

### CSE 587 HW # 3

**Name:** Daniel Nazareth

**UbitName:** dnazaret

**Email:** [dnazaret@buffalo.edu](mailto:dnazaret@buffalo.edu)

**SUMMARY:** The project uses the Apache PIG and HIVE frameworks each layered on Hadoop 2.6.0 to quickly and efficiently compute the most and least volatile stocks from small, medium and large sized NASDAQ stock datasets.

**RATIONALE FOR USING PIG/HIVE:** Pig and Hive are both frameworks sitting on top of Hadoop. The biggest reason for using either framework over traditional Hadoop MapReduce code (in say Java or Python) is that Pig Lating and HiveQL statements encapsulate huge chunks of MapReduce functionality into very simple syntax and readable code. In general Pig Lating and HiveQL scripts will use approximately 5% of the lines of code taken by traditional Java mapReduce to accomplish the same task. This means that development and QA time is reduced tremendously and data scientists are free to program short scripts at a high level without worrying too much about what's going on "under the hood". The only downside to this is that Pig/Hive scripts can run marginally slower than the corresponding Java mapreduce code.

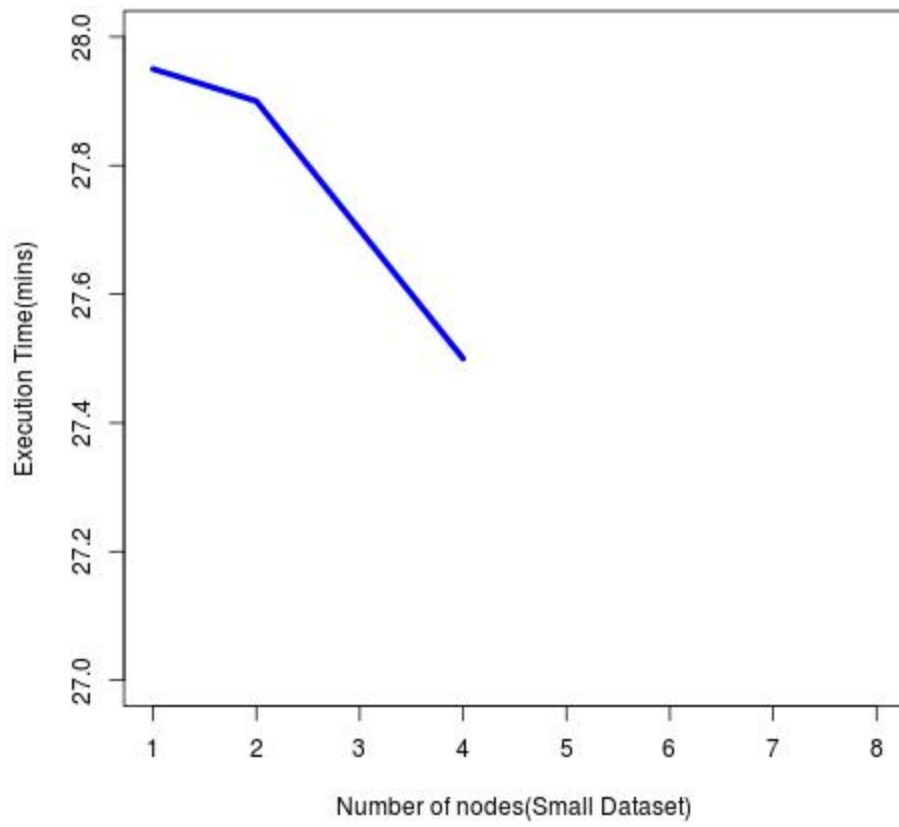
**To summarise, Pig/Hive offer the following benefits:**

- Shorter, simpler and much more readable code.
- Shorter development and QA cycles leading to quicker releases to production.
- Data scientist is free to program and think at a high level without worrying about low-level, specific details of the mapreduce.

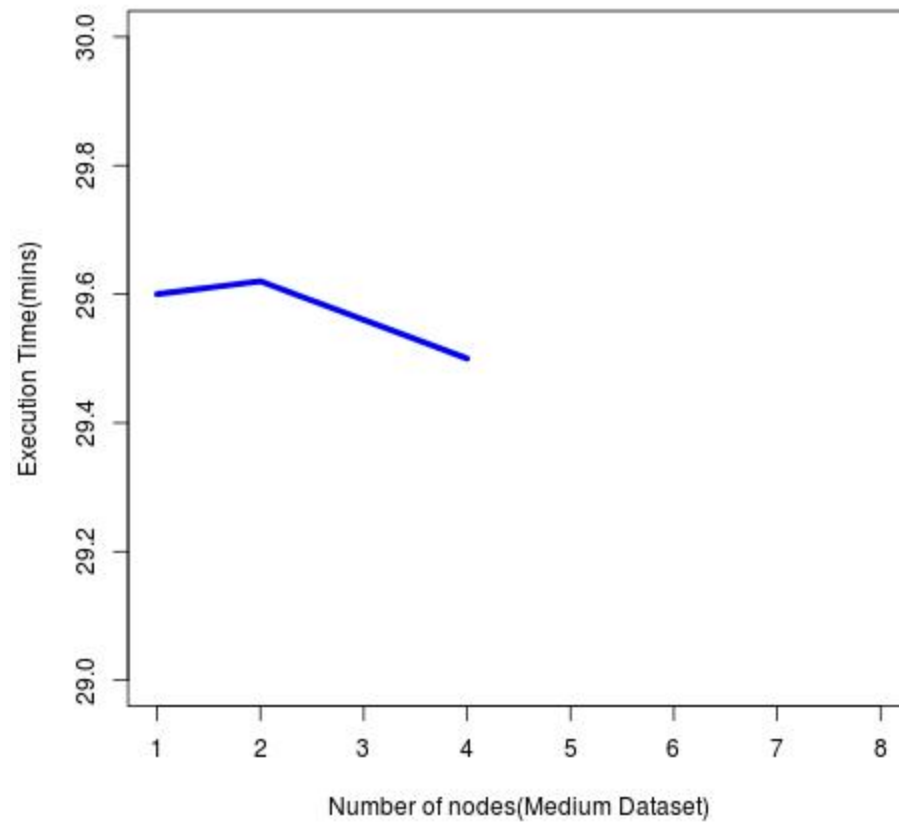
**RESULTS:** The results are summarised in tabular and graphical form below for each of the small, medium and large datasets. This can be verified by running the attached code at the command line as well.

#### **PIG RESULTS**

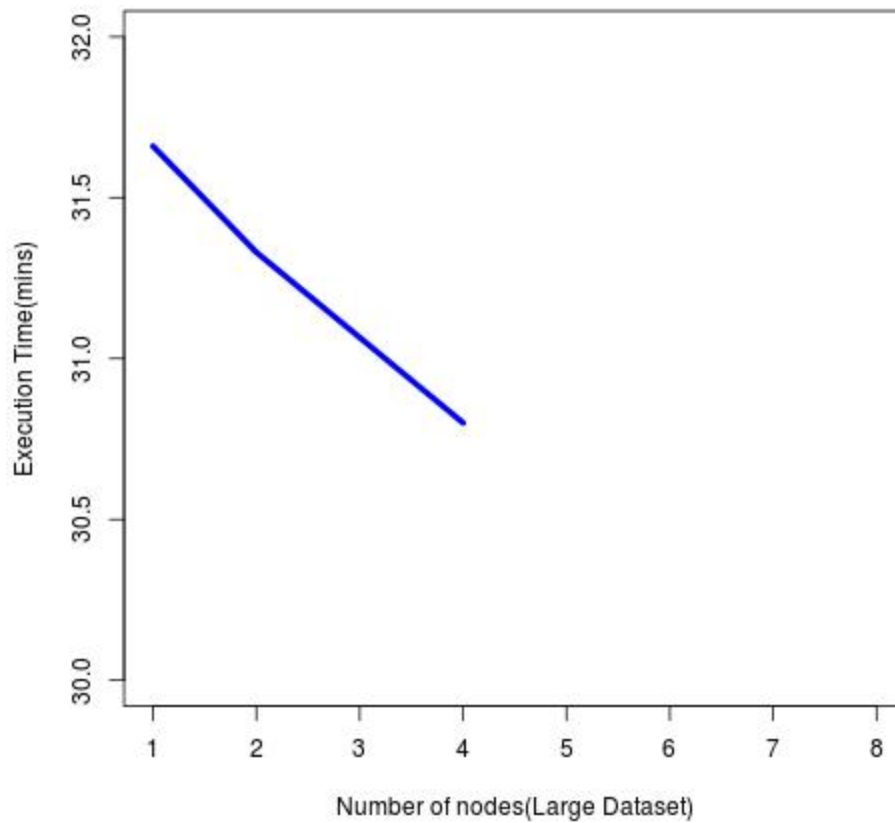
<b><u>PROBLEM SIZE</u></b>	<b><u>EXECUTION TIME(1 NODE,12 CORES)</u></b>	<b><u>EXECUTION TIME(2 NODES,12 CORES)</u></b>	<b><u>EXECUTION TIME(4 NODES,12 CORES)</u></b>
<b><i>Small</i></b>	27 mins,59 seconds	27 mins, 54 seconds	27 mins,37 seconds
<b><i>Medium</i></b>	29 mins,44 seconds	29 mins,45 seconds	29 mins,31 seconds
<b><i>Large</i></b>	31 mins, 46 seconds	31 mins,21 seconds	30 mins,49 seconds



**MEDIUM DATASET:**



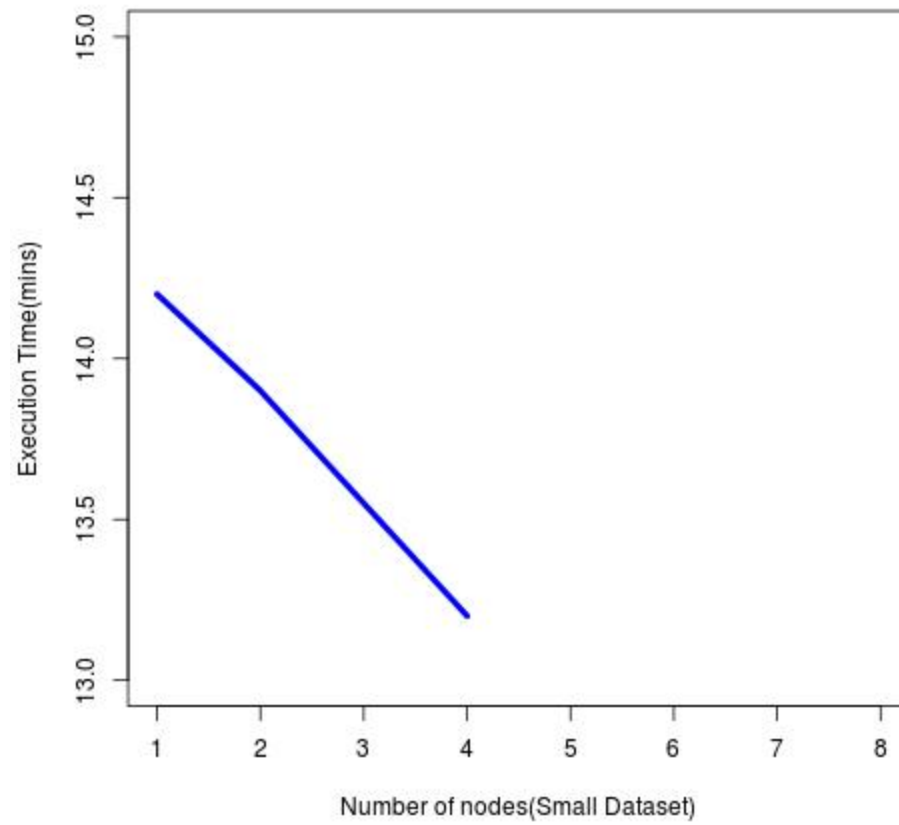
**LARGE DATASET:**



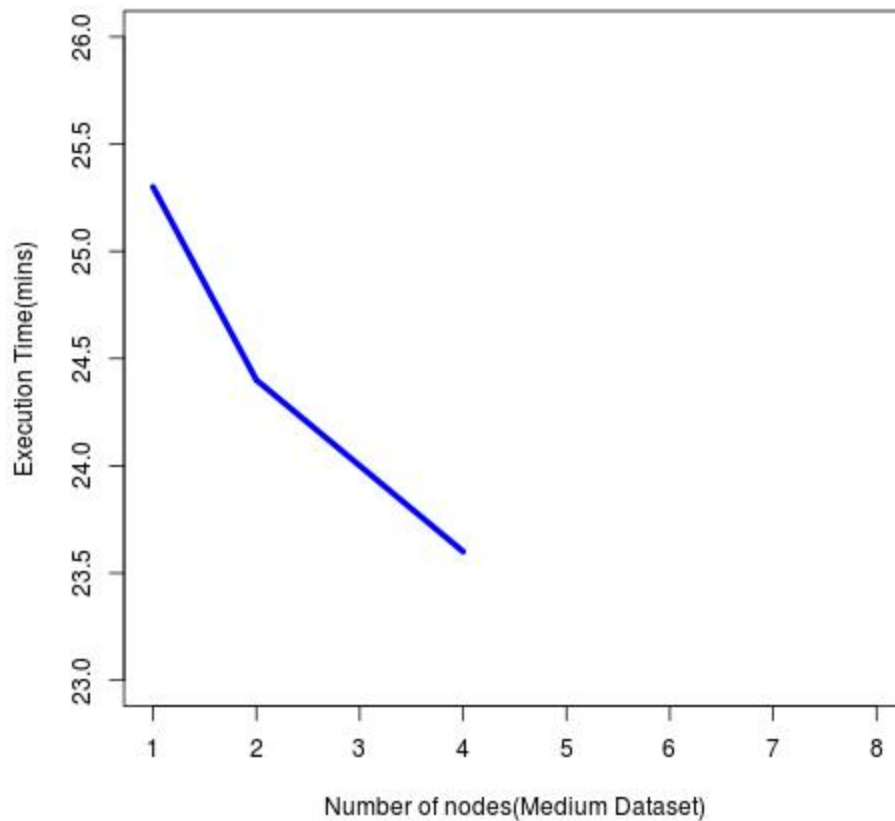
### HIVE RESULTS

<u>PROBLEM SIZE</u>	<u>EXECUTION TIME(1 NODE,12 CORES)</u>	<u>EXECUTION TIME(2 NODES,12 CORES)</u>	<u>EXECUTION TIME(4 NODES,12 CORES)</u>
<b>Small</b>	14 mins,14 seconds	13 mins, 56 seconds	13 mins,15 seconds
<b>Medium</b>	25 mins,17 seconds	24 mins,22 seconds	23 mins,36 seconds
<b>Large</b>	Memory Related Exception: Garbage Collection Overhead Exceeded.		

**SMALL DATASET GRAPH/SCALING:**



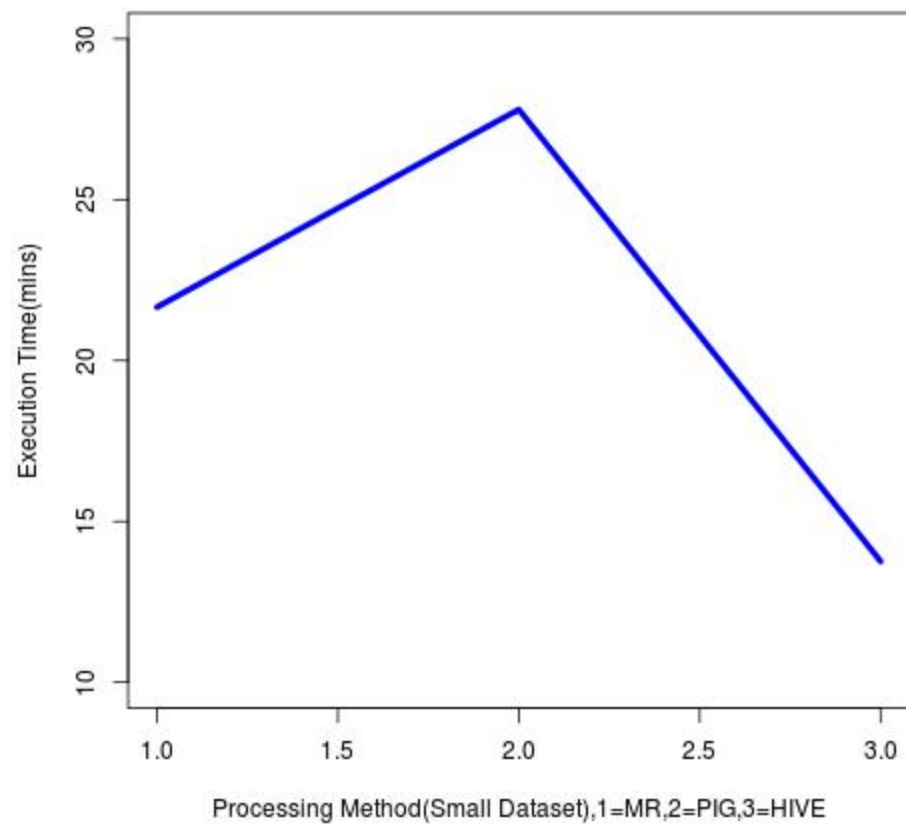
**MEDIUM DATASET GRAPH/SCALING:**



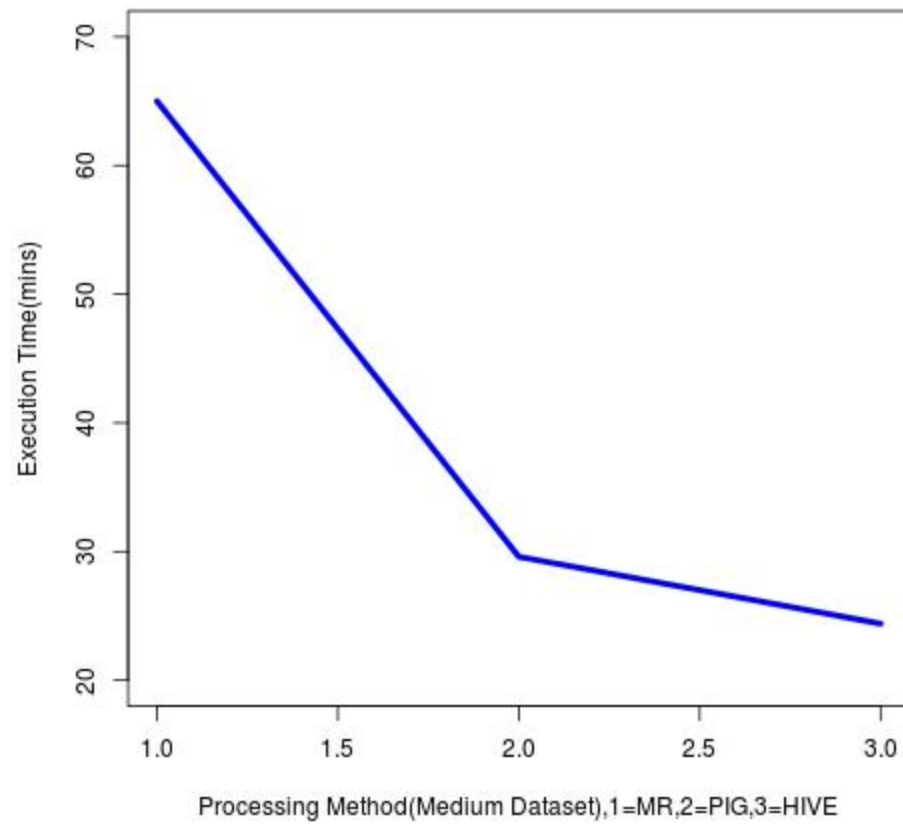
**COMPARISON OF MAP REDUCE/PIG/HIVE:** We see from the below tabular and graphical comparison that for small datasets, all 3 methods are closely matched with Hive being the runaway winner. For medium and large datasets Pig and Hive in particular offer significantly superior performance. However Map Reduce performance can be significantly improved by assigning more computing nodes-this cannot be done for Pig/Hive.

<u>PROBLEM SIZE</u>	<u>AVERAGE EXECUTION TIME(MAP REDUCE)</u>	<u>AVERAGE EXECUTION TIME(PIG)</u>	<u>AVERAGE EXECUTION TIME(HIVE)</u>
<b>Small</b>	21.66 minutes	27.8 mins	13.76 minutes
<b>Medium</b>	65 mins	29.6 mins	24.4 minutes
<b>Large</b>	120 mins	31.2 mins	Not Available

**SMALL GRAPHICAL COMPARISON:**

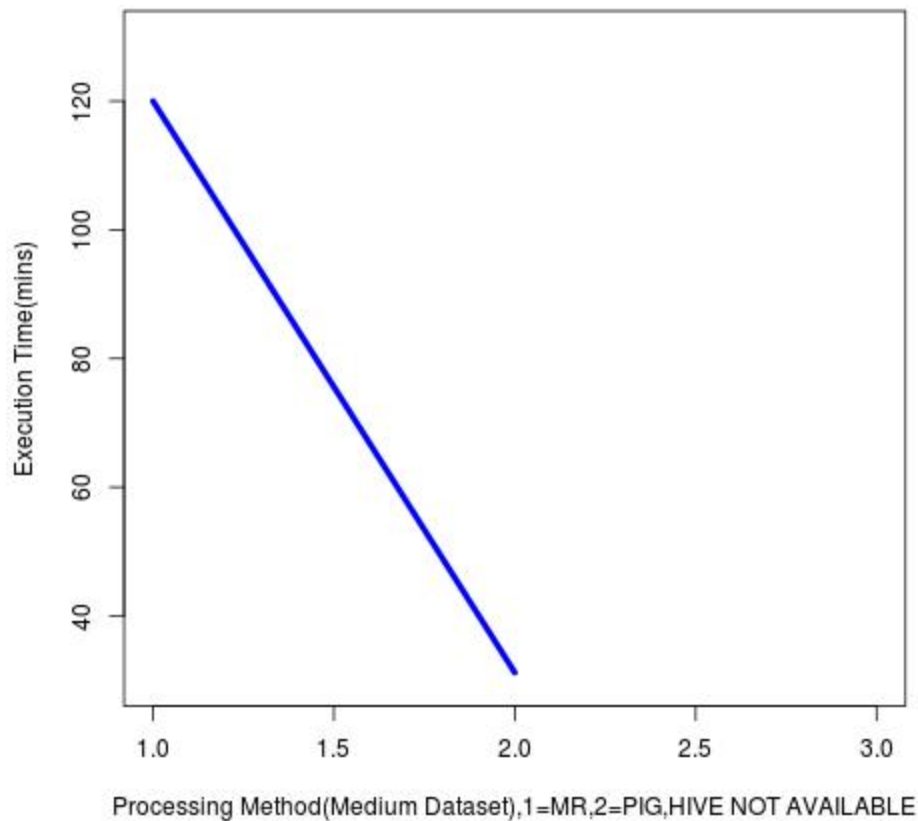


**MEDIUM GRAPHICAL COMPARISON:**



**LARGE GRAPHICAL COMPARISON:**





### **NOTES/ISSUES FACED:**

- We observe that additional nodes makes little or no difference to execution times for both Pig and Hive but a huge difference for Map Reduce. Times scale by small factors over dataset sizes but barely at all over number of nodes allocated. This would seem to indicate that for a given Hive or Pig script, computational efficiency cannot be reduced beyond a certain point, no matter how many resources are allocated.
- Memory overhead seemed to be excessive for Hive large dataset and hence job could not complete on CCR. Possible improvements to code efficiency could possibly help remedy this problem.

### **REFERENCES**

<https://cwiki.apache.org/confluence/display/Hive/Tutorial>

[http://www.tutorialspoint.com/hive/hive\\_installation.htm](http://www.tutorialspoint.com/hive/hive_installation.htm)

[http://chimera.labs.oreilly.com/books/1234000001811/ch05.html#foreach\\_basic](http://chimera.labs.oreilly.com/books/1234000001811/ch05.html#foreach_basic)

[https://pig.apache.org/docs/r0.7.0/piglatin\\_ref2.html](https://pig.apache.org/docs/r0.7.0/piglatin_ref2.html)