Cài đặt Apache Pig

1. Download Apache Pig

wget https://dlcdn.apache.org/pig/pig-0.17.0/pig-0.17.0.tar.gz

2. Giải nén và đổi tên thư mục pig\$ tar -xzf pig-0.17.0.tar.gz

```
$ mv pig-0.17.0 pig
```

```
hadoopminhchau@minhchau-master:~$ ls

apache-hive-3.1.3-bin.tar.gz hadoop-3.3.2.tar.gz output_dir sample.txt units.jar

db-derby-10.15.2.0-bin.tar.gz hadoop-core-1.2.1.jar pig test.sh

hadoop-core-1.2.1.jar.zip

hadoop hive ProcessUnits.java

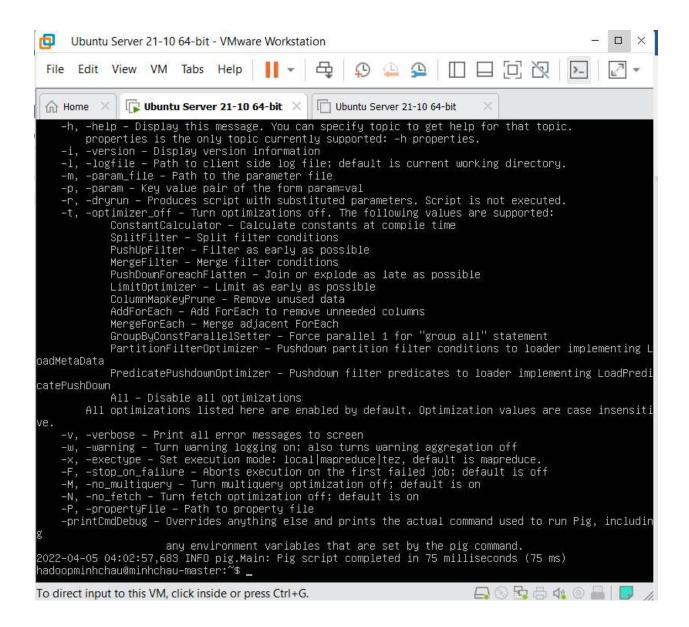
hadoopminhchau@minhchau-master:~$
```

3. Khai báo biến môi trường

- \$ export PIG HOME=/home/hadoopminhchau/pig
- \$ export PATH=\$PATH:\$PIG_HOME/bin
- \$ export PIG CLASSPATH=\$HADOOP CONF DIR

4. Check Pig

\$ pig -help



5. Thuc thi Pig

3. Thực thi Tig						
	Local	Tez Local	Spark	MapReduce	Tez Mode	Spark
	Mode	Mode	Local	Mode		Mode
			Mode			
Interactive	Yes	Experimental	Yes	Yes		
Mode						
Batch	Yes	Experimental	Yes	Yes		
Mode						

Pig có 06 chế độ thực thi

Local Mode

Chạy trên máy đơn. Tất cả các tập tin được cài đặt và chạy trên hệ thống file và máy chủ
cục bộ.

```
$ pig -x local
```

Tez Local Mode

• Tương tự chế độ cục bộ, ngoại trừ bên trong Pig sẽ gọi công cụ thực thi Tez.

```
$ pig -x tez local
```

• Spark Local Mode

• Tương tự chế độ cục bộ, ngoại trừ bên trong Pig sẽ gọi công cụ thực thi Spark.

```
$ pig -x spark local
```

• Map Reduce Mode

 Chạy ở chế độ mapreduce, cần truy cập cụm Hadoop và HDFS. Đây là chế độ chạy mặc định, có thể không cần cung cấp thông số -x

```
$ pig
```

Hoăc

\$ pig -x mapreduce

• Tez Mode

• Cần truy cập cụm Hadoop và hệ thống HDFS

```
$ pig -x tez
```

• Spark Mode

 Cần truy cập Spark, cụm Yarn hoặc Mesos và hệ thống file HDFS. Pig Script chạy trên Spark có thể tận dụng tính năng cấp phát động (dynamic allocation). Tính năng này có thể được kích hoạt bằng cách bật spark.dynamicAllocation.enabled.

```
$ pig -x spark
```

6. Interactive Mode

 Có thể chạy Pig ở chế độ Interactive sử dụng Grunt Shell. Gọi Grunt Shell sử dụng lệnh Pig và nhập các câu lệnh Pig Latin và thực thi.

Ví dụ sau lấy tất cả các ID người dùng từ file /etc/passwd

Copy file passwd ra thư mục làm việc cục bộ

\$ cp /etc/passwd /home/hadoopminhchau

```
nadoopminhchau@minhchau–master:~$ cp /etc/passwd /home/hadoopminhchau/
nadoopminhchau@minhchau–master:~$ ls
                                                                                                                    sample.txt
                                                                                                                    test.sh
                                                                                  pig_1650382433801.log
                                                                                  pig_1650382485427.log
                                            passwd
                                                                                   ProcessUnits.java
                                            ~$
nadoopminhchau@minhchau−master:
Khởi đông Grunt Shell (trong cục bộ hoặc hadoop) và nhập các câu lệnh Pig trực tiếp
$ grunt> A = load '/home/hadoopminhchau/passwd' using
PigStorage(':');
$ grunt> B = foreach A generate $0 as id;
$ grunt> dump B;
hadoopminhchau@minhchau–master:~$ pig –x local
2022–04–19 16:26:00,596 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
2022–04–19 16:26:00,596 INFO pig.ExecTypeProvider: Picked LOCAL as the ExecType
2022–04–19 16:26:00,638 [main] INFO org.apache.pig.Main – Apache Pig version 0.17.0 (r1797386) o
iled Jun 02 2017, 15:41:58
2022–04–19 16:26:00,638 [main] INFO org.apache.pig.Main – Logging error messages to: /home/hadoo
nhchau/pig_1650385560636.log
2022–04–19 16:26:00,655 [main] INFO org.apache.pig.impl.util.Utils – Default bootup file /home/h
opminhchau/.pigbootup not found
2022–04–19 16:26:00,735 [main] INFO org.apache.hadoop.conf.Configuration.deprecation – mapred.jo
racker is deprecated. Instead, use mapreduce.jobtracker.address
2022–04–19 16:26:00,737 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngi
– Connecting to hadoop file system at: file:///
2022–04–19 16:26:00,808 [main] INFO org.apache.hadoop.conf.Configuration.deprecation – io.bytes.
.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2022–04–19 16:26:00,822 [main] INFO org.apache.pig.PigServer – Pig Script ID for the session: PI
efault-97cd2578-2b38-46e3-91b3-f03b9abf487f
2022–04–19 16:26:00,822 [main] WARN org.apache.pig.PigServer – ATS is disabled since yarn.timeli
grunt> A = load '/home/hadoopminhchau/passwd' using PigStorage(':');
grunt> A = load '/home/hadoopminhchau/passwd' using PigStorage(':');
2022–04–19 16:26:17,649 [main] INFO org.apache.hadoop.conf.Configuration.deprecation – io.bytes.
.checksum is deprecated. Instead, use dfs.bytes-per-checksum
```

grunt> B = foreach A generate \$0 as id;

grunt> dump B;

```
Total input paths to process : 1
(root)
(daemon)
(bin)
(sys)
(sync)
(games)
(man)
(lp)
(mail)
(news)
(uucp)
(proxy)
(www–data)
(backup)
(list)
(irc)
(gnats)
(nobody)
(_apt)
(systemd–timesync)
(systemd-network)
(systemd-resolve)
(messagebus)
(pollinate)
(sshd)
(syslog)
(uuidd)
(tcpdump)
(landscape)
(usbmux)
(systemd-coredump)
(minhchau)
(1xd)
(hadoopminhchau)
grunt>
```

7. Batch Mode

Có thể chạy Pig ở chế độ batch mode sử dụng Pig Script và chạy pig ở chế độ cục bộ hoặc Hadoop

Ví dụ: chạy lại ví dụ trên ở chế độ Batch Mode

Tạo file Pig Script id.pig

\$ vim id.pig

hadoopminhchau@minhchau–master:~\$ vim id.pig

```
Thực thi file script
```

```
$ pig -x local id.pig
```

```
Success!
Job Stats (time in seconds):
JobId
       Maps Reduces MaxMapTime
                                        MinMapTime
                                                         AvgMapTime
                                                                         MedianMapTime
                                                                                          MaxReduce
        MinReduceTime AvgReduceTime
                                        MedianReducetime
                                                                 Alias
                                                                         Feature Outputs
job_local1540184028_0001
                                                                         n/a
                                                 n/a
                                                         n/a
                                                                 n/a
       MAP_ONLY
                        file:///home/hadoopminhchau/id.out,
A,B
Input(s):
Successfully read 35 records from: "file:///home/hadoopminhchau/passwd"
Successfully stored 35 records in: "file:///home/hadoopminhchau/id.out"
Counters:
Total records written : 35
Total bytes written : O
Spillable Memory Manager spill count : O
Total bags proactively spilled: O
Total records proactively spilled: O
Job DAG:
job_local1540184028_0001
                                     org.apache.hadoop.metrics2.impl.MetricsSystemImpl – JobTrack
2022–04–19 16:38:24,852 [main] WARN
metrics system already initialized!
2022-04-19 16:38:24,855 [main] WARN
                                     org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTrack
metrics system already initialized!
2022–04–19 16:38:24,857 [main] WARN
metrics system already initialized!
                                     org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTrack
2022–04–19 16:38:24,870 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer
pReduceLauncher – Success!
2022–04–19 16:38:24,888 [main] INFO org.apache.pig.Main – Pig script completed in 2 seconds and
milliseconds (2258 ms)
hadoopminhchau@minhchau–master:~$ _
```

Xuất kết quả

```
hadoopminhchau@minhchau-master:~$ ls

apache-hive-3.1.3-bin.tar.gz hadoop-core-1.2.1.jar.zip pig test.sh

db-derby-10.15.2.0-bin.tar.gz hive pig_1650382433801.log pig_1650382433801.log pig_1650382433801.log pig_1650382485427.log units

hadoop id.pig pig_1650382485427.log pig_1650382485427.log units.jar

hadoop-core-1.2.1.jar passwd sample.txt

hadoopminhchau@minhchau-master:~$ cd id.out/
hadoopminhchau@minhchau-master:~/id.out$ ls

part-m-00000 _SUCCESS

hadoopminhchau@minhchau-master:~/id.out$ cd
hadoopminhchau@minhchau-master:~$ cat id.out/part-m-00000
```

```
hadoopminhchau@minhchau–master:~$ cat id.out/part–m–00000
root
daemon
bin
sys
sync
games
man
1p
mail
news
uucp
proxy
www–data
backup
list
irc
gnats
nobody
_apt
systemd-timesync
systemd-network
systemd-resolve
messagebus
pollinate
sshd
syslog
uuidd
tcpdump
tss
landscape
usbmux
systemd-coredump
minhchau
1xd
hadoopminhchau
hadoopminhchau@minhchau–master:~$
```

```
$ hadoop fs -mkdir pigdata
$ hadoop fs -put /home/hadoopminhchau/passwd pigdata
$ pig
grunt> A = load './pigdata/passwd' using PigStorage (':');
```

grunt> B = foreach A generate \$0 as id; grunt> dump B;

8. Chạy Pig trên mapreduce

```
hadoopminhchau@minhchau–master:~$ hadoop fs –mkdir pigdata
hadoopminhchau@minhchau–master:~$ hadoop fs –put /home/hadoopminhchau/passwd pigdata
hadoopminhchau@minhchau–master:~$ pig
2022–04–19 16:45:28,814 INFO pig.ExecTypeProvider: Trying ExecType : LOCAL
2022–04–19 16:45:28,815 INFO pig.ExecTypeProvider: Trying ExecType : MAPREDUCE
2022–04–19 16:45:28,815 INFO pig.ExecTypeProvider: Picked MAPREDUCE as the ExecType
2022–04–19 16:45:28,862 [main] ĪNFO org.apache.pig.Main – Apache Pig version 0.17.0 (r1797386) c
iled Jun 02 2017, 15:41:58
2022–04–19 16:45:28,862 [main] INFO org.apache.pig.Main – Logging error messages to: /home/hadoo
nhchau/pig_1650386728851.log
2022–04–19 16:45:28,882 [main] INFO org.apache.pig.impl.util.Utils – Default bootup file /home/h
opminhchau/.pigbootup not found
2022–04–19 16:45:29,134 [main] INFO org.apache.hadoop.conf.Configuration.deprecation – mapred.jo
racker is deprecated. Instead, use mapreduce.jobtracker.address
2022–04–19 16:45:29,134 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngi
- Connecting to hadoop file system at: hdfs://minhchau-master:9000
2022–04–19 16:45:29,507 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngi
- Connecting to map–reduce job tracker at: minhchau–master:9001
2022–04–19 16:45:29,521 [main] INFO org.apache.pig.PigServer – Pig Script ID for the session: PI
efault-f4b14b41-a01e-45a1-82a0-89ea8d1f8093
2022–04–19 16:45:29,521 [main] WARN org.apache.pig.PigServer – ATS is disabled since yarn.timeli
service.enabled set to false
grunt> A = load './pigdata/passwd' using PigStorage(':');
grunt> B = foreach A generate $0 as id;
grunt> dump B;
```

```
Total input paths to process: 1
(root)
(daemon)
(bin)
(sys)
(sunc)
(games)
(man)
(lp)
(mail)
(news)
(uucp)
(proxy)
(www-data)
(backup)
(list)
(irc)
(gnats)
(nobody)
(_apt)
(systemd-timesync)
(systemd-network)
(systemd-resolve)
(messagebus)
(pollinate)
(sshd)
(syslog)
(uuidd)
(tcpdump)
(tss)
(landscape)
(usbmux)
(systemd-coredump)
(minhchau)
(1xd)
(hadoopminhchau)
grunt> _
```