

# Hortonworks Tutorial 6

## Optimization on Pig

In this tutorial, we'll focus on how to optimize and improve the performance of Pig on **TEZ**.

### Task 1. Dataset

We will use the same tables used in Tutorial 5: `basketball_master.csv`, `basketball_players.csv` and `basketball_teams.csv`. Upload these files to HDFS, once checking their existence.

Recall that **TEZ** aims to optimise the performance and speed. By default, Pig uses the MapReduce framework, but we can configure Pig to use the **TEZ**.

### Task 2. Speed Improvements in Pig

1. Create a new Pig script named "t5". Then copy the pig script below and paste it in editor:

```
players_raw = LOAD '/tmp/FIT5148/basketball_players.csv'
USING PigStorage(',');

players = FILTER players_raw BY $1 > 0;

player_points = FOREACH players GENERATE $0
AS playerID, $1 AS year, $8 AS points;

grp_data = GROUP player_points BY (year);

max_points = FOREACH grp_data GENERATE group as grp,
MAX(player_points.points) AS max_points;

join_max_points = JOIN max_points
BY ($0, max_points), player_points BY (year, points);

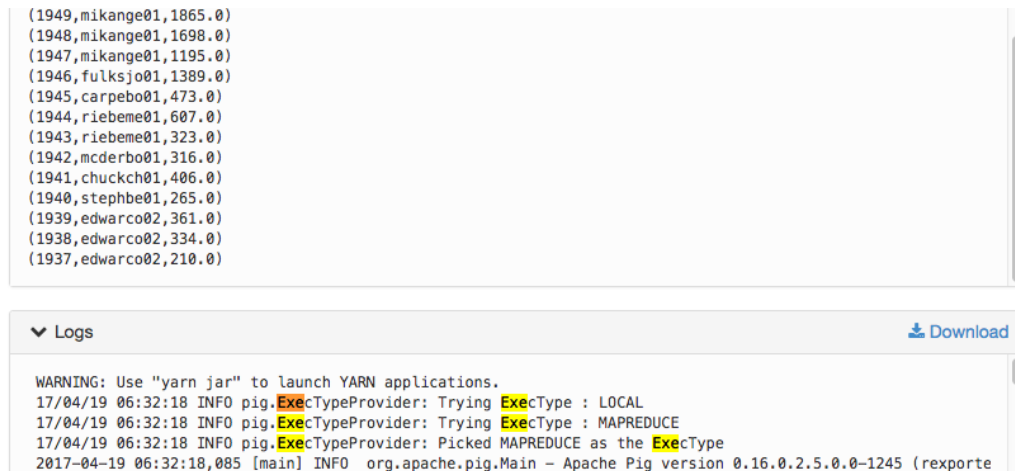
join_data = FOREACH join_max_points GENERATE $0
AS year, $2 AS playerID, $1 AS points;

data_order = ORDER join_data BY year desc;

DUMP data_order;
```

2. Execute the pig script using the MapReduce (the default). It'll take few minutes to finish the job. After the job is completed, examine the "ExecType" in Logs where the default mode is mentioned as "MAPREDUCE":
3. This time, we run the same Pig script using "TEZ". Check the "Execute on Tez" box, and Execute it. Observe the response time (much faster or slower?). Also, check the difference in "Logs" (see Fig. 2)
4. You can always check the executed scripts from the "History" tab and the time taken to execute/complete them.

### Task 3. Using the Command Line

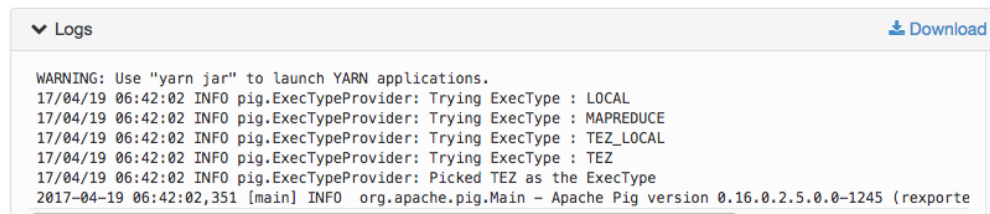


```
(1949,mikange01,1865.0)
(1948,mikange01,1698.0)
(1947,mikange01,1195.0)
(1946,fulksjo01,1389.0)
(1945,carpebo01,473.0)
(1944,riebeme01,607.0)
(1943,riebeme01,323.0)
(1942,mcderbo01,316.0)
(1941,chuckch01,406.0)
(1940,stephbe01,265.0)
(1939,edwarco02,361.0)
(1938,edwarco02,334.0)
(1937,edwarco02,210.0)
```

▼ Logs [Download](#)

```
WARNING: Use "yarn jar" to launch YARN applications.
17/04/19 06:32:18 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
17/04/19 06:32:18 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
17/04/19 06:32:18 INFO pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2017-04-19 06:32:18,085 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0.2.5.0.0-1245 (rexporte
```

Figure 1: The default ExecType of a Pig query



▼ Logs [Download](#)

```
WARNING: Use "yarn jar" to launch YARN applications.
17/04/19 06:42:02 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
17/04/19 06:42:02 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
17/04/19 06:42:02 INFO pig.ExecTypeProvider: Trying ExecType : TEZ_LOCAL
17/04/19 06:42:02 INFO pig.ExecTypeProvider: Trying ExecType : TEZ
17/04/19 06:42:02 INFO pig.ExecTypeProvider: Picked TEZ as the ExecType
2017-04-19 06:42:02,351 [main] INFO org.apache.pig.Main - Apache Pig version 0.16.0.2.5.0.0-1245 (rexporte
```

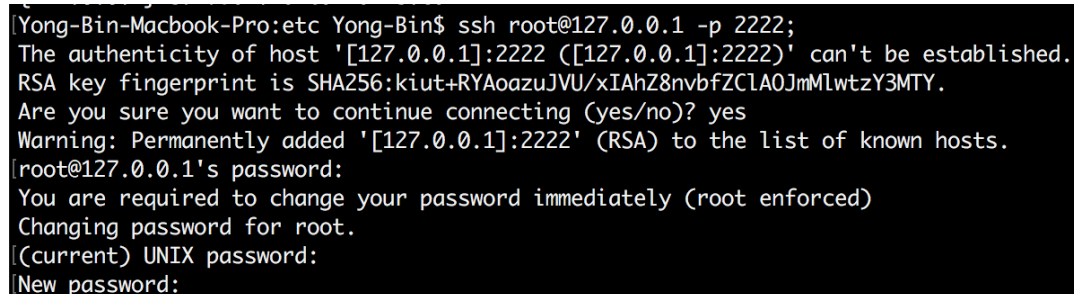
Figure 2: The ExecType: Tez

To use "Tez" with "Pig" we have to use "SSH" via the VM instead of **Ambari**. Open your terminal (mac and linux) or puTTY (windows). Type the following command to access the Sandbox through SSH:

```
ssh root@127.0.0.1 -p 2222;
```

If you use the Windows Command, run "putty -ssh root@127.0.0.1 2222" on it. If you access this URL for the first time, there is a prompt box, "PuTTY Security Alert". Read it and click "Yes".

If you login to the sandbox, you will see the following message (see Fig. 3): **"Are you sure you want to continue connecting (yes/no)?"** Enter "yes" to proceed. You will be immediately asked to change the password. Enter current password as "hadoop", and change the password. Make sure you remember the new password because you will be using it the next time you login.



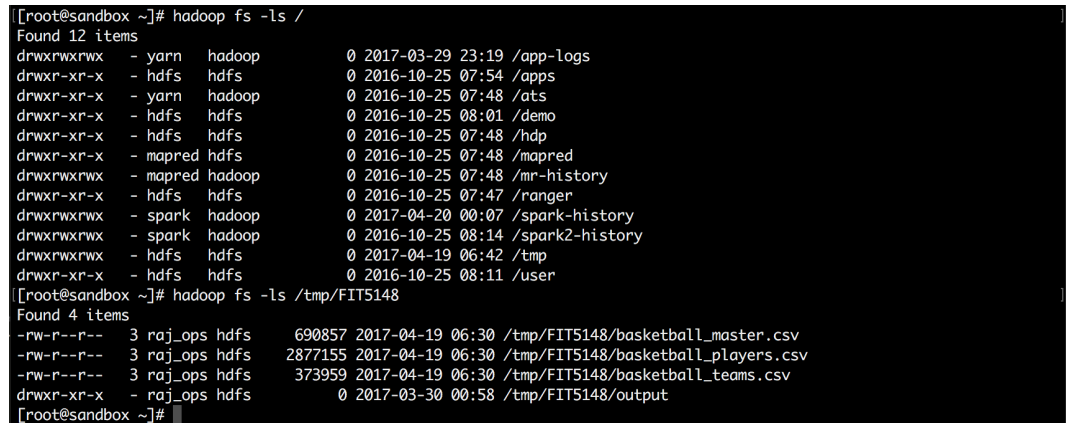
```
Yong-Bin-Macbook-Pro:etc Yong-Bin$ ssh root@127.0.0.1 -p 2222;
The authenticity of host '[127.0.0.1]:2222 ([127.0.0.1]:2222)' can't be established.
RSA key fingerprint is SHA256:kiut+RYAoazuJVU/xIAhZ8nbfZC1A0JmMltzY3MTY.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added '[127.0.0.1]:2222' (RSA) to the list of known hosts.
root@127.0.0.1's password:
You are required to change your password immediately (root enforced)
Changing password for root.
(current) UNIX password:
New password:
```

Figure 3: The ExecType: Tez

Let's check current files available on HDFS with the commands:

```
hadoop fs-ls/
hadoop fs-ls/tmp/FIT5148
```

The results will be as Fig. 4



```
[root@sandbox ~]# hadoop fs -ls /
Found 12 items
drwxrwxrwx - yarn   hadoop      0 2017-03-29 23:19 /app-logs
drwxr-xr-x - hdfs   hdfs      0 2016-10-25 07:54 /apps
drwxr-xr-x - yarn   hadoop      0 2016-10-25 07:48 /ats
drwxr-xr-x - hdfs   hdfs      0 2016-10-25 08:01 /demo
drwxr-xr-x - hdfs   hdfs      0 2016-10-25 07:48 /hdp
drwxr-xr-x - mapred hdfs      0 2016-10-25 07:48 /mapred
drwxrwxrwx - mapred hadoop      0 2016-10-25 07:48 /mr-history
drwxr-xr-x - hdfs   hdfs      0 2016-10-25 07:47 /ranger
drwxrwxrwx - spark  hadoop      0 2017-04-20 00:07 /spark-history
drwxrwxrwx - spark  hadoop      0 2016-10-25 08:14 /spark2-history
drwxrwxrwx - hdfs   hdfs      0 2017-04-19 06:42 /tmp
drwxr-xr-x - hdfs   hdfs      0 2016-10-25 08:11 /user
[root@sandbox ~]# hadoop fs -ls /tmp/FIT5148
Found 4 items
-rw-r--r--  3 raj_ops hdfs      690857 2017-04-19 06:30 /tmp/FIT5148/basketball_master.csv
-rw-r--r--  3 raj_ops hdfs    2877155 2017-04-19 06:30 /tmp/FIT5148/basketball_players.csv
-rw-r--r--  3 raj_ops hdfs    373959 2017-04-19 06:30 /tmp/FIT5148/basketball_teams.csv
drwxr-xr-x  - raj_ops hdfs      0 2017-03-30 00:58 /tmp/FIT5148/output
[root@sandbox ~]#
```

Figure 4: Sandbox data files

First we will run Pig without using Tez. You need to save the pig script below in Word Notepad or Mac TextEdit with the name "script1.pig":

```
players_raw = LOAD '/tmp/FIT5148/basketball_players.csv'
USING PigStorage(',');
players = FILTER players_raw BY $1 > 0;
player_points = FOREACH players
GENERATE $0 AS playerID, $1 AS year, $8 AS points;
grp_data = GROUP player_points BY (year);
max_points = FOREACH grp_data
GENERATE group as grp, MAX(player_points.points) AS max_points;
join_max_points = JOIN max_points
BY ($0, max_points), player_points BY (year, points);
join_data = FOREACH join_max_points
GENERATE $0 AS year, $2 AS playerID, $1 AS points;
data_order = ORDER join_data BY year desc;
DUMP data_order;
```

After creating the "script1.pig" file, copy and paste this file to your Dropbox (or any cloud based storage space that you have such that it can be downloaded later using its URL). If you are using Dropbox, make sure the URL for the file ends with ?dl=1 Then type this command in SSH:

```
wget https://www.dropbox.com/YOURFOLDERPATH/script1.pig?dl=1-Oscript1.pig
```

If you put the file in your OneDrive folder, go to the web interface of OneDrive (e.g. <https://onedrive.live.com/>). Right-click on the file and choose "Embed". Then, click the "Generate" button on the right-bottom frame of the page. After that, copy the url tagged by "src" in the "<iframe>". Make changes in the URL by replacing 'embed' with 'download'. After making changes, you URL will be look like:

```
https://onedrive.live.com/download?cid=2C4B95628EC34191&resid=2C4B95628EC34191%21112&authkey=AGGrdQ8QWjKvzMk.
```

Note that the "cid", "resid" and "authkey" should be differently shown on your page.

Now, You have done. The above modified link is the direct link of your file.

Then, type the command in SSH:

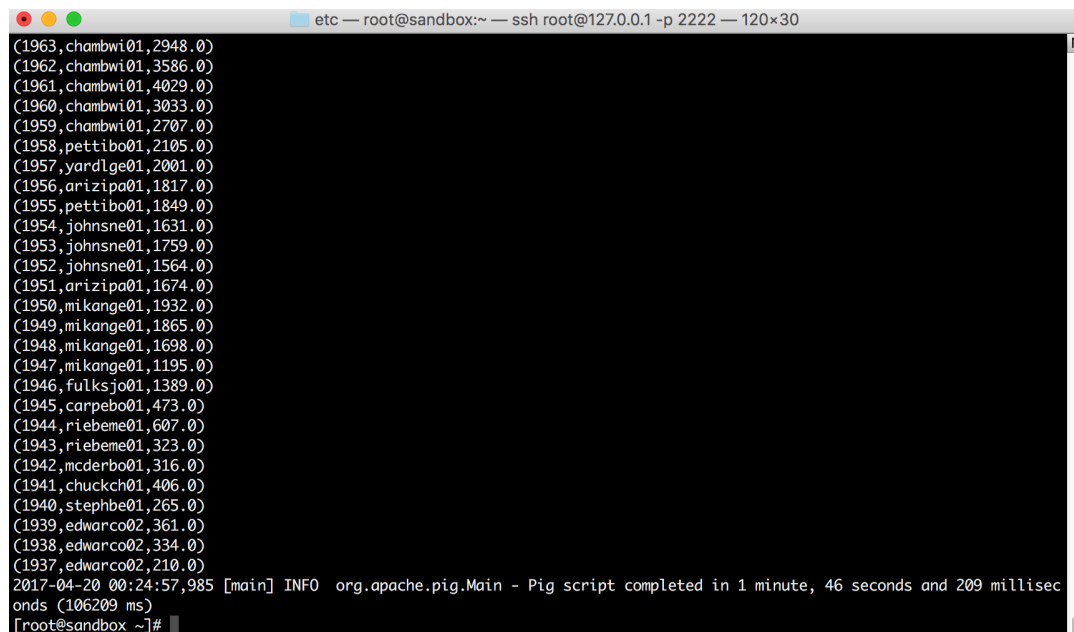
```
wget fileURL -O script1.pig
```

The "fileURL" must be replaced with the actual data URL that you can download.

Execute the Pig scripts in the file with the MapReduce (the default) in the command line by typing the following command in SSH:

```
pig script1.pig
```

You will see the result of the job with the time taken (see Fig. 5).



```
(1963,chambwi01,2948.0)
(1962,chambwi01,3586.0)
(1961,chambwi01,4029.0)
(1960,chambwi01,3033.0)
(1959,chambwi01,2707.0)
(1958,pettibo01,2105.0)
(1957,yardlge01,2001.0)
(1956,arizipa01,1817.0)
(1955,pettibo01,1849.0)
(1954,johnsne01,1631.0)
(1953,johnsne01,1759.0)
(1952,johnsne01,1564.0)
(1951,arizipa01,1674.0)
(1950,mikange01,1932.0)
(1949,mikange01,1865.0)
(1948,mikange01,1698.0)
(1947,mikange01,1195.0)
(1946,fulksjo01,1389.0)
(1945,carpebo01,473.0)
(1944,riebeme01,607.0)
(1943,riebeme01,323.0)
(1942,mcderbo01,316.0)
(1941,chuckch01,406.0)
(1940,stephbe01,265.0)
(1939,edwarco02,361.0)
(1938,edwarco02,334.0)
(1937,edwarco02,210.0)
2017-04-20 00:24:57,985 [main] INFO org.apache.pig.Main - Pig script completed in 1 minute, 46 seconds and 209 milliseconds (106209 ms)
[root@sandbox ~]#
```

Figure 5: The result of the job, script1.pig

#### Task 4. Exercise: Optimize a Pig script on Tez

Now, we try to perform optimization and run the same Pig scripts with Tez. Take a look at this URL to know how to optimize a Pig script, and then execute "script1.pig" on Tez: [https://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.4.2/bk\\_dataintegration/content/ch\\_running-pig-tez.html](https://docs.hortonworks.com/HDPDocuments/HDP2/HDP-2.4.2/bk_dataintegration/content/ch_running-pig-tez.html)

The completion time should be reduced when using Tez compared to MapReduce (Note: your time can be slightly different from ours but running Pig with Tez should show some speed improvement). Find the execution time in the log in SSH and compare it with the time when the script was executed with MapReduce.