



UTILIZING HADOOP BASED DISTRIBUTED FILE SYSTEM FOR RELATIONAL DATA BASE

Arron Mesi.C¹, Gopika Vijayan², Libin Rose.P³, Muthu.M⁴ and Hema.R⁵

1,2,3,4 Final Year B.E, 5 Assistant Professor,

CSE, Annai Vailankanni College of Engineering,

Email: arronmesi@gmail.com, gopikavijayan94@gmail.com, libinrose10@gmail.com

muthujegatheesh@gmail.com, hema25me@gmail.com

[1] ABSTRACT

Now a days there is an exponential growth of information and data in all the fields such as industrial areas, social networking, business areas, etc.. This paper describes the effective storage methodology to handle big volume of data. Also it deals with the effective utilization of distributed file system for analytical purposes by using specialized tools. HDFS is used to store the large data sets rather than the relational data base. HIVE is used to process those large data sets effectively. The paper can be used for various application area which deals with huge volume of data. The files in the database are stored in HDFS as it is fault tolerant and the file is availed easily for the users. To process the data from HDFS we use a processing technique Hive which is a platform used to develop SQL type scripts to do Map Reduce operations.

Key Words: BigData, HDFS, Hive, Sql, HQL

I. INTRODUCTION

As Big Data requirements shift to the general public, it is important to ensure that the methodologies and techniques which are worked well on small area. In this paper we developed a prototype for utilizing the Hadoop based distributed file system by replacing the relational data storage. It also works well for analytical purposes by retriving data from hdfs by using the tool callled Hive.

To process the data from HDFS we use a processing technique and program model called

as map reduce. In order to make the processing easier queries are written to process the data and store it in HDFS HiveQL is used for writing queries which is equivalent to the program written in map reduce. Hive is a platform used to develop SQL type scripts to do Map Reduce operations. The Hive Query Language (HiveQL) is a query language for Hive to process and analyze structured data in a Metastore. Hive is a data warehouse infrastructure tool to process structured data in Hadoop. It resides on top of Hadoop to summarize Big Data, and makes querying and analyzing easy.

II. HADOOP

Hadoop is an Apache open source framework written in java that allows distributed processing of large datasets across clusters of computers using simple programming models. A Hadoop frame-worked application works in an environment that provides distributed storage and computation across clusters of computers. Hadoop is designed to scale up from single server to thousands of machines, each offering local computation and storage.

1.1 Hadoop Architecture

Hadoop framework includes following four modules:

- **Hadoop Common:** These are Java libraries and utilities required by other Hadoop modules. These libraries provide file system and OS level abstractions and contains the necessary Java files and scripts required to start Hadoop.

- **Hadoop YARN:** This is a framework for job scheduling and cluster resource management.
- **Hadoop Distributed File System (HDFS™):** A distributed file system that provides high-throughput access to application data.
- **Hadoop MapReduce:** This is YARN-based system for parallel processing of large data sets.

1.2 HDFS Architecture

Given below is the architecture of a Hadoop File System.

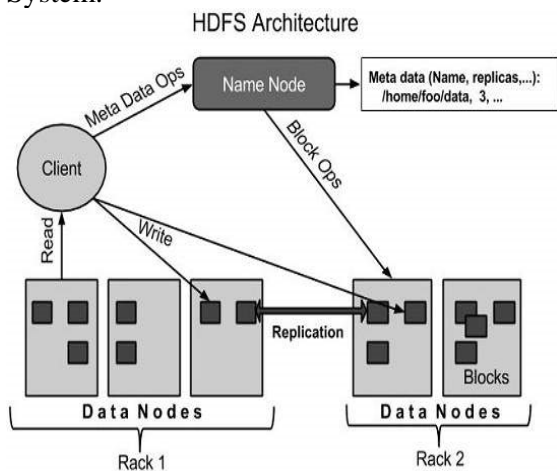


Fig (1.2) – HDFS Architecture

HDFS follows the master-slave architecture and it has the following elements.

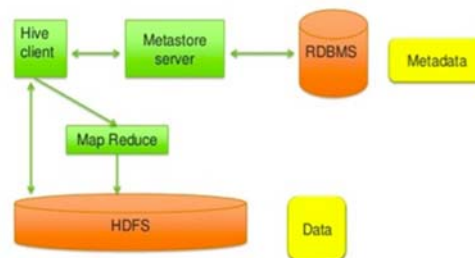
- **Name node**
The name node is the commodity hardware that contains the GNU/Linux operating system and the name node
- **Data node**
The data node is a commodity hardware having the GNU/Linux operating system and data node software. For every node (Commodity hardware/System) in a cluster, there will be a data node.

III . HIVE

Hive is a data warehouse infrastructure tool to process structured data in Hadoop. It resides on top of Hadoop to summarize Big Data, and makes querying and analyzing easy.

Hive architecture

What are we trying to protect here ?



Initially Hive was developed by Face book, later the Apache Software Foundation took it up and developed it further as an open source under the name Apache Hive. It is used by different companies. For example, Amazon uses it in Amazon Elastic Map Reduce.

IV. IMPLEMENTATION

In order to implement the proposed methodology, it requires the configuration of Hadoop and Hive. After the configurations are over, create the table of contents and load the data into the table to the hdfs.

```
student@host17:~$ cd hadoop-2.7.1
student@host17:~/hadoop-2.7.1$ start-dfs.sh
16/03/24 11:42:54 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java class
cable
Starting namenodes on [host17]
host17: namenode running as process 6047. Stop it first.
192.168.1.19: datanode running as process 7896. Stop it first.
Starting secondary namenodes [0.0.0.0]
0.0.0.0: secondarynamenode running as process 6943. Stop it first.
16/03/24 11:43:11 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java class
cable
student@host17:~/hadoop-2.7.1$ start-yarn.sh
starting yarn daemons
resourceanagener running as process 7144. Stop it first.
192.168.1.19: nodenanager running as process 7276. Stop it first.
student@host17:~/hadoop-2.7.1$ jps
9136 Jps
6047 NameNode
7144 ResourceManager
7896 DataNode
7276 NodeManager
6943 SecondaryNameNode
student@host17:~/hadoop-2.7.1$
```

Fig (4.1) Starting Hdfs

```
hadoop@hadoop1:~$ cd /home/student/apache-hive-1.2.1-bin
student@hadoop1:~/apache-hive-1.2.1-bin$ bin/hive
Logging initialized using configuration in jar:file:/home/student/apache-hive-1.2.1-bin/lib/hive-common-1.2.1.jar!/hive-log4j.properties
hive> create database details;
hive> use details;
hive> show databases;
hive> use details;
hive> create table student(reg_no int, name string, age string, sex int);
hive> insert into table student values(901, 'nitesh', 'm', '41');
query id = student_352802d2111a2_cff099c5-9c68-46ff-8b79-15c28604e711
Total jobs = 1
Launching job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Starting job = job_146279438867_0003, tracking url = http://hadoop1:8080/jobs/application_146279438867_0003
hive>
MapReduce job information for Stage-1: number of mappers: 1; number of reducers: 0
2016-03-24 11:11:12,498 Stage-1 map = 0%, reduce = 0%
2016-03-24 11:11:19,124 Stage-1 map = 100%, reduce = 0%, cumulative CPU 2.92 sec
MapReduce Total cumulative CPU time: 2 seconds 938 msec
hive> job = job_146279438867_0003
Job 1 is selected by condition resolver.
Stage-3 is filtered out by condition resolver.
Stage-4 is filtered out by condition resolver.
Loading data for HDFS://192.168.1.15:8020/user/hive/warehouse/details_db/student/
hive> create table student2(reg_no int, name string, age string, sex int);
Logging data to table details.student
Table details.student state: (numFiles=1, totalSize=4, rowDeltaSize=1)
MapReduce job launched.
Stage-Stage1 Map: 1 Cumulative CPU: 2.92 sec HDFS Read: 4024 HDFS Write: 0
3 SUCCESS
Total MapReduce CPU Time Spent: 2 seconds 938 msec
hive>
hive> insert into table student values(902, 'litesh', 'm', '41');
query id = student_352802d2111a2_cff099c5-9c68-46ff-8b79-15c28604e711
Total jobs = 1
Launching job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Starting job = job_146279438867_0004, tracking url = http://hadoop1:8080/jobs/application_146279438867_0004
```

Fig (4.2) Starting HIVE

```
hadoop@hadoop1:~$ cd /home/student/apache-hive-1.2.1-bin
student@hadoop1:~/apache-hive-1.2.1-bin$ bin/hive
Logging initialized using configuration in jar:file:/home/student/apache-hive-1.2.1-bin/lib/hive-common-1.2.1.jar!/hive-log4j.properties
hive> use details;
hive> show tables;
hive> use libhiv;
hive> insert into table student2(name string);
hive> insert into table student2 values('nitesh');
hive> insert into table student2 values('litesh');
hive> insert into table student2 values('nitesh', totalSize=4);
hive> select * from student2;
hive>
hive> use libhiv;
hive> create table student3(name string);
hive> insert into table student3 values('nitesh');
hive> insert into table student3 values('litesh');
hive> insert into table student3 values('nitesh', totalSize=9);
hive> select * from student3;
hive>
hive> use libhiv;
hive> create table student4(name string);
hive> insert into table student4 values('nitesh');
hive> insert into table student4 values('litesh');
hive> insert into table student4 values('nitesh', totalSize=9);
hive> select * from student4;
hive>
hive> use libhiv;
hive> create table student5(name string);
hive> insert into table student5 values('nitesh');
hive> insert into table student5 values('litesh');
hive> insert into table student5 values('nitesh', totalSize=9);
hive> select * from student5;
hive>
```

Fig (4.2) Loading data into the table

V. CONCLUSION

This paper describes the effective storage methodology to handle big volume of data. In this paper we proposed and implemented a prototype for utilizing the hadoop distributed file system for relational data base with analytical tool called Hive. It is best suited for evaluating and analysing the large volume of data.

VI. FUTURE ENHANCEMENT

In our proposed model we have implemented the system with Hadoop and Hive which is a popular SQL interface for batch processing using Hadoop. Until now Mapreduce was the only execution engine in the hadoop ecosystem and hive queries could only run on MapReduce. Now alternative execution engine to mapreduce are available such as Apache Spark.

Spark is relatively new to the Hadoop ecosystem, but its adoption has been meteoric. It is an open-source data analytics cluster computing framework and is built outside of Hadoop's two stage MapReduce paradigm but runs on top of HDFS. So in future we can integrate our proposed with Spark and the data processing can be made even more easier.