# Hadoop Environment Ready for Developer Basic Verification

## ▼ TOC

## Overview

因為 Hadoop Cluster 乃分散式系統，配置上較為複雜，需要處理相依性。為確認 Hadoop 叢集環境元件已經正確配置完成，因此在開發前，需請系統建置者執行元件測試；以確保測試通過，就可以交接給開發者使用。

## Components of Hadoop Ecosystem

以下截圖皆為 Cloudera Hadoop (CDP) 之執行結果畫面，僅供參考。

### ▼ HDFS

- 查看資料夾

```
# 查看資料夾
hdfs dfs -ls /
```

```
[root@devm01 ~]# hdfs dfs -ls /
Found 7 items
drwxr-xr-x   - hbase hbase              0 2022-10-04 15:45 /hbase
drwxr-xr-x   - hdfs  supergroup         0 2022-09-29 16:55 /ranger
drwxrwxr-x   - solr  solr               0 2022-09-29 16:55 /solr-infra
drwxrwxrwt   - hdfs  supergroup         0 2022-09-29 17:47 /tmp
drwxr-xr-x   - hdfs  supergroup         0 2022-10-05 15:02 /user
drwxr-xr-x   - hdfs  supergroup         0 2022-09-29 16:55 /warehouse
drwxr-xr-x   - hdfs  supergroup         0 2022-09-29 16:55 /yarn
[root@devm01 ~]#
```
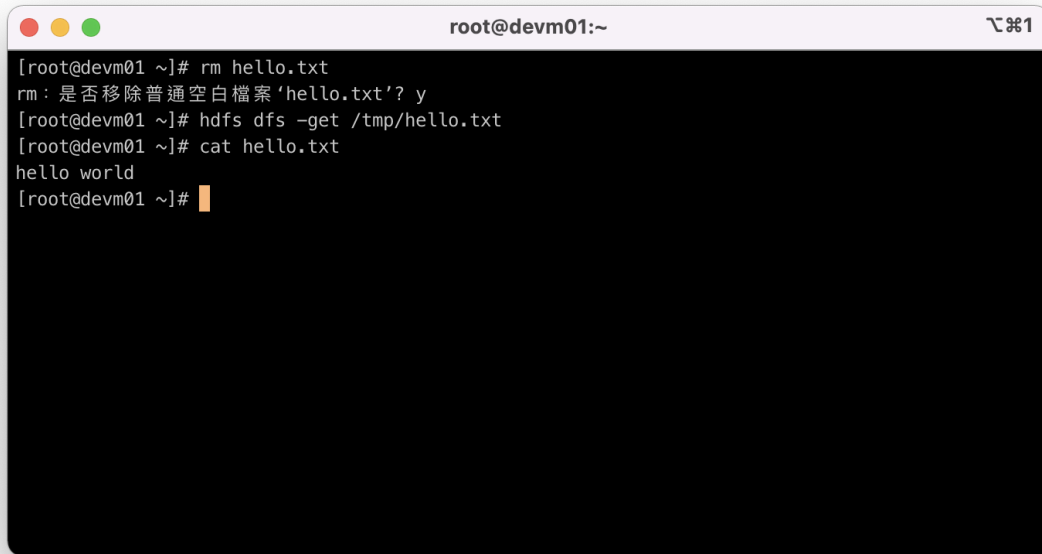
- 測試丟入檔案

```
echo "hello world" >> hello.txt
hdfs dfs -put hello.txt /tmp
hdfs dfs -ls /tmp
hdfs dfs -cat /tmp/hello.txt
```



```
[root@devm01 ~]# echo "hello world" >> hello.txt
[root@devm01 ~]# hdfs dfs -put hello.txt /tmp
[root@devm01 ~]# hdfs dfs -ls /tmp
Found 4 items
d---------   - hdfs        supergroup         0 2022-10-06 09:55 /tmp/.cloudera_health_monitori
ng_canary_files
-rw-r--r--   3 athemaster supergroup        12 2022-10-06 09:55 /tmp/hello.txt
drwx-wx-wx   - hive        supergroup         0 2022-10-03 17:24 /tmp/hive
drwxrwxrwt   - mapred      hadoop             0 2022-10-03 17:50 /tmp/logs
[root@devm01 ~]# hdfs dfs -cat /tmp/hello.txt
hello world
[root@devm01 ~]#
```

- 取出檔案

```
rm hello.txt
hdfs dfs -get /tmp/hello.txt
cat hello.txt
```

```
[root@devm01 ~]# rm hello.txt
rm：是否移除普通空白檔案 'hello.txt'? y
[root@devm01 ~]# hdfs dfs -get /tmp/hello.txt
[root@devm01 ~]# cat hello.txt
hello world
[root@devm01 ~]#
```

## ▼ Hive

- beeline

```
beeline
show databases

create external table if not exists default.users (
  name string,
  age int
);

insert into default.users values
 ("Alice", 16),  ("Bob", 18),
 ("Cindy", 53),  ("Joe", 28),
 ("Apple", 5),    ("John", 33),
 ("Roger", 45),  ("Kevin", 17),
 ("Gordon", 30), ("Paul", 33),
 ("Sandy", 39),  ("Alex", 70);

create table if not exists default.users_internal (
  name string,
  age int
);

insert into default.users_internal values
 ("Alice", 16),  ("Bob", 18),
 ("Cindy", 53),  ("Joe", 28),
 ("Apple", 5),    ("John", 33),
 ("Roger", 45),  ("Kevin", 17),
 ("Gordon", 30), ("Paul", 33),
 ("Sandy", 39),  ("Alex", 70);
```

```
select * from default.users;
select * from default.users_internal;
```

show databases;



建立 external table



寫入 external table

```
●●●                                        root@devm01:~                                          ⌥⌘1
INFO  : Executing command(queryId=hive_20221006102811_873fd2b4-0be5-4310-8406-03aec9b093a9): insert into default.users values
("Alice", 16),  ("Bob", 18),
("Cindy", 53),  ("Joe", 28),
("Apple", 5),   ("John", 33),
("Roger", 45),  ("Kevin", 17),
("Gordon", 30), ("Paul", 33),
("Sandy", 39),  ("Alex", 70)
INFO  : Query ID = hive_20221006102811_873fd2b4-0be5-4310-8406-03aec9b093a9
INFO  : Total jobs = 1
INFO  : Launching Job 1 out of 1
INFO  : Starting task [Stage-1:MAPRED] in serial mode
INFO  : Subscribed to counters: [] for queryId: hive_20221006102811_873fd2b4-0be5-4310-8406-03aec9b093a9
INFO  : Session is already open
INFO  : Dag name: insert into default.users values
("Ali...70) (Stage-1)
INFO  : Tez session was closed. Reopening...
INFO  : Session re-established.
INFO  : Session re-established.
INFO  : Status: Running (Executing on YARN cluster with App id application_1665023095293_0001)

INFO  : Status: DAG finished successfully in 9.74 seconds
INFO  :
INFO  : Query Execution Summary
INFO  : ----------------------------------------------------------------------------------------------
INFO  : OPERATION                      DURATION
INFO  : ----------------------------------------------------------------------------------------------
INFO  : Compile Query                     5.96s
INFO  : Prepare Plan                      0.67s
INFO  : Get Query Coordinator (AM)        0.05s
INFO  : Submit Plan                      16.26s
INFO  : Start DAG                         0.22s
INFO  : Run DAG                           9.74s
INFO  : ----------------------------------------------------------------------------------------------
INFO  :
INFO  : Task Execution Summary
INFO  : ----------------------------------------------------------------------------------------------
INFO  :   VERTICES     DURATION(ms)   CPU_TIME(ms)    GC_TIME(ms)   INPUT_RECORDS   OUTPUT_RECORDS
INFO  : ----------------------------------------------------------------------------------------------
INFO  :     Map 1         3036.00         4,780            140              3               1
INFO  :   Reducer 2         69.00           400              0              1               0
INFO  : ----------------------------------------------------------------------------------------------
```

建立 internal_table

```
0: jdbc:hive2://devm01.devops.com:2181,devm02> create table if not exists default.users_internal (
. . . . . . . . . . . . . . . . . . . . .>    name string,
. . . . . . . . . . . . . . . . . . . . .>    age int
. . . . . . . . . . . . . . . . . . . . .> );
INFO  : Compiling command(queryId=hive_20221021153851_c8fdd21b-efc1-424b-81fb-ded932275f86): create table if not exists default.users_internal
(
name string,
age int
)
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Created Hive schema: Schema(fieldSchemas:null, properties:null)
INFO  : Completed compiling command(queryId=hive_20221021153851_c8fdd21b-efc1-424b-81fb-ded932275f86); Time taken: 0.033 seconds
INFO  : Executing command(queryId=hive_20221021153851_c8fdd21b-efc1-424b-81fb-ded932275f86): create table if not exists default.users_internal
(
name string,
age int
)
INFO  : Starting task [Stage-0:DDL] in serial mode
INFO  : Completed executing command(queryId=hive_20221021153851_c8fdd21b-efc1-424b-81fb-ded932275f86); Time taken: 0.153 seconds
INFO  : OK
No rows affected (0.301 seconds)
0: jdbc:hive2://devm01.devops.com:2181,devm02>
```

寫入 internal_table

```
INFO  : TaskCounter_Map_1_INPUT__dummy_table:
INFO  :     INPUT_RECORDS_PROCESSED: 4
INFO  :     INPUT_SPLIT_LENGTH_BYTES: 1
INFO  : TaskCounter_Map_1_OUTPUT_Reducer_2:
INFO  :     ADDITIONAL_SPILLS_BYTES_READ: 0
INFO  :     ADDITIONAL_SPILLS_BYTES_WRITTEN: 0
INFO  :     ADDITIONAL_SPILL_COUNT: 0
INFO  :     DATA_BYTES_VIA_EVENT: 0
INFO  :     OUTPUT_BYTES: 140
INFO  :     OUTPUT_BYTES_PHYSICAL: 189
INFO  :     OUTPUT_BYTES_WITH_OVERHEAD: 149
INFO  :     OUTPUT_LARGE_RECORDS: 0
INFO  :     OUTPUT_RECORDS: 1
INFO  :     SPILLED_RECORDS: 0
INFO  : TaskCounter_Reducer_2_INPUT_Map_1:
INFO  :     FIRST_EVENT_RECEIVED: 54
INFO  :     INPUT_RECORDS_PROCESSED: 1
INFO  :     LAST_EVENT_RECEIVED: 54
INFO  :     NUM_FAILED_SHUFFLE_INPUTS: 0
INFO  :     NUM_SHUFFLED_INPUTS: 1
INFO  :     SHUFFLE_BYTES: 165
INFO  :     SHUFFLE_BYTES_DECOMPRESSED: 149
INFO  :     SHUFFLE_BYTES_DISK_DIRECT: 165
INFO  :     SHUFFLE_BYTES_TO_DISK: 0
INFO  :     SHUFFLE_BYTES_TO_MEM: 0
INFO  :     SHUFFLE_PHASE_TIME: 112
----------------------------------------------------------------------------------------------
        VERTICES      MODE        STATUS   TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------------------
Map 1 .......... container   SUCCEEDED      1        1         0        0        0       0
Reducer 2 ...... container   SUCCEEDED      1        1         0        0        0       0  anaged/hive/users_internal
----------------------------------------------------------------------------------------------
VERTICES: 02/02  [==========================>>] 100%  ELAPSED TIME: 10.50 s
----------------------------------------------------------------------------------------------
12 rows affected (20.696 seconds)
0: jdbc:hive2://devm01.devops.com:2181,devm02> █
```

查詢 external table

```
0: jdbc:hive2://devm01.devops.com:2181,devm02> select * from default.users;
INFO  : Compiling command(queryId=hive_20221021154402_b70cfa33-5302-455e-86b4-485941a9ffed): select * from default.users
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Created Hive schema: Schema(fieldSchemas:[FieldSchema(name:users.name, type:string, comment:null), FieldSchema(name:users.age, type:int
, comment:null)], properties:null)
INFO  : Completed compiling command(queryId=hive_20221021154402_b70cfa33-5302-455e-86b4-485941a9ffed); Time taken: 0.078 seconds
INFO  : Executing command(queryId=hive_20221021154402_b70cfa33-5302-455e-86b4-485941a9ffed): select * from default.users
INFO  : Completed executing command(queryId=hive_20221021154402_b70cfa33-5302-455e-86b4-485941a9ffed); Time taken: 0.006 seconds
INFO  : OK
+--------------+-------------+
| users.name   | users.age   |
+--------------+-------------+
| Alice        | 16          |
| Bob          | 18          |
| Cindy        | 53          |
| Joe          | 28          |
| Apple        | 5           |
| John         | 33          |
| Roger        | 45          |
| Kevin        | 17          |
| Gordon       | 30          |
| Paul         | 33          |
| Sandy        | 39          |
| Alex         | 70          |
+--------------+-------------+
12 rows selected (0.147 seconds)
0: jdbc:hive2://devm01.devops.com:2181,devm02> █
```

查詢 internal table

```
0: jdbc:hive2://devm01.devops.com:2181,devm02> select * from default.users_internal;
INFO  : Compiling command(queryId=hive_20221021154223_f5aa69ae-7af4-4478-8b37-30b3e4f60d82): select * from default.users_internal
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Created Hive schema: Schema(fieldSchemas:[FieldSchema(name:users_internal.name, type:string, comment:null), FieldSchema(name:users_inte
rnal.age, type:int, comment:null)], properties:null)
INFO  : Completed compiling command(queryId=hive_20221021154223_f5aa69ae-7af4-4478-8b37-30b3e4f60d82); Time taken: 0.143 seconds
INFO  : Executing command(queryId=hive_20221021154223_f5aa69ae-7af4-4478-8b37-30b3e4f60d82): select * from default.users_internal
INFO  : Completed executing command(queryId=hive_20221021154223_f5aa69ae-7af4-4478-8b37-30b3e4f60d82); Time taken: 0.007 seconds
INFO  : OK
+----------------------+---------------------+
| users_internal.name  | users_internal.age  |
+----------------------+---------------------+
| Alice                | 16                  |
| Bob                  | 18                  |
| Cindy                | 53                  |
| Joe                  | 28                  |
| Apple                | 5                   |
| John                 | 33                  |
| Roger                | 45                  |
| Kevin                | 17                  |
| Gordon               | 30                  |
| Paul                 | 33                  |
| Sandy                | 39                  |
| Alex                 | 70                  |
+----------------------+---------------------+
12 rows selected (0.419 seconds)
0: jdbc:hive2://devm01.devops.com:2181,devm02>
```

- JDBC

使用 jdbc-kerberos-0.0.1-SNAPSHOT.jar 進行測試

- 先確認 application.properties 使用 hive 的設定 (圖片僅供參考)

```
#spring.datasource.url=jdbc:impala://devw01.devops.com:21050/default;AuthMech=1;KrbHostFQDN=devw01.devops.com;KrbServiceName=impala;
spring.datasource.url=jdbc:hive2://devm03.devops.com:10000;AuthMech=1;KrbHostFQDN=devm03.devops.com;KrbServiceName=hive;
#spring.datasource.driver-class-name=com.cloudera.impala.jdbc.Driver
spring.datasource.driver-class-name=com.cloudera.hive.jdbc.HS2Driver
spring.jpa.database-platform=org.hibernate.dialect.SQLServerDialect
```

- 啟動 web 程式 (圖片啟動指令僅供參考)

```
java -jar jdbc-kerberos-0.0.1-SNAPSHOT.jar
```

- 另開 Terminal，透過 CURL 指令呼叫 (圖片啟動指令僅供參考)

```
curl --location --request POST 'http://<ip>:8080/api/query' \
--data-raw ''
```



## ▼ Impala

- Impala-shell (圖片僅供參考)

```
impala-shell -i <host_name>

select * from default.users;
select * from default.users_internal;
```

查詢 external table



- JDBC

使用 jdbc-kerberos-0.0.1-SNAPSHOT.jar 進行測試

- 先確認 application.properties 使用 impala 的設定 (圖片僅供參考)



- 啟動 web 程式 (圖片啟動指令僅供參考, log 過多提供部分截圖)

```
java -jar jdbc-kerberos-0.0.1-SNAPSHOT.jar
```

- 另開 Terminal，透過 CURL 指令呼叫 (圖片啟動指令僅供參考)

```
curl --location --request POST 'http://<ip>:8080/api/query' \
--data-raw ''
```



## ▼ Spark

▼ Scala

- 執行計算 pi (jar 檔路徑僅供參考，可能會因為版本有差異)

```
spark-submit --class org.apache.spark.examples.SparkPi \
  --master yarn-client \
  /opt/cloudera/parcels/CDH/lib/spark/examples/jars/spark-examples_2.11-2.4.7.7.1.7.0-551.jar 1000
```

- 使用 HWC (指令僅供參考，可能會因為版本有差異)

```
spark-shell \
--master yarn \
--conf spark.sql.extensions="com.hortonworks.spark.sql.rule.Extensions" \
--conf spark.kryo.registrator=com.qubole.spark.hiveacid.util.HiveAcidKyroRegistrator \
--conf spark.sql.hive.hiveserver2.jdbc.url="jdbc:hive2://devm01.devops.com:2181,devm02.devops.com:2181,devm03.devops.com:2181/defau
--conf spark.sql.hive.hiveserver2.jdbc.url.principal="hive/_HOST@AMDEV.COM" \
--conf spark.datasource.hive.warehouse.read.mode=DIRECT_READER_V2
```

- 使用 HWC 查詢 Internal table 資料

```
import com.hortonworks.hwc.HiveWarehouseSession
import com.hortonworks.hwc.HiveWarehouseSession._
val hive = HiveWarehouseSession.session(spark).build()
hive.sql("select * from default.users_internal").show
```

```
scala> import com.hortonworks.hwc.HiveWarehouseSession
import com.hortonworks.hwc.HiveWarehouseSession

scala> import com.hortonworks.hwc.HiveWarehouseSession._
import com.hortonworks.hwc.HiveWarehouseSession._

scala> val hive = HiveWarehouseSession.session(spark).build()
hive: com.hortonworks.spark.sql.hive.llap.HiveWarehouseSessionImpl = com.hortonworks.spark.sql.hive.llap.HiveWarehouseSessionImpl@22ab40c1

scala> hive.sql("select * from default.users_internal").show
22/10/21 16:08:18 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:08:19 WARN llap.HiveWarehouseSessionImpl:  com.qubole.spark.hiveacid.HiveAcidAutoConvertExtension will be replaced by com.hortonwor
ks.spark.sql.rule.Extensions.Please switch to com.hortonworks.spark.sql.rule.Extensions.
22/10/21 16:08:20 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
Hive Session ID = 18449d1b-9501-43a4-80b4-6ce130931b9f
22/10/21 16:08:22 INFO rule.HWCSwitchRule: using DIRECT_READER_V2 extension for reading
22/10/21 16:08:22 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:08:24 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:08:24 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
+-----+---+
| name|age|
+-----+---+
|Alice| 16|
|  Bob| 18|
|Cindy| 53|
|  Joe| 28|
|Apple|  5|
| John| 33|
|Roger| 45|
|Kevin| 17|
|Gordon| 30|
| Paul| 33|
|Sandy| 39|
| Alex| 70|
+-----+---+

scala>
```

- 讀取資料 (HDFS)

```
val df = spark.read.format("csv").load("/tmp/hello.txt")
df.show
```

```
scala> val df = spark.read.format("csv").load("/tmp/hello.txt")
df: org.apache.spark.sql.DataFrame = [_c0: string]

scala> df.show
+-----------+
|        _c0|
+-----------+
|hello world|
+-----------+
```

- 讀取資料 (External table)

```
spark.sql("select * from default.users").show
```



- 寫入資料至 HDFS

```
val df = spark.sql("select * from default.users")
df.write.format("parquet").mode("overwrite").save("/tmp/user")
```



▼ PySpark
- 使用 HWC (指令僅供參考，可能會因為版本有差異)

```
pyspark \
--jars /opt/cloudera/parcels/CDH-7.1.7-1.cdh7.1.7.p0.15945976/jars/hive-warehouse-connector-assembly-1.0.0.7.1.7.0-551.jar \
--py-files /opt/cloudera/parcels/CDH/lib/hive_warehouse_connector/pyspark_hwc-1.0.0.7.1.7.0-551.zip \
--conf spark.sql.extensions="com.hortonworks.spark.sql.rule.Extensions" \
--conf spark.kryo.registrator="com.qubole.spark.hiveacid.util.HiveAcidKyroRegistrator" \
--conf spark.sql.hive.hiveserver2.jdbc.url="jdbc:hive2://devm01.devops.com:2181,devm02.devops.com:2181,devm03.devops.com:2181/defau
--conf spark.sql.hive.hiveserver2.jdbc.url.principal="hive/_HOST@AMDEV.COM" \
--conf spark.security.credentials.hiveserver2.enabled=true \
--conf spark.datasource.hive.warehouse.read.mode=DIRECT_READER_V2
```

```
[leslie@devm01 ~]$ pyspark \
> --jars /opt/cloudera/parcels/CDH-7.1.7-1.cdh7.1.7.p0.15945976/jars/hive-warehouse-connector-assembly-1.0.0.7.1.7.0-551.jar \
> --py-files /opt/cloudera/parcels/CDH/lib/hive_warehouse_connector/pyspark_hwc-1.0.0.7.1.7.0-551.zip \
> --conf spark.sql.extensions="com.hortonworks.spark.sql.rule.Extensions" \
> --conf spark.kryo.registrator="com.qubole.spark.hiveacid.util.HiveAcidKyroRegistrator" \
> --conf spark.sql.hive.hiveserver2.jdbc.url="jdbc:hive2://devm01.devops.com:2181,devm02.devops.com:2181,devm03.devops.com:2181/default;service
DiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2" \
> --conf spark.sql.hive.hiveserver2.jdbc.url.principal="hive/_HOST@AMDEV.COM" \
> --conf spark.security.credentials.hiveserver2.enabled=true \
> --conf spark.datasource.hive.warehouse.read.mode=DIRECT_READER_V2
Python 2.7.5 (default, Oct 14 2020, 14:45:30)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-44)] on linux2
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/10/21 16:20:21 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:20:24 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:20:28 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:20:28 WARN ipc.Client: Exception encountered while connecting to the server : org.apache.hadoop.ipc.RemoteException(org.apache.hado
op.ipc.StandbyException): Operation category READ is not supported in state standby. Visit https://s.apache.org/sbnn-error
22/10/21 16:20:36 WARN cluster.YarnSchedulerBackend$YarnSchedulerEndpoint: Attempted to request executors before the AM has registered!
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 2.4.7.7.1.7.0-551
      /_/

Using Python version 2.7.5 (default, Oct 14 2020 14:45:30)
SparkSession available as 'spark'.
>>>
```

- 使用 HWC 查詢 Internal table 資料

```
from pyspark_llap import HiveWarehouseSession
from pyspark.sql.functions import col
hive = HiveWarehouseSession.session(spark).build()
hive.sql("select * from default.users_internal").show()
```

```
>>> from pyspark_llap import HiveWarehouseSession
>>> from pyspark.sql.functions import col
>>> hive = HiveWarehouseSession.session(spark).build()
>>> hive.sql("select * from default.users_internal").show()
22/10/21 16:24:36 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:24:37 WARN llap.HiveWarehouseSessionImpl:  com.qubole.spark.hiveacid.HiveAcidAutoConvertExtension will be replaced by com.hortonwor
ks.spark.sql.rule.Extensions.Please switch to com.hortonworks.spark.sql.rule.Extensions.
22/10/21 16:24:38 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
Hive Session ID = d7d91867-ee9d-4cda-8d49-63b614702bc2
22/10/21 16:24:40 INFO rule.HWCSwitchRule: using DIRECT_READER_V2 extension for reading
22/10/21 16:24:40 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:24:42 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
22/10/21 16:24:42 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
+------+---+
| name|age|
+------+---+
| Alice| 16|
|   Bob| 18|
| Cindy| 53|
|   Joe| 28|
| Apple|  5|
|  John| 33|
| Roger| 45|
| Kevin| 17|
|Gordon| 30|
|  Paul| 33|
| Sandy| 39|
|  Alex| 70|
+------+---+
```

- 讀取資料 (HDFS)

```
df = spark.read.format('csv').load('/tmp/hello.txt')
df.show()
```

```
>>> df = spark.read.format('csv').load('/tmp/hello.txt')

>>> df.show()
+-----------+
|        _c0|
+-----------+
|hello world|
+-----------+

>>>
```

- 讀取資料 (External table)

```
spark.sql("select * from default.users").show()
```

```
>>> spark.sql("select * from default.users").show()
+------+---+
|  name|age|
+------+---+
| Alice| 16|
|   Bob| 18|
| Cindy| 53|
|   Joe| 28|
| Apple|  5|
|  John| 33|
| Roger| 45|
| Kevin| 17|
|Gordon| 30|
|  Paul| 33|
| Sandy| 39|
|  Alex| 70|
| Alice| 16|
|   Bob| 18|
| Cindy| 53|
|   Joe| 28|
| Apple|  5|
|  John| 33|
| Roger| 45|
| Kevin| 17|
+------+---+
only showing top 20 rows
```

- 寫入資料至 HDFS

```
df = spark.sql('select * from default.users')
df.write.format('parquet').mode('overwrite').save('/tmp/user')
```

```
>>> df = spark.sql('select * from default.users')
>>> df.write.format('parquet').mode('overwrite').save('/tmp/user')
>>> 22/10/21 16:34:49 WARN conf.HiveConf: HiveConf of name hive.masking.algo does not exist
```

## ▼ Python 3

```
python3 --version
```

### ▼ Java 8

```
java -version
javac -version
```

## Executor List

1. bin/application.properties

2. bin/jdbc-kerberos-0.0.1-SNAPSHOT.jar

3. README.md

## Reference

1. jdbc-doc/Cloudera-JDBC-Driver-for-Apache-Hive-Install-Guide.pdf

2. jdbc-doc/Cloudera-JDBC-Driver-for-Impala-Install-Guide.pdf

3. source_code (jdbc-kerberos)