

FINISHED ▶ ❌ 🔍 ⚙

- Took 0 sec. Last updated by anonymous at February 22 2022, 1:23:37 AM

FINISHED ▶ ❌ 🔍 ⚙

FINISHED ▶ 其他命令

```
-- show table meta data
--=====
```

DESCRIBE FORMATTED WSS...

Took 0 sec. Last updated by anonymous on February 17 2022, 7:58:58 PM. (undetected)

```
-- counts the number of rows from wdi_gs table
```

FINISHED ▶ ⌂ ⌕ ⌕

```
-- counts the number of rows from wdi_gs table
```

FINISHED ▶ ❌ 📄

count
21759408

Book 34 sec. Last updated by anonymous at February 17 2022, 7:42:55 PM. Counted

FINISHED ▶ ⌵ ⌶ ⌷

Took 0 sec. Last updated by anonymous at February 22, 2022, 1:23:54 AM.

• Query that loads data from wdi_gs table to wdi_csv_text table.

FINISHED ▶ 🔍 📄

Took 0 sec. Last updated by anonymous at February 22 2020, 1:20:22 AM.

• Check HDFS file size for wdi_csv_text file FINISHED

took 0 sec. Last updated by anonymous on February 22, 2020, 1:28:17 AM.

```
hdfs dfs -ls -h /user/garslan/hive/hd/wd_csv_test
```

FINISHED

▶ ◀ ⌂

```
Found 5 items
-rw-r--r-- 2 garslan hadoop 385.0 K 2022-02-18 06:34 /user/garslan/hive/hd/wd_csv_test/000000_0
-rw-r--r-- 2 garslan hadoop 385.8 K 2022-02-18 06:33 /user/garslan/hive/hd/wd_csv_test/000001_0
-rw-r--r-- 2 garslan hadoop 386.0 K 2022-02-18 06:34 /user/garslan/hive/hd/wd_csv_test/000002_0
-rw-r--r-- 2 garslan hadoop 385.6 K 2022-02-18 06:34 /user/garslan/hive/hd/wd_csv_test/000003_0
-rw-r--r-- 2 garslan hadoop 218.4 K 2022-02-18 06:34 /user/garslan/hive/hd/wd_csv_test/000004_0
```

Task 3 was last updated by anonymous at February 18 2022, 5:27:08 AM.

```
hsp
hdfs dfs -du -s -h /user/garshan/hive/hdfs/hdfs_test
1.7 G  3.4 G  /user/garshan/hive/hdfs/hdfs_test
Task 5 sec. Last updated by anonymous on February 18 2020, 1:57:59 AM.
```

```
SELECT count(countryName) as count FROM wsl_csv_test
```

FINISHED D [0] 0

```
DMP0 : Compiling command[QueryId=2028218609312_d5dc6d-c-90f0-4948-b2e6-476cd4ee731]: SELECT count(countryName) as count FROM wsl_csv_test
DMP0 : Concurrency mode is disabled, not creating a lock manager
DMP0 : Semantic Analysis Completed (retires = 1x)
DMP0 : Returning Hive schema: Schema{fieldschema=[fieldschema{name(count, type:string, comment:null), properties=null}]
DMP0 : Compiling Command[QueryId=2028218609312_d5dc6d-c-90f0-4948-b2e6-476cd4ee731]: Time taken: 0.262 seconds
DMP0 : Concurrency mode is disabled, not creating a lock manager
DMP0 : Executing Command[QueryId=2028218609312_d5dc6d-c-90f0-4948-b2e6-476cd4ee731]: SELECT count(countryName) as count FROM wsl_csv_test
DMP0 : Query ID = hive_2028218609312_d5dc6d-c-90f0-4948-b2e6-476cd4ee731
DMP0 : Total Spgs = 1
DMP0 : Launching Sps 1 out of 1
DMP0 : Starting task [Stage-1:MAPRED] in serial mode
DMP0 : Submitted to cluster(s) {hive_queryId=hive_2028218609312_d5dc6d-c-90f0-4948-b2e6-476cd4ee731}
DMP0 : Task session hasn't been created yet. Opening session
DMP0 : Dag name = com.dummy.as ...wsl_csv_test [Stage-1]
DMP0 : Status: Running (Executing on YARN cluster with App id application_1645146646832_805)
```



• Clear filesystem cache and execute the count query again.

```

echo 3 | sudo tee /proc/sys/vm/drop_caches

We trust you have received the usual lecture from the Local System
Administrator. It usually boils down to these three things:

#1) Respect the privacy of others.
#2) Think before you type.
#3) With great power comes great responsibility.

sudo: no tty present and no explicit program specified

ExitValue: 1

```

Hive vs Bash FINISHED ▶ 🔍 📄 🌐

bash took 3 seconds to execute, while Hive took 20 seconds. The total time of the Hive technique was longer. Because of the cost of parsing queries, developing an implementation of execution plan, and performing Hadoop Map Reduce task.

Tool Used: Last updated by anonymous at February 22 2022, 1:24:07 AM.

2. Create a Table with OpenCSV SerDe

Task 0 sec. Last updated by anonymous at February 22 2022, 1:56:46 AM

hive

ERROR

-- 2a. Create wdl_opencsv_gs source table (load 65 data with OpenCSVSerde)

CREATE EXTERNAL TABLE wdl_opencsv_gs
(year INT, countryname STRING, countrycode STRING,
indicatorname STRING, indicatorcode STRING, indicatorvalue
FLOAT)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
LOCATION 'gs://hive1-data-ops-us-east-1-bucket/wdl_2015'
TBLPROPERTIES ("tableheader.line.count"="1")

Error while processing statement: FAILED: Execution error, return code 1 from org.apache.hadoop.hiveql.exec.DDLTask. Already exists function(message:Table hive.default.wdl_opencsv_gs already exists)

Task 1 sec. Last updated by anonymous at February 18 2022, 1:57:09 AM

hive

FINISHED

-- 2b. Create wdl_opencsv_test destination table (upload data with HDFS localfile)

CREATE EXTERNAL TABLE wdl_opencsv_test
(year INT, countryname STRING, countrycode STRING,
indicatorname STRING, indicatorcode STRING, indicatorvalue
FLOAT)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
STORED AS TEXTFILE
LOCATION 'hdfs://user/garima/hive/wdl/wdl_csv_test'

Query executed successfully. Affected rows: 1

Task 0 sec. Last updated by anonymous at February 18 2022, 1:57:19 AM

hive

FINISHED

-- 2c. A HiveQL which load data from wdl_opencsv_gs --
INSERT OVERWRITE Table wdl_opencsv_test
SELECT * FROM wdl_opencsv_gs
-- 23f44ba-9d2b-57b0d88b8a11

INFO : Compiling command(query=hive_20220218064717_2460751f...)
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldschema=[fieldschema=year, type=string, comment=null], fieldschema=countryname, type=string, comment=null), fieldschema=countrycode, type=string, comment=null), fieldschema=indicatorname, type=string, comment=null), fieldschema=indicatorcode, type=string, comment=null), fieldschema=indicatorvalue, type=string, comment=null], properties=null)
Query executed successfully. Affected rows: 1

Task 2 min 12 sec. Last updated by anonymous at February 18 2022, 1:58:30 AM

hive

ABORT

-- 2d. Verifying the data parsing has been successful --
SELECT distinct(indicatorcode)
FROM wdl_opencsv_test
ORDER BY indicatorcode
LIMIT 20

INFO : Compiling command(query=hive_20220218064831_8a228462...)
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldschema=[fieldschema=year, type=string, comment=null], fieldschema=countryname, type=string, comment=null), fieldschema=countrycode, type=string, comment=null), fieldschema=indicatorname, type=string, comment=null), fieldschema=indicatorcode, type=string, comment=null), fieldschema=indicatorvalue, type=string, comment=null], properties=null)
Query cancelled

Task 1 min 47 sec. Last updated by anonymous at February 18 2022, 1:57:12 AM

hive

FINISHED

-- Comparison of execution time between wdl_opencsv_test and wdl_csv_test
-- Usage of opencsvSerde makes execution slower than LazyInputSerDe

SELECT count(countryname) FROM wdl_opencsv_test

INFO : Compiling command(query=hive_20220222090632_c3a77cda-49f8-48e2-86f6-1bd686457288):

SELECT count(countryname) FROM wdl_opencsv_test

INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldschema=[fieldschema=year, type=bigint, comment=null], properties=null)
INFO : Completed compiling command(query=hive_20220222090632_c3a77cda-49f8-48e2-86f6-1bd686457288); Time taken: 2.874 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
INFO : Executing command(query=hive_20220222090632_c3a77cda-49f8-48e2-86f6-1bd686457288):

21759408

21759408

Task 1 min 48 sec. Last updated by anonymous at February 22 2022, 12:08:21 AM (outbreak)

OpenCSVSerde limitaion

Task 0 sec. Last updated by anonymous at February 22 2022, 1:24:52 AM

Compare metadata of two file

Task 0 sec. Last updated by anonymous at February 22 2022, 1:24:57 AM

hive

DESCRIBE FORMATTED wdl_opencsv_test

col_name	data_type	comment
# col_name	data_type	comment
year	string	from deserializer
countryname	string	from deserializer
countrycode	string	from deserializer
indicatorname	string	from deserializer
indicatorcode	string	from deserializer
indicatorvalue	string	from deserializer
	null	null

Task 1 sec. Last updated by anonymous at February 22 2022, 1:31:27 AM (outbreak)

SQL

DESCRIBE FORMATTED wdi_opencv_text

FINISHED

Settings

col_name	data_type	comment
# col_name	data_type	comment
year	int	
countryname	string	
countrycode	string	
indicatorname	string	
indicatorcode	string	
indicatorvalue	float	
	null	null

Task 1 ran. Last updated by anonymous at February 22 2022, 12:11:49 AM (updated)

SQL

-- Create a view on top of wdi_opencv_text to cast specific columns to correct data type.
DROP VIEW IF EXISTS wdi_opencv_text_view;
CREATE VIEW IF NOT EXISTS wdi_opencv_text_view
AS
SELECT cast(year as INT64), countryname, countrycode, indicatorname, indicatorcode, cast(indicatorvalue as Float) FROM wdi_opencv_text

Query executed successfully. Affected rows : -1

Query executed successfully. Affected rows : -1

Task 1 ran. Last updated by anonymous at February 22 2022, 12:12:38 AM

SQL

-- Check if the data types casted properly
DESCRIBE FORMATTED wdi_opencv_text_view

FINISHED

Settings

col_name	data_type	comment
# col_name	data_type	comment
year	int	
countryname	string	
countrycode	string	
indicatorname	string	
indicatorcode	string	
indicatorvalue	float	
	null	null

Task 2 ran. Last updated by anonymous at February 22 2022, 12:12:40 AM (updated)

2015 Canada GDP Growth HQL

Task 2 ran. Last updated by anonymous at February 22 2022, 1:23:05 AM

Find 2015 GDP growth (annual %) for Canada.

Task 2 ran. Last updated by anonymous at February 22 2022, 1:23:13 AM

SQL

SELECT year ,countryname , indicatorvalue AS GDP_growth_value
FROM wdi_opencv_text_view
WHERE indicatorname LIKE "GDP growth" AND year =2015 AND countryname = "Canada"

FINISHED

Settings

INFO : Compiling command(queryTohive_2022022205321321_0a7e4095-c07f-474e-8f30-5337e6d3bf77):
SELECT year ,countryname , indicatorvalue AS GDP_growth_value
FROM wdi_opencv_text_view
WHERE indicatorname LIKE "GDP growth" AND year =2015 AND countryname = "Canada"

INFO : Currency mode is disabled, not creating a lock manager
INFO : Semantic Analysis Completed (retrial = false)
INFO : Returning Hive schema: Schema(fieldSchemas:[fieldSchema(name=year, type=int, comment=null), fieldSchema(name=countryname, type=string, comment=null), fieldSchema(name=gdp_growth_value, type=float, comment=null)], properties=null)
INFO : Completed compiling command(queryTohive_2022022205321321_0a7e4095-c07f-474e-8f30-5337e6d3bf77); Time taken: 0.373 seconds
INFO : Currency mode is disabled, not creating a lock manager
INFO : Executing command(queryTohive_2022022205321321_0a7e4095-c07f-474e-8f30-5337e6d3bf77):
SELECT year ,countryname , indicatorvalue AS GDP_growth_value
FROM wdi_opencv_text_view
WHERE indicatorname LIKE "GDP growth" AND year =2015 AND countryname = "Canada"

INFO : Query ID = hive_2022022205321321_0a7e4095-c07f-474e-8f30-5337e6d3bf77
INFO : Total jobs = 1

Settings

year	countryname	gdp_growth_value
2015	Canada	1.0782658

Took 0 sec. Last updated by anonymous at February 22 2022, 1:25:19 AM.

Took 0 sec. Last updated by anonymous at February 22 2022, 1:25:24 AM

[illegible]

year	countryname	gdp_growth_value
2015	Canada	1.0782688

- Optimizing HQL query using columnar file

[illegible]

Highest GDP Growth

Task 0 sec. Last updated by anonymous at February 22 2022, 1:03:58 AM.

```
Spark
--Finding the highest GDP growth (w/ GDP_HTF_V0.30) year for each country.
SELECT w.year, w.countryname, w.indicatorvalue AS gdp_growth_value
FROM
(
  SELECT max(indicatorvalue) as value, countryname
  FROM wdt_cw_pargset
  WHERE indicatorcode = 'w.gdp.htf.v0.30' AND indicatorvalue != 0
  GROUP BY countryname )t
WHERE (SELECT max(wt_cw_pargset.w ON t.value = w.indicatorvalue AND t.countryname=w.countryname
ORDER BY gdp_growth_value DESC);

INFO : Compiling command[query2=task_202202200207_0c40209c-6020-41a8-995c-d7c6a0d55702]; Time taken: 22.900 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
```

year	countryname	gdp_growth_value
1965	Oman	109.82993
1997	Equatorial Guinea	149.97296
1997	Liberia	106.27961
2012	Libya	104.486786
1996	Bosnia and Herzegovina	88.957664
1990	Iraq	57.81783
1974	Kiribati	45.362753
1974	Gabon	39.487095

Task 7 sec. Last updated by anonymous at February 22 2022, 1:03:12 AM (outdated)

```
Spark Job FINISHED
Spark SQL
--executing the query using spark to compare execution time
SELECT w.year, w.countryname, w.indicatorvalue AS gdp_growth_value
FROM
(
  SELECT max(indicatorvalue) as value, countryname
  FROM wdt_cw_pargset
  WHERE indicatorcode = 'w.gdp.htf.v0.30' AND indicatorvalue != 0
  GROUP BY countryname )t
WHERE (SELECT max(wt_cw_pargset.w ON t.value = w.indicatorvalue AND w.countryname=t.countryname
ORDER BY gdp_growth_value DESC);

INFO : Compiling command[query2=task_202202200207_0c40209c-6020-41a8-995c-d7c6a0d55702]; Time taken: 22.900 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager
```

year	countryName	GDP_growth_value
1965	Oman	109.82993
1997	Equatorial Guinea	149.97296
1997	Liberia	106.27961
2012	Libya	104.486786
1996	Bosnia and Herzegovina	88.957664
1990	Iraq	57.81783
1974	Kiribati	45.362753
1974	Gabon	39.487095

Task 7 sec. Last updated by anonymous at February 22 2022, 1:03:01 AM (outdated)

Sort GDP by country and year

Task 0 sec. Last updated by anonymous at February 22 2022, 1:04:42 AM.

```
Spark
-- Retrieve GDP growth for all countries sorted by country name and year.
SELECT countryname, year, indicatorcode, indicatorvalue
FROM wdt_cw_pargset
WHERE indicatorcode = 'w.gdp.htf.v0.30'
ORDER BY countryname, year;

INFO : Compiling command[query2=task_202202200207_17659308-cf98-4a08-bc39-8861c0d4d7d1]; Time taken: 0.132 seconds
INFO : OK
INFO : Concurrency mode is disabled, not creating a lock manager

SELECT countryname, year, indicatorcode, indicatorvalue
FROM wdt_cw_pargset
WHERE indicatorcode = 'w.gdp.htf.v0.30'
ORDER BY countryname, year;

INFO : Semantic analysis completed (retail = false)
INFO : Returning new schema: Schema(FieldSchema(FieldSchemaName(countryname, type:string, comment:null), FieldSchemaName(year, type:int, comment:null), FieldSchemaName(indicatorcode, type:string, comment:null), FieldSchemaName(indicatorvalue, type:float, comment:null)), properties:null)
INFO : Completed compiling command[query2=task_202202200207_17659308-cf98-4a08-bc39-8861c0d4d7d1]; Time taken: 0.132 seconds
INFO : Concurrency mode is disabled, not creating a lock manager
```


countryname				year		indicatorcode	
Afghanistan				1993		NY.GDP.MKTP.KD.ZG	
Afghanistan				1994		NY.GDP.MKTP.KD.ZG	
Afghanistan				1995		NY.GDP.MKTP.KD.ZG	
Afghanistan				1996		NY.GDP.MKTP.KD.ZG	
Afghanistan				1997		NY.GDP.MKTP.KD.ZG	
Afghanistan				1998		NY.GDP.MKTP.KD.ZG	
Afghanistan				1999		NY.GDP.MKTP.KD.ZG	
Afghanistan				2000		NY.GDP.MKTP.KD.ZG	