# A Comparative Analysis of Apache hive based on MapReduce and Impala based on distributed query engine

Prasad Mitkari[1], Rutuja Banswal [2], Akshata Sirsat[3], Prof. Daivashala Deshmukh[4]

B.Tech Student, Department of Computer Science & Engineering,  Dr. BAMU, Maharashtra Institute of Technology, Aurangabad (MS), India

B.Tech Student, Department of Computer Science & Engineering, Dr. BAMU, Maharashtra Institute of Technology, Aurangabad (MS), India

B.Tech Student, Department of Computer Science & Engineering, Dr. BAMU, Maharashtra Institute of Technology, Aurangabad (MS), India

Assistant Professor, Department of Computer Science & Engineering, Dr. BAMU, Maharashtra Institute of Technology, Aurangabad (MS), India

**ABSTRACT:** Big Data realm moves around 5 Vs- volume, velocity, variety, value and veracity. Analyzing and Storing huge amount of data available in different formats which is increasing with huge velocity to gain values out and it is itself a big deal. Faster and Quick query in the Big Data is very important for getting the valuable information to improve the system performance. To get this goal, number of research institutions and internet companies develop some tools which are respectively Hive (based on MapReduce), Impala (based on Distributed query engine). In this paper, we compare these two-type of query tools and we find the time taken by Impala and Hive to execute the same query. In this paper we contrast the relative merits of two technologies called Apache Hive and Cloudera's Impala. In efficiency and performance, total time taken by both for their execution and their requirements. Basically we focus on performance analysis of both Apache Hive and Impala using Cloudera Distribution Including Apache Hadoop (CDH) and movielens datasets.

**KEYWORDS**: Hadoop, Hive, Impala, Framework, MapReduce, Analysis

## I. INTRODUCTION

Big data is a large amount of drastically increasing real world data. Data is increasing every minute, every second in a day. This large amount of data is used to derive knowledge from the raw data, with this in mind, big data is the computing strategy and technology that are used to handle large datasets. There are five aspects which define the big data such as volume, velocity, variety, veracity, and value. Big Data is usually difficult to store and process using traditional data management systems. Apache Software Foundation introduced processing tools to solve processing challenges, which are used to solve big data related problems and used to derive the patterns and trends from data.

**HIVE**: Apache Hive is a tool, which built on top of Apache Hadoop framework. At initial phase, it was being developed by Facebook. Later it was given to Apache Software Foundation and developed further as open-source software called as Apache Hive. Hive is mainly used for processing unstructured data. This processing is done in three main steps summarization, query, and analysis. Hive is primarily used for data mining purpose.  The interaction with the Hive is done by SQL like query language called HiveQL. Due to its HiveQL, it is very familiar and user-friendly, fast, scalable, and extensible.
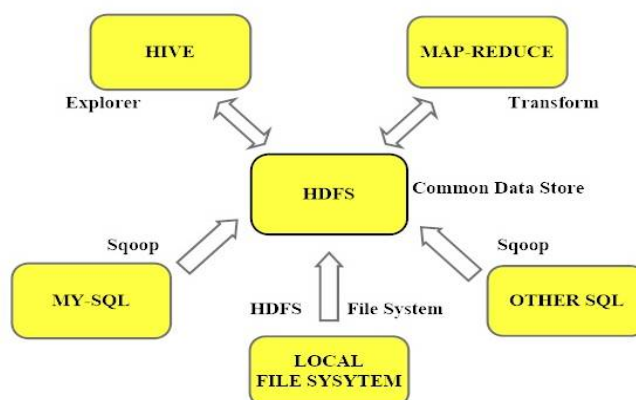
Figure 1: Hive architecture

**IMPALA**: Impala uses Massive Parallel Processing (MPP) SQL Query Engine and it implements a distributed architecture based on daemon processes. Impala is the best choice when the requirement is a quick result in real time. The intermediate result is stored in In-Memory. Thus, query execution is very fast when compared to other tools.
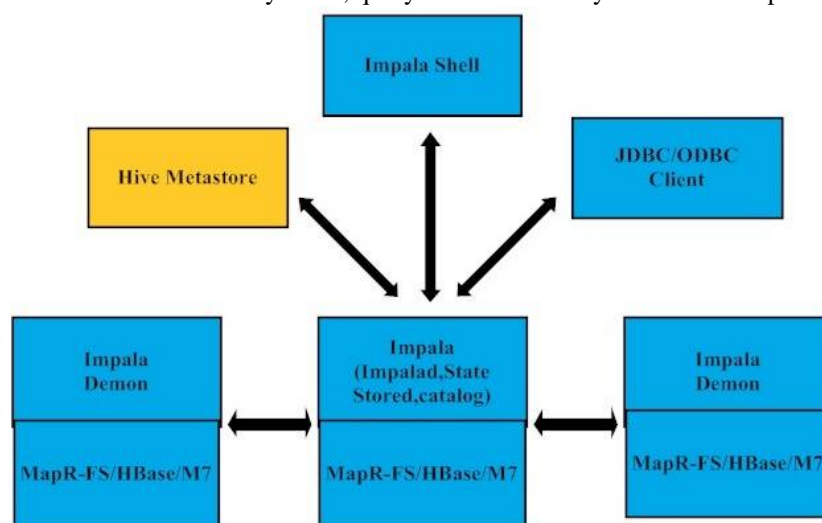


Figure 2: Impala architecture

## II. WORKFLOW OF HIVE AND IMPALA

Cloudera's Impala's performance is more efficient in terms of execution time as well as the complexity of queries. Apache Hive is Infrastructure developed on Hadoop Framework for analyzing and processing data. Basically, Hive is the front end to parse the SQL queries, designs logical plans that execute in the background by MapReduce and Tez this takes comparatively more time than Impala.

Cloudera's Impala does not require data to be moved or transformed as it uses Hive's megastore, it can query Hive's tables directly. Impala does not use MapReduce to execute the query. The Impala daemon executes the process on each node to plan queries and, coordinates between them and query execution engine. As Impala uses parallel processing it responds to query very fast using massively parallel processing. Each node accepts query and planner maps these requests to parallel fragment after that the coordinator starts the execution on the name node of the cluster. The Network systems are highly multithreaded thus each node runs efficiently in a cluster.

Apache Hive generates queries at compile-time whereas Impala does runtime code generation for big loops. MapReduce takes more time for running at full capacity. So, in Hive's query suffers cold start problem on the other hand Impala daemon processes are started at boot time. Hive is fault tolerant whereas Impala does not support fault tolerance if query execution fails Impala needs to start again the execution.

For performance analysis of Apache Hive and Cloudera's Impala, we have used movielens dataset. These are some queries we have executed on movielens dataset. Here we have tables namely ratings and movies.

**Attributes List:**
movies :- movie_id: int, name: chararray, genres: chararray.
ratings :- user_id: int, movie_id: int, ratings: float, timestamp: int.

| **Case 1:** Movies in Particular (1997) year : | |
|---|---|
| The query is to find movie_id and its name released in 1997 from movies table. | **Query:**<br>SELECT DISTINCT movie_id, name<br>>FROM movies<br>>WHERE name LIKE '%1997%'; |

**HIVE:**

**IMPALA:**



| **Case 2:** Particular User Ratings. | |
|---|---|
| The query joins two table ratings and movies on the movie_id field to find movie id, name and its rating for a specific user. | **Query:**<br>SELECT r.movie_id, name, ratings<br>> FROM movie.ratings r<br>> JOIN movie.movies m<br>> ON (r.movie_id=m.movie_id)<br>> WHERE user_id=10564; |

**HIVE:**

**IMPALA:**



| CASE 3: Finding MAX and MIN ratings from 663452198 records. | |
|---|---|
| The query finds maximum and minimum rating from entire ratings table which contain 663452198 records. | **Query:** <br> SELECT MAX(ratings),MIN(ratings) <br> >from ratings; |

**HIVE:**

**IMPALA:**



| Case 4: Movies by Genres | |
|---|---|
| The query is to find the total number of count of the highest rating. | **Query:**<br>SELECT COUNT(ratings)<br>>FROM ratings<br>>WHERE ratings=5; |

**HIVE:**

**IMPALA:**



| Case 5: Pure Comedy Movies With Lowest Ratings. | |
|---|---|
| The query is to find movies having genres as 'Comedy' with the lowest rating. It joins two tables ratings and movies on the movie_id field that displays movie id, name, ratings and genres. | **Query:**<br>>SELECT r.movie_id, name,ratings,genres<br>>FROM movie.ratings r<br>>JOIN movie.movies m<br>>ON (r.movie_id=m.movie_id)<br>>WHERE genres LIKE 'Comedy' AND ratings=0.5; |

**HIVE:**

**IMPALA:**



## III. PERFORMANCE ANALYSIS OF HIVE WITH IMPALA FOR ABOVE CASES

| Impala | 1.18 | 5.11 | 13.20 | 6.35 | 9.08 |
|--------|-------|--------|---------|---------|---------|
| Hive | 48.587 | 97.669 | 102.226 | 137.997 | 122.397 |

## IV. CONCLUSION

Cloudera's Impala has many advantages over Apache Hive. Considering the performance of both, Cloudera's Impala is always preferable to developers when it comes to analyzing HDFS or HBase data because it does not require moving this data, Impala uses Hive's metastore. As Cloudera's Impala is written in C/C++ the advantage is it takes less time to execute the query but it has one demerit that it is not suitable for every file format especially for a file written in java. Although Cloudera's Impala is fast when it comes to upgradation project where compatibility is as important as speed, Apache Hive would nudge. Thus we conclude that Cloudera's Impala's performance is better, time efficient than Apache Hive.

## REFERENCES

[1]. https://www.cloudera.com/documentation/enterprise/5-8-x/topics/impala_file_formats.html
[2]. https://www.linkedin.com/pulse/20140910142911-22744472-why-is-impala-faster-than-hive
[3]. https://www.slideshare.net/hadooparchbook/impala-architecture-presentation
[4]. https://www.dezyre.com/article/impala-vs-hive-difference-between-sql-on-hadoop-components/180
[5]. https://grouplens.org/datasets/movielens
[6]. Big Data: An Introduction, ARD-IJEET, ISSN- 2320-8821, Volume 5, Issue 1
[7]. https://en.wikipedia.org/wiki/Apache_Hive
[8]. https://impala.apache.org/

## BIOGRAPHY

1. Ms. Prasad Mitkari is pursuing his bachelor's degree from Maharashtra Institute of Technology. He is a student of final year, computer science and engineering department. Recently, He completed training on Cloudera's Big Data Course from MIT's Big Data Academy and also doing his final year project in big data area.

2. Ms. Rutuja Banswal is pursuing her bachelor's degree from Maharashtra Institute of Technology. She is a student of final year, computer science and engineering department. Recently, she completed training on Cloudera's Big Data Course from MIT's Big Data Academy and also doing her final year project in big data area.

3. Ms. Akshata Sirsat is pursuing her bachelor's degree from Maharashtra Institute of Technology. She is a student of final year, computer science and engineering department. Recently, she completed training on Cloudera's Big Data Course from MIT's Big Data Academy and also doing her final year project in big data area.

4. Prof. Daivashala Deshmukh is an assistant professor in Maharashtra Institute of Technology. Other than academics she is a coordinator and instructor for Big Data Academy in Computer science and Engineering department, Maharashtra Institute of Technology. She completed training of Cloudera's and HortornWork's Big Data Course. Her current research area of interest is Big Data.