

# Hive

DIKSHA YADAV<sup>1,\*</sup>, +

<sup>1</sup> School of Informatics and Computing, Bloomington, IN 47408, U.S.A.

\* Corresponding authors: yadavd@umail.iu.edu

+ HID - S17-IR-2044

Paper-002, April 3, 2017

Hive is an open source data warehousing solution which is built on top of Hadoop. It structures data into understandable and conventional database terms like tables, columns, rows and partitions. It supports HiveQL queries which have structure like SQL queries. HiveQL queries are compiled to map reduce jobs which are then executed by Hadoop. Hive also contains Metastore which includes schemas and statistics which is useful in query compilation, optimization and data exploration.

© 2017 <https://creativecommons.org/licenses/>. The authors verify that the text is not plagiarized.

**Keywords:** Hive, Hadoop, HiveQL, SQL

<https://github.com/diksha2112/sp17-i524/tree/master/paper2/S17-IR-2044/report.pdf>

This review document is provided for you to achieve your best. We have listed a number of obvious opportunities for improvement. When improving it, please keep this copy untouched and instead focus on improving report.tex. The review does not include all possible improvement suggestions and for each comment you may want to check if it applies elsewhere in the document.

Please include the correct link to your report on the cloudmesh user's Github account, not *diksha2112*'s.

Your paper does not meet the length requirements for the assignment. Moreover, what you have included is more like a random selection of tidbits about Hive than a coherent paper. You need to rewrite the paper.

Assessment: Major revisions required.

## INTRODUCTION

The reason behind development of Hive is making it easier for end users to use Hadoop. Map reduce

MapReduce

programs were required to be developed by users for simple to complex tasks. It lacked expressiveness like query language.

Not the case. MapReduce has more expressiveness than Hive, that is, you can express more programs with MapReduce than with Hive. Hive simply makes common use cases of MapReduce easier to write in a familiar, SQL-like environment.

So, it was a time consuming and difficult task for end users to

use Hadoop. For solving this problem Hive was built in January 2007 and open sourced in August 2008. Hive is an open source data warehousing solution which is built on top of Hadoop. It structures data into understandable and conventional database terms like tables, columns, rows and partitions. It supports HiveQL queries which have structure like SQL queries. HiveQL queries are compiled to map reduce

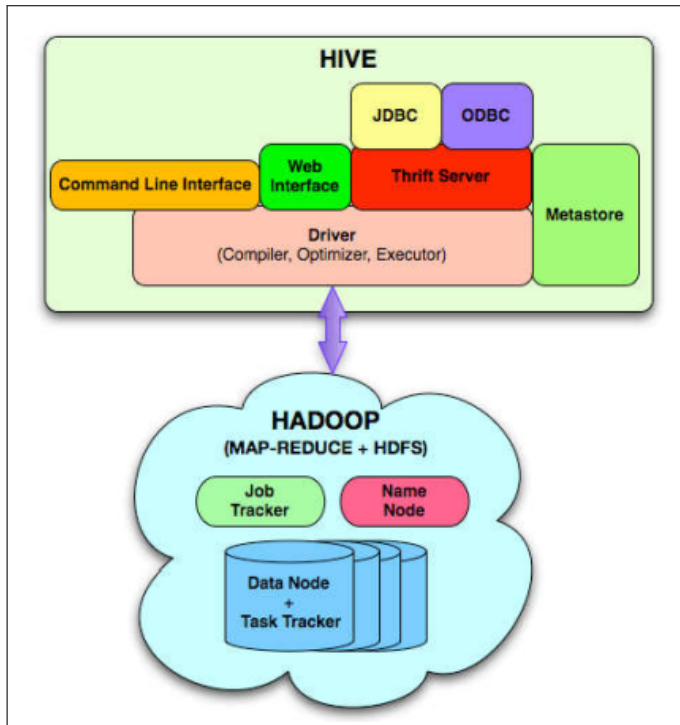
Term

jobs which are then executed by Hadoop. Hive also contains Metastore which includes schemas and statistics which is useful in query compilation, optimization and data exploration[1]

## ARCHITECTURE

This section is incomplete. There is very little detail about the different parts of the architecture, there is no motivation for why they exists, or explanation of how they fit together. They are not even formatted properly in the text. The architecture figure is not referenced anywhere in the text. Please revisit.

Hive architecture includes, Database-It consists of tables created by the user. Metastore-It contains information about the system. It can be accessed by different components as and when needed. Interfaces-User interface and Application programming interface both are present in hive. Driver-manage HiveQL statements at every stage. Query compiler-It compiles HiveQL queries to acyclic graphs (directed) representing map reduce tasks. Execution Engine- It executes the tasks generated by the compiler. Hive Server- It provide JDBC/ODBC server and thrift interface[2]



**Fig. 1.** Hive Architecture

Did you create this figure? If not, you need to provide a reference to where you got it from.

## HIVEQL QUERY FORMAT

```
SELECT [ALL | DISTINCT] select-expr, select-expr, ...
FROM table-reference
[WHERE where-condition]
[GROUP BY col-list]
[HAVING having-condition]
[CLUSTER BY col-list | [DISTRIBUTE BY col-list] [SORT BY col-list]]
[LIMIT number]
[3]
```

This is out of scope for the paper. In an overview paper like this you don't need to provide the format for a single SELECT statement. Anyone can look this up in the docs.

## SYSTEM REQUIREMENTS

Hive is cross platform. So, It does not need any specific operating system to work.

But what are the system requirements?

## COMPARISON OF HIVE WITH OTHER TRADITIONAL DATABASES

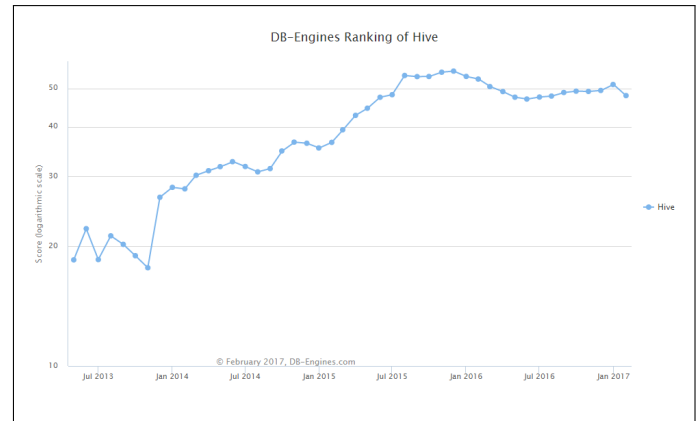
Traditional databases follow schema on write approach while Hive follows schema on read approach. In schema on write, databases checks at load time if the data follows the table representation given by user while in schema on read approach, it is checked at run time only. This saves the time for hive to load the data when traditional databases takes longer time[4]

Is this the only or even most important point to compare Hive against other platforms on?

## POPULARITY OF HIVE

The popularity of hive increasing with time. This can be proved by the following plot made by DB Engines Ranking. It ranks database management systems according to their status and popularity. Following plot shows popularity of hive with time.

Why is it important that Hive's popularity has increased over time? In addition, this is not the proper way to cite a source.



**Fig. 2.** Hive Popularity

[5]

## RESOURCES FOR LEARNING HIVE

Someone new to hive can start learning it by going through the following links in sequence: Install Hive[https://www.edureka.co/blog/apache-hive-installation-on-ubuntu?utm\\_source=quora&utm\\_medium=crosspost&utm\\_campaign=social-media-edureka-ab](https://www.edureka.co/blog/apache-hive-installation-on-ubuntu?utm_source=quora&utm_medium=crosspost&utm_campaign=social-media-edureka-ab)

Hive Tutorial[https://www.edureka.co/blog/hive-tutorial/?utm\\_source=quora&utm\\_medium=crosspost&utm\\_campaign=social-media-edureka-ab](https://www.edureka.co/blog/hive-tutorial/?utm_source=quora&utm_medium=crosspost&utm_campaign=social-media-edureka-ab)

Top Hive commands with examples[https://www.edureka.co/blog/hive-commands-with-examples?utm\\_source=quora&utm\\_medium=crosspost&utm\\_campaign=social-media-edureka-ab](https://www.edureka.co/blog/hive-commands-with-examples?utm_source=quora&utm_medium=crosspost&utm_campaign=social-media-edureka-ab)

## ACKNOWLEDGEMENT

I am also grateful to Dr. Gregor von Laszewski for providing the appropriate paper template.

## CONCLUSION

Since Hive is making the use of Hadoop easier for users, its popularity is increasing with time.

The conclusion needs to summarize the paper. I don't know what to make of this sentence about Hive's popularity.

## REFERENCES

- [1] "Hive," Web Page, online; accessed 24-March-2017. [Online]. Available: [https://en.wikipedia.org/wiki/Apache\\_Hive](https://en.wikipedia.org/wiki/Apache_Hive)
- [2] "Hive article," Web Page, online; accessed 24-March-2017. [Online]. Available: <https://docs.treasuredata.com/articles/hive>
- [3] A. T. J. S. N. J. Z. S. P. C. S. A. H. L. P. W. R. Murthy, "Hive-a warehousing solution over map reduce framework," Paper, online; accessed 23-March-2017. [Online]. Available: <http://laser.inf.ethz.ch/2013/material/breitman/additionalpercent20reading/hive.pdf>
- [4] —, "Hive-a petabyte scale datawarehouse using hadoop," Paper. [Online]. Available: <http://infolab.stanford.edu/~ragho/hive-icde2010.pdf>
- [5] "Ranking hive," Web Page, online; accessed 20-March-2017. [Online]. Available: [http://db-engines.com/en/ranking\\_trend/system/Hive](http://db-engines.com/en/ranking_trend/system/Hive)