

The page features three large, semi-transparent blue circles of varying sizes. The largest circle is at the top center, a medium-sized one is in the middle, and the largest one is at the bottom right. Thin blue lines connect the top-left corners of these circles, forming a zig-zag pattern across the page.

BigDataBench Subset II

User's Manual

1 Cloud OLTP

We use YCSB to run database basic operations. And, we provide the HBase to run operations for each operation

HBase_Write

To Prepare

1. tar zxvf Central_ones_as_representatives.tar.gz
2. cd \$hbase
3. bin/hbase shell
create 'usertable','f1','f2','f3'

To run

1. cd /Central_ones_as_representatives/Cloud_OLTP_Write/ycsb-0.1.4
2. bin/ycsb load hbase -P workloads/workloadc -p threads=<thread-numbers> -p columnfamily=<family> -p recordcount=<recordcount-value> -p hosts=<hostip> -s>load.dat

A few notes about this command:

- <thread-number> : the number of client threads, this is often done to increase the amount of load offered against the database.
- <family> : In Hbase case, we used it to set database column. You should have database usertable with column family before running this command. Then all data will be loaded into database usertable with column family.
- <recordcount-value>: the total records for this benchmark. For example, when you want to load 10GB data you should set it to 10000000.
- <hostip> : the IP of the hbase's master node.

2 Hadoop-version

Wordcount

To prepare

1. tar zxvf Central_ones_as_representatives.tar.gz
2. cd Central_ones_as_representatives/Hadoop_wordcount
3. sh genData_MicroBenchmarks.sh

To run

sh run_MicroBenchmarks.sh

Naive Bayes

To prepare

1. tar zxvf Central_ones_as_representatives.tar.gz
2. cd Central_ones_as_representatives /Hadoop_Bayes
3. tar zxvf mahout-distribution-0.6.tar.gz
4. export Marginal_ones_as_representatives/Hadoop_Bayes/mahout-distribution-0.6
and then you can run it

Basic command-line usage:

1. cd BigDataBench_V3.0_Hadoop_Hive /E-commerce
2. sh genData_naivebayes.sh

To run

sh run_naivebayes.sh

3 Spark-version

Sort, Grep

(If you use not one machine you must download the spark on each machines, and must download in the right way)

To prepare

1. tar zxvf Central_ones_as_representatives.tar.gz
2. cd Central_ones_as_representatives/Spark_sort_grep
3. sh genData_MicroBenchmarks.sh
4. sh sort-transfer.sh

(before to run the sort you should put the sort-transfer file in your hadoop and then to

run sort-transfer.sh,
the sort-transfer you can find in Central_ones_as_representatives)

To run

when you chose sort like this

```
./run-bigdatabench cn.ac.ict.bigdatabench.Sort <master> <data_file> <save_file>  
[<slices>]
```

parameters:

<master>: URL of Spark server, for example: spark://172.16.1.39:7077

<data_file>: the HDFS path of input data, for example: /test/data.txt

<save_file>: the HDFS path to save the result

[<slices>]: optional, times of number of workers

when you chose grep like this

```
./run-bigdatabench cn.ac.ict.bigdatabench.Grep <master> <data_file> <keyword>  
<save_file> [<slices>]
```

parameters:

<master>: URL of Spark server, for example: spark://172.16.1.39:7077

<data_file>: the HDFS path of input data, for example: /test/data.txt

<keyword>: the keyword to filter the text

<save_file>: the HDFS path to save the result

[<slices>]: optional, times of number of workers

PageRank

To prepare

1. tar zxvf Central_ones_as_representatives.tar.gz
2. cd /Central_ones_as_representatives/Spark_Pagerank/Pagerank
3. sh genData_PageRank.sh

To run

```
./run-bigdatabenchorg.apache.spark.examples.PageRank  
<master> <file> <number_of_iterations> <save_path> [<slices>]
```

parameters:

#<master>: URL of Spark server, for example: spark://172.16.1.39:7077

#<file>: the HDFS path of input data, for example: /test/data.txt

#<number_of_iterations>: number of iterations to run the algorithm

#<save_path>: path to save the result

#[<slices>]: optional, times of number of workers

Kmeans

The Kmeans program we used is obtained from Mahout.

(If you use not one machine you must download the spark on each machines, and must download in the right way)

To prepare

1. tar zxvf Central_ones_as_representatives.tar.gz
2. cd Central_ones_as_representatives/Spark_Kmeans/Kmeans
3. sh genData_Kmeans.sh

To run

```
./run-bigdatabench org.apache.spark.mllib.clustering.KMeans <master> <input_file>  
<k> <max_iterations> [<runs>]
```

parameters:

#<master>: URL of Spark server, for example: spark://172.16.1.39:7077

#<input_file>: the HDFS path of input data, for example: /test/data.txt

#[<k>]: number of centers

#<max_iterations>: number of iterations to run the algorithm

#[<runs>]: optional, level of parallelism to split computation into

4 Hive-version

Hive_Aggregation_AVG

To prepare

Generat Table_data

1. tar xzf Central_ones_as_representatives.tar.gz
2. cd /Central_ones_as_representatives/BigDataGeneratorSuite/Table_datagen/e-com
3. java -XX:NewRatio=1 -jar pdgf.jar -l demo-schema.xml -l demo-generation.xml -c
-s -sf <parameter>
4. mkdir BigDataBench/BigDataGeneratorSuite/Table_datagen/
5.mv

```
Central_ones_as_representatives/BigDataGeneratorSuite/Table_datagen/e-com/output  
/BigDataBench/BigDataGeneratorSuite/Table_datagen
```

Greate table

1. cd \$HIVE_HOME/bin
2. ./hive
create database bigdatabench;
use bigdatabench;
create table bigdatabench_dw_order(order_id int,buyer_id int,create_date string)

```
ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE;
```

```
create table bigdatabench_dw_item(item_id int,order_id int,goods_id int,goods_number double,goods_price double,goods_amount double) ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE;
load data local inpath 'BigDataBench/BigDataGeneratorSuite/Table_datagen/output/OS_ORDER.txt' overwrite into table bigdatabench_dw_order;
load data local inpath 'BigDataBench /BigDataGeneratorSuite/Table_datagen/output/OS_ORDER_ITEM.txt' overwrite into table bigdatabench_dw_item;
create table item_temp as select ORDER_ID from bigdatabench_dw_item;
```

To run

```
cd Central_ones_as_representatives/Hive_Aggregation_AVG
sh BigOP-e-commerce-aggregationAVG.sh
(For ease of use, we recommend that you use a local mysql server to store metadata)
```

Hive_Aggregation_Query

To prepare

Generat Table_data

1. tar xzf Central_ones_as_representatives.tar.gz
2. cd /Central_ones_as_representatives/BigDataGeneratorSuite/Table_datagen/e-com
3. java -XX:NewRatio=1 -jar pdgf.jar -l demo-schema.xml -l demo-generation.xml -c -s - sf <parameter>
4. mkdir BigDataBench/BigDataGeneratorSuite/Table_datagen/
5. mv Central_ones_as_representatives/BigDataGeneratorSuite/Table_datagen/e-com/output /BigDataBench/BigDataGeneratorSuite/Table_datagen

Greate table

1. cd \$HIVE_HOME/bin
2. ./hive
create database bigdatabench;
use bigdatabench;
create table bigdatabench_dw_order(order_id int,buyer_id int,create_date string) ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE;

create table bigdatabench_dw_item(item_id int,order_id int,goods_id

```
int,goods_number double,goods_price double,goods_amount double)
ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS
TEXTFILE;
load data local inpath
'BigDataBench/BigDataGeneratorSuite/Table_datagen/output/OS_ORDER.txt'
overwrite into table bigdatabench_dw_order;
load data local inpath 'BigDataBench
/BigDataGeneratorSuite/Table_datagen/output/OS_ORDER_ITEM.txt' overwrite
into table bigdatabench_dw_item;
create table item_temp as select ORDER_ID from bigdatabench_dw_item;
```

To run

```
cd Central_ones_as_representatives/Hive_AggregationQuery
sh e-commerce-aggregation.sh
```

(For ease of use, we recommend that you use a local mysql server to store metadata)

Hive_TPC_DS_query13

To prepare

Generate data

1. cd Central_ones_as_representatives/Hive_TPC_DS_query13
2. ./dsdgen -scale <data-size> -dir <data-directory>

A few notes about this command:

- <data-size> : the number of the data size, the unite is GB.
 - <data-directory> :the path where your data want to put.
- ```
./mvdata.sh <data-directory>
```

Create the TPC-ds tables

1. Firstly, you should put TPC-DS data into the hdfs.
2. cd Central\_ones\_as\_representatives/Hive\_TPC\_DS\_query13
3. sh create-tpcds-tables.sh

### **To run**

```
cd Central_ones_as_representatives/Hive_TPC_DS_query13
sh query13.sh
```

# 5 Shark

## Cross,Union

### To prepare

Generate data

1. tar xzf Central\_ones\_as\_representatives.tar.gz
- 2.cd / Central\_ones\_as\_representatives /BigDataGeneratorSuite/Table\_datagen/e-com
3. java -XX:NewRatio=1 -jar pdgf.jar -l demo-schema.xml -l demo-generation.xml -c -s - sf <parameter>

Upload the text files in

Central\_ones\_as\_representatives /BigDataGeneratorSuite/Table\_datagen/e-com /output/ to HDFS and make sure these files in different pathes.

Create tables

1. cd Central\_ones\_as\_representatives / Shark\_cross\_union
2. shark

```
create external table bigdatabench_dw_item(item_id int,order_id int,goods_id int,goods_number double,goods_price double,goods_amount double) ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE LOCATION 'path to OS_ODER_ITEM.txt';
```

```
create external table bigdatabench_dw_order(order_id int,buyer_id int,create_date string) ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE LOCATION 'path to OS_ORDER.txt';
```

```
create table item_temp as select ORDER_ID from bigdatabench_dw_item;
```

### To run

```
cd Central_ones_as_representatives/Shark_cross_union
```

edit free\_m.sh to make sure it runs correctly.

```
shark -f crossproduct.sql
```

```
shark -f union.sql
```

# 6 Impala

## Orderby

### To prepare

Generate data

1. tar xzf Central\_ones\_as\_representatives.tar.gz



2. cd /Central\_ones\_as\_representatives/BigDataGeneratorSuite/Table\_datagen/e-com
3. java -XX:NewRatio=1 -jar pdgf.jar -l demo-schema.xml -l demo-generation.xml -c -s - sf <parameter>

Upload the text files in

Marginal\_ones\_as\_representatives/BigDataGeneratorSuite/Table\_datagen/e-com /output/ to HDFS and make sure these files in different pathes

Create tables

1. cd \$HIVE\_HOME/bin

2. ./hive

```
create external table bigdatabench_dw_item(item_id int,order_id int,goods_id int,goods_number double,goods_price double,goods_amount double) ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE LOCATION 'path to OS_ODER_ITEM.txt';
```

```
create external table bigdatabench_dw_order(order_id int,buyer_id int,create_date string) ROW FORMAT DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE LOCATION 'path to OS_ORDER.txt';
```

```
create table item_temp as select ORDER_ID from bigdatabench_dw_item;
```

### To run

cd Central\_ones\_as\_representatives/Impala\_aggregationMax

edit free\_m.sh and impala-restart.sh to make sure them run correctly.

sh runMicroBenchmark.sh

## Impala\_TPC\_DS\_query

### To prepare

Generate data

1. cd Central\_ones\_as\_representatives/Impala\_TPC\_DS\_query

2. ./dsdgen -scale <data-size> -dir <data-directory>

A few notes about this command:

- <data-size> : the number of the data size, the unite is GB.
  - <data-directory> :the path where your data want to put.
- ./mvdata.sh <data-directory>

Create the TPC-ds tables

1. Firstly, you should put TPC-DS data into the hdfs.

2. cd Marginal\_ones\_as\_representatives/Hive\_TPC\_DS\_query3

3. sh create-tpcds-tables.sh

### To run

cd Central\_ones\_as\_representatives/Impala\_TPC\_DS\_query

sh runTPC-DS.sh