# Hive, PIG, Hadoop benchmark results

June 18, 2009


We have set up a 10 node cluster to test the performance of Hive, PIG and Hadoop based on the data and queries from a SIGMOD 09 paper:

> A. Pavlo, E. Paulson, A. Rasin, D. J. Abadi, D. J. Dewitt, S. Madden, and M. Stonebraker, "A Comparison of Approaches to Large-Scale Data Analysis," in SIGMOD '09: Proceedings of the 2009 ACM SIGMOD International Conference, 2009

Because the latest PIG works with Hadoop 0.18.x, we used Hadoop version 0.18.3. We used PIG truck version 786346. For Hive, we used hive trunk version 786346. Information on how we installed and set up those systems can be found in README.

Aside from the basic settings, we also configure Hive, PIG and Hadoop to use Gzip to compress both the intermediate data and final output. Combiners are enabled in PIG.

All the test data and queries are pulled from the SIGMOD 09 paper. The amount of the data per node is scaled accordingly so that the queries would not take more than 30 minutes. We tested the first four of the five queries experimented in their paper. The last query requires writing particular mappers and could not be done with SQL style queries. As a result, that query will not give us any hints on the performance of Hive and PIG, so we ignored it. The queries are listed in their website (http://database.cs.brown.edu/projects/mapreduce-vs-dbms/), including two selects, one aggregation and one join. For the join query, we slightly changed the date parameters so that we get enough output records. Details on how to generate the data can also be found in README.

For any particular query, Hadoop, hive and PIG are configured to use the same number of mappers and the same number of reducers if possible. The numbers are also included in the following table. All the data are pre-loaded into Hadoop Distributed File System (HDFS) before running the benchmarks. Basically, all the loads and stores in those queries happen in HDFS.

The timings are averaged over 3 runs. In each run, we record the time of the entire command, including query compilation time, table creation time, and query execution time and so on.

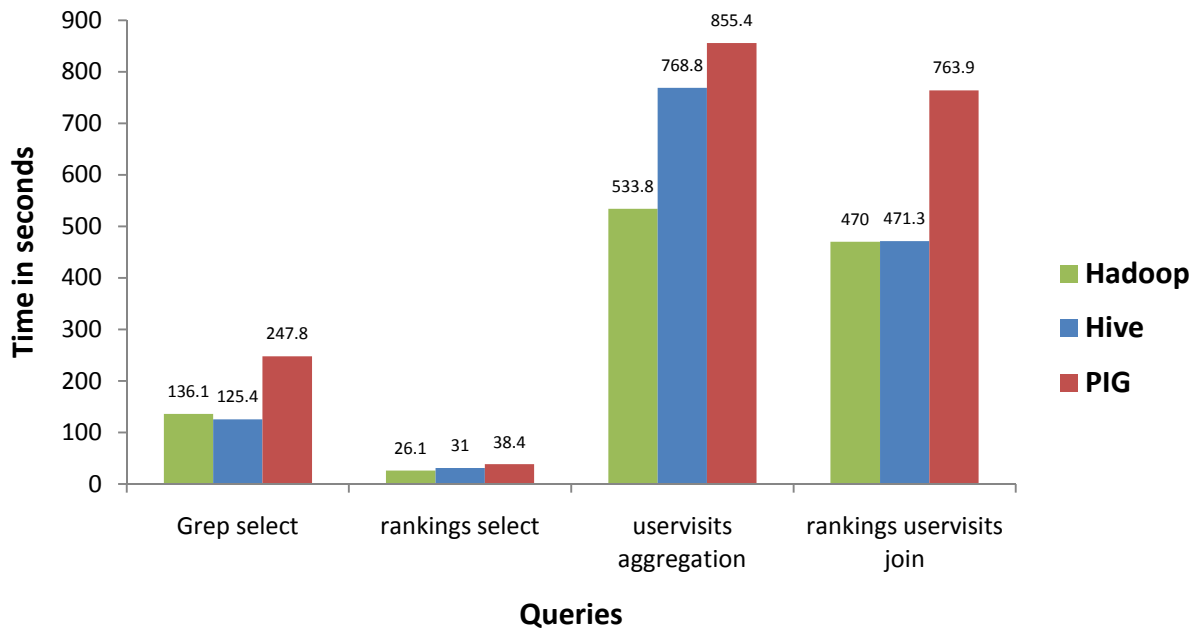The final timing result is shown in Figure 1.

# Hive, PIG and Hadoop benchmark



900
800
700

600

Time in seconds

500
400
300
200
100
0

|          | Grep select | rankings select | uservisits aggregation | rankings uservisits join |
| Hadoop | 136.1 | 26.1 | 533.8 | 470 |
| Hive | 125.4 | 31 | 768.8 | 471.3 |
| PIG | 247.8 | 38.4 | 855.4 | 763.9 |

**Queries**

Figure 1  Hive, PIG and Hadoop benchmark timing result

Here are more details about the queries and the Hadoop jobs.

| Select query1 | SELECT * FROM grep WHERE field like '%XYZ%'; |
|---|---|
| Data | The table "grep" has two columns: (key STRING, field STRING). It has 500,000,000 rows and takes 50 GB disk space. |
| #Mappers #Reducers | The query is finished in a single job with 380 mappers and 0 reducers on all of Hive, PIG and Hadoop. |
| Timing | Hadoop took 136 seconds. Hive took 125 seconds. PIG took 248 seconds. Hive took 8% less time than Hadoop. PIG took 82% more time than Hadoop and 98% more time than hive |

Table 1 "Grep" select query

| Select query2 | SELECT pageRank, pageURL FROM rankings WHERE pageRank > 10; |
|---|---|
| Data | The table "rankings" has three columns: (pageRank INT, pageURL STRING, avgDuration INT). It has 56289700 rows and takes 3.3 GB disk space. |
| #Mappers #Reducers | The query is finished in a single job with 30 mappers and 0 reducers on all of Hive, PIG and Hadoop. |
| Timing | Hadoop took 26 seconds. Hive took 31 seconds. PIG took 38 seconds. Hive took 19% more time than Hadoop. PIG took 46% more time than Hadoop and 23% more time than hive. |

Table 2 "rankings" select query

| Aggregation query | SELECT sourceIP, SUM(adRevenue) FROM uservisits GROUP BY sourceIP; |
|---|---|
| Data | The table "uservisits" has 9 columns: (sourceIP STRING,destURL STRING,visitDate STRING,adRevenue DOUBLE,userAgent STRING,countryCode STRING,languageCode STRING,searchWord STRING,duration INT ). It has 465000000 rows and takes 60GB disk space. |
| #Mappers #Reducers | The query is finished in a single job with 450 mappers and 60 reducers on all of Hive, PIG and Hadoop. |
| Timing | Hadoop took 534 seconds. Hive took 769 seconds. PIG took 855 seconds. Hive took 44% more time than Hadoop. PIG took 60% more time than Hadoop and 11% more time than hive. |

**Table 3 "uservisits" aggregation query**

| Join query | SELECT sourceIP, AVG(pageRank), SUM(adRevenue) FROM rankings R JOIN (SELECT * FROM uservisits UV WHERE UV.visitDate > '1999-01-01' AND UV.visitDate < '2000-01-01') NUV ON (R.pageURL = NUV.destURL) GROUP BY sourceIP; |
|---|---|
| Data | The tables "rankings" and "uservisits" are the same as in previous queries. |
| #Mappers #Reducers | The query is finished in two jobs on Hive and PIG and three jobs on Hadoop. The first job, which takes the majority of the time, has 480 mappers and 60 reducers on all Hive, PIG and Hadoop. The second job on PIG has 120 mappers and 60 reducers. The second job on both Hive and Hadoop has 60 mappers and 60 reducers. The third job on Hadoop has 60 mappers and 1 reducer. |
| Timing | Hadoop took 470 seconds. Hive took 471 seconds. PIG took 764 seconds. Hive took 0.2% more time than Hadoop. PIG took 63% more time than Hadoop and 62% more time than hive. |

**Table 4 "rankings" and "uservisits" join query**

If you have any questions or suggestions about this benchmark, please email to Yuntao Jia (yjia@facebook.com).