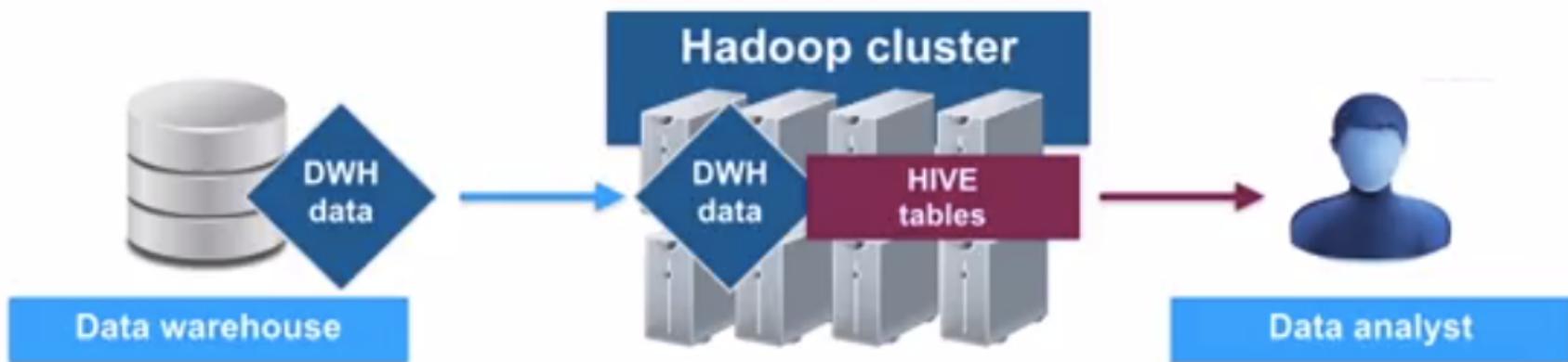
Use Case #5: Data Consolidation

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



Use Case #6: Move reports to Hadoop

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



Import Commands

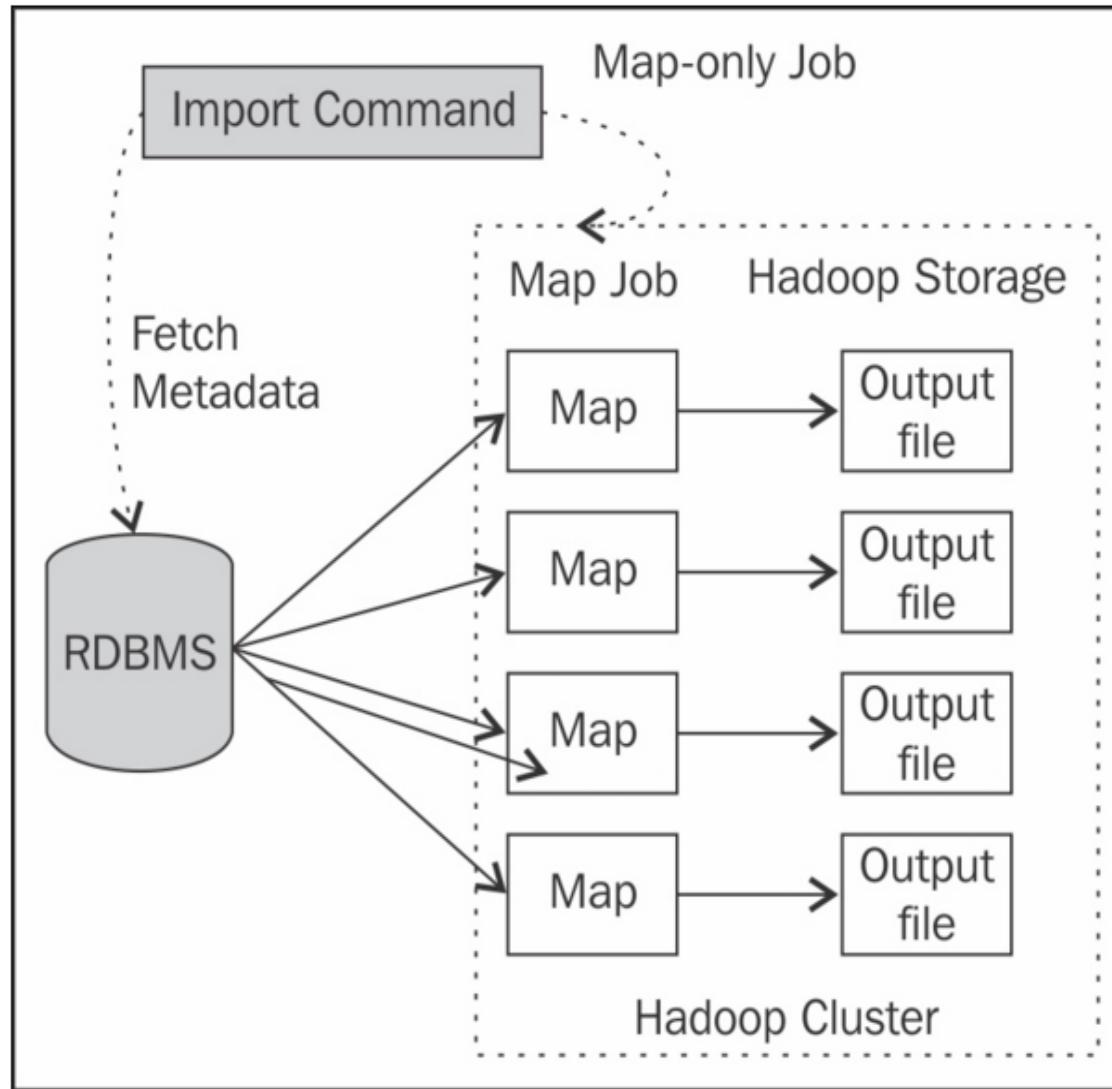
--connect <jdbc-uri>	Specifies the server or database to connect to. It also specifies the port. For example: --connect jdbc:mysql://host:port/databaseName
--connection-manager <class-name>	Specifies the connection manager class name.
--driver <class-name>	Specifies the fully qualified name of the JDBC 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.

Architecture of the import process



Incremental import

Parameter/argument	Description
--check-column <i><column-name></i>	The value of this column is used to determine the rows to be imported during the import process.
--incremental <i><incremental-type></i>	Specifies the type of incremental mode. Possible values are <code>append</code> and <code>lastmodified</code> .
--last-value <i><value></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.

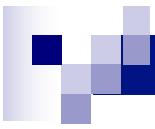
```
bin/sqoop import -connect jdbc:mysql://localhost:3306/db1 -username root -password password -table student -target-dir /user/abc/student -columns "student_id,address,name" --incremental lastmodified -last-value "2012-11-06 19:01:35"--check-column col4
```

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<dir></code>	The location of input files in HDFS.
<code>--table <table-name></code>	Name of the output table (the RDBMS table).
<code>-m,--num-mappers <n></code>	Refers to the number of map tasks.

Export Commands

--update-mode <mode>	Specifies how updates are performed when new rows are found with non-matching keys in the database. Legal values for the mode include updateonly (default) and allowinsert .
--update-key <col-name>	The value of this column is used to identify the records that a user wants to update during the update mode. Use a comma-separated list of columns if there is more than one column.
--staging-table <staging-table-name>	Specifies the name of the staging table. The staging table is used to stage the data before inserting it into the destination table.
--clear-staging-table	This argument is used to clean the data from the staging table.



Hands-On: Loading Data from RDBMS to Hadoop

Running MySQL Docker

- **UbuntuCmd:** sudo docker pull mysql
- **UbuntuCmd:** sudo docker run --name mysqlServer -e MYSQL_ROOT_PASSWORD=itkm1t -p 3306:3306 -d mysql
- **UbuntuCmd:** sudo docker exec -it mysqlServer bash

```
root@f1922a70e09c:/# 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 itkmitl_db;
mysql> USE itkmitl_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

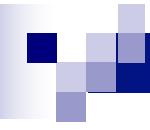
```
mysql> SELECT * FROM country_tbl;
```

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

UbuntuCmd: sudo docker inspect mysqlServer | grep IPAddress

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

Restart the Cloudera docker with linking to the MySQL Docker

UbuntuCmd: `sudo docker run --hostname=quickstart.cloudera --privileged=true --link mysqlServer:mysqldb -t -i -p 8888:8888 cloudera/quickstart /usr/bin/docker-quickstart`

If both of these Dockers are up and running, you can find out the internal IP address of each of them by running this command. This gets the IP for mysqlServer.

Case I: Importing data from MySQL to HDFS

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

The screenshot shows the Hue File Browser interface. The top navigation bar includes links for 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, there is a sidebar with "ACTIONS" and "INFO" sections, each containing several options like View as binary, Edit file, Download, View file location, and Refresh. The "INFO" section lists the contents of the file, which are the data rows imported from MySQL.

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

Case II: Importing data from MySQL to HBase

```
$sqoop import --connect jdbc:mysql://172.17.0.7/itkmittl_db --  
username root --password itkmittl --table country_tbl --hbase-  
table country --column-family hbase_country_cf --hbase-row-key id --  
hbase-create-table -m 1
```

Start HBase

```
$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. 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. A filter section shows "hbase_country_cf:" selected. The results table displays three rows of data:

row_key	hbase_country_cf: country
1	USA
2	CANADA
3	Mexico

Each row is numbered (1, 2, 3) and contains a "Filter Columns/Families" button, a checked "All" checkbox, and a "Sort By ASC" dropdown.

Case III: Importing data from MySQL to Hive Table

```
$sqoop import --connect jdbc:mysql://172.17.0.7/itkmittl_db --  
username root --password itkmittl --table country_tbl --hive-  
import --hive-table country -m 1
```

The screenshot shows the Apache Hue interface, specifically the File Browser. The top navigation bar includes links for Home, Query Editors, Data Browsers, Workflows, Search, and Security. Below the navigation is a toolbar with icons for file operations like View as binary, Edit file, Download, View file location, and Refresh. The main area is titled "File Browser" and displays a list of files under the path "/user/hive/warehouse/country/part-m-00000". The list contains the following data:

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

A red oval highlights the path in the breadcrumb navigation bar: / user / hive / warehouse / country / part-m-00000.

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)
```

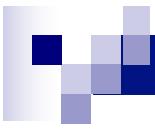
Running from Hue: Beewax

The screenshot shows the Hue interface with the 'Query Editor' tab selected. In the query editor area, a single line of SQL code is entered:

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

Below the query, there are several action buttons: 'Execute' (highlighted in blue), 'Save as...', 'Explain', 'Format', and 'or create a New query'. The results section is active, displaying the following 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

```
$wget https://s3.amazonaws.com/imcbucket/data/country.txt
$hdfs dfs -put country.txt /user/cloudera/input/
$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

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