

# An Overview - Google File System (GFS) and Hadoop Distributed File System (HDFS)

Snehsudha Popatrao Dhage<sup>1</sup>, Tanpure Renuka Subhash<sup>2</sup>, Rutuja Vilas Kotkar<sup>2</sup>, Prachi Dattatray Varpe<sup>2</sup>, Sonali Sanjay Pardeshi<sup>1</sup>

<sup>1</sup>Loknete Ramdas Patil Dhumal Arts science and commerce college, Rahuri, Ahmednagar, India

<sup>2</sup>PIRENS Institute of Computer Technology Loni (bk), Ahmednagar, India

## Publication Info

### Article history:

Received : 16 February 2020

Accepted : 23 May 2020

### Keywords:

GFS, HDFS, MapReduce, Apache Hadoop, File System

### \*Corresponding author:

Snehsudha Popatrao Dhage

e-mail: snehsudhadhage@gmail.com

com

## Abstract

Each day huge quantity of data is produced by means of text, audio, image, and video format via many sources like social media, YouTube, websites and so on. The most important part is storing this generated data in the file system. To store this intensive data cluster-based storage systems are required. Hadoop Distributed File System (HDFS) and Google File System (GFS) are the two major file system which are used to store the data. This paper gives an overview of GFS and HDFS. Comparison of these file system is also discussed.

## 1. INTRODUCTION

The major problem of many organizations now a days is storing the large amount of data as well as processing it. The search engines are handling data in petabytes per week. Even there are many research and advisement systems datasets are handled by websites like Yahoo [1]. In addition to this there are many ecommerce companies that are handling large amount of data and needs high storage. The You Tube is having lots of videos streaming per week it also needs to store this intensive data. The first requirement of these application is they need cluster-based storage system which are reliable, extremely offered and scalable.

The HDFS and GFS are used to store the huge quantity of data. In 1990's Google File System was developed by Google. There are thousands of storage systems utilized by the GFS built from low-cost service components. These systems offers the many users petabytes of storage for their varied requirements [2]. A key apprehension of the GFS designers was consistency of a system showing to hardware failures, application errors, system software errors, and human errors.

The GFS design important aspects are:

- Reliability and scalability
- File size from GB to TB
- Operations like appending file, random write operation to file.
- Sequential read operations
- Users are not concern about the response time they

process the data in bulk.

As per the analysis some design choices were completed:

- The file is divided or segmented into big chunks
- Automatic file append operation is implemented that permits several applications to append in the same file concurrently.
- The clusters are built around the high bandwidth. The low latency interconnection network is avoided. The control flow is isolated from the data flow. The high band width data flow is planned using the pipelining the data transfer over Transmission Control Protocol (TCP) connections for decreasing the response time. The network topology is exploited by sending data to the neighbouring node in the network.
- Client site caching is eliminated to improve the performance.
- The critical operations are channelized to guarantee consistency through a master monitoring the full system
- Minimize master's involvement in file access operations to avoid hot-spot contention and to ensure scalability.
- Backing effective checkpointing and reckless retrieval mechanisms.
- The effective support is given for garbage collection.

Google has developed MapReduce programming model [3][4] for meeting the fast-growing demands of their web search indene has developed Xing process. GFS

system is used to perform the MapReduce computations [5].

Hadoop is open source framework developed by inspiring from the GFS and MapReduce success [6]. It is used to build large clusters. HDFS used the distributed file system design. HDFS run on service hardware and it is extremely fault tolerant. Hadoop file system earlier developed by the Yahoo but later on it became open source framework. HDFS stores a large amount of data and it offers faster and easier access. The data is stored in the file and the files are stored on many machines. The files are stored in a redundant way to protect the system from data loss in case of failure. Hadoop Distributed File System prepares applications accessible to parallel processing

The features of HDFS are:

- Appropriate for distributed storage and processing.
- To interact with the HDFS Hadoop offers a command interface.
- It has two built in servers called namenode and datanode. These servers assist the user to check the status of the clusters.
- It streams the access to file system data.
- HDFS offers the file permissions and authentication facility.

## 2. GOOGLE FILE SYSTEM

Google File System (GFS) is nothing but it is clusters of computers. The cluster is formed by the network of computers. Every cluster has thousands of machines. There are three main components in each google file system and they are chunk servers, clients, and master servers. Figure 1 shows the general architecture of GFS.

The client makes a file request like retrieve the file, manipulate the existing file, and create a new file. A client can be a computer or it can be a computer application. We say that a client is a customer of the GFS.

The master server works as a coordinator for the cluster. The operational logs are maintained by the master server. The operational logs contain the data of master cluster's activity track. There is minimum service interruption caused due to an operation log. In the event of master server crash, the server that has monitors the operational log is substituted at the place of crashed server. The metadata tracking is done by the master server. The chunks are described by the metadata. The metadata tells the master server that chunk belongs to which file and where in the overall file they fit. All the chunks in the clusters are polled by the master server at the time of startup. The contents of the inventory are sent as reply to master server by the chunk server. Since that moment on, the chunk location within the clusters are tracked by the master servers. At one time, per cluster there is one active master server. There are multiple copies exists for master server to

deal with the failure. The one master server for clusters may result in bottleneck. To overcome this problem GFS keeps very small messages sent and receive by master server.

The chunk servers actually handle the data not master server. They are the workhorses of GFS. system. The 64 MB file chunks are stored in chunk servers. They do not send chunks to the master server. Client request the chunks to the chunk server and chunk servers directly send the requested chunks. Each chunk is copied by the GFS many times and kept on diverse chunk servers. Each copy is known as a replica. GFS by default creates 3 replicas of each chunk. users can do the changes in the setting, and they can make additional replicas as they wanted.

## 3. HADOOP DISTRIBUTED FILE SYSTEM

The HDFS uses the master slave architecture. The datanode, namenode and the block are the main components of the HDFS. Figure 2 shows the architecture of HDFS. The HDFS namespace contains the files and directory hierarchy. The inode represents the files and directories on the namenode by inode, which record qualities such as access time, modifications, permissions, and disk space quotas.

The namenode is nothing but the service hardware. The namenode comprises the GNU or Linux operating system and the name node software. The namenode software can run on service hardware. The system which has namenode is a master server. The master server manages the namespaces of the files. It regulates the file access. The file operations like renaming a file, opening and closing a file or directory are executed by the master server.

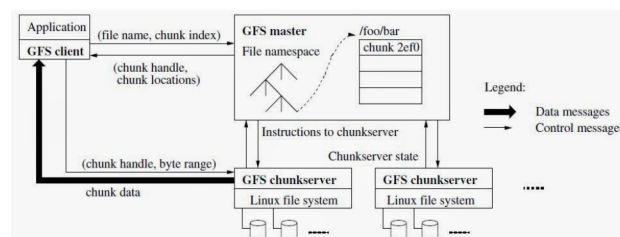


Figure 1: General Architecture of GFS [7]

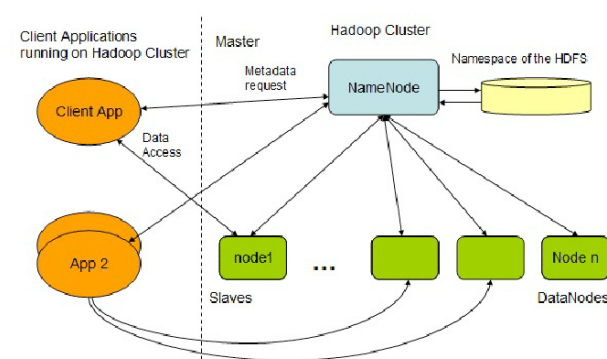


Figure 2: HDFS Architecture [8]

**Table 1:** Comparison between HDFS and GFS

	<i>HDFS</i>	<i>GFS</i>
Platform	Cross Platform	Linux
Developer	Yahoo. now it is an open source framework	Google
Servers	Name node, data node	Master node and chunk server
Hardware used	Commodities hardware	Commodities hardware
Read and write	Write once and reading can be done many times.	It is a multiple read and write model
File deletion	The files deleted are renamed into a specific folder and after that the file will be removed thru garbage	The deleted file is renamed in hidden namespace, it cannot be reclaimed instantly. The file will be deleted after 3days if the file is not in use.
Network Stack Issue	No	Yes
Availability	Journal, edit log	Operational log
Other operations	only append is possible	random file writes possible

The datanode also has the commodity hardware, software and the operating systems GNU or Linux. For each node in the cluster there is a datanode. These nodes accomplish the storage of their system. Datanode as per client's request performs the read and write operations on the file system. Data nodes also performs the operations like creation of block, deletion and replication of the blocks according to the orders of the namenode.

The user data is kept in the file system of the HDFS. The files are divided into one or additional segments and are stored in a separate data node. Segments are known as blocks. The block is the least quantity of data is written and read by the HDFS. The defaulting block size is 64MB. This block size can be changed according to the requirement.

Table 1 shows the contrast between the HDFS and GFS.

#### 4. CONCLUSION

This paper has given an overview of the Hadoop Distributed File System and Google File System as well as these file systems are also compared. The google is using its own File system that is GFS. The HDFS is inspired from the GFS. Both the file systems are using the master slave architecture. The GFS works on the Linux platform on the other hand the HDFS works on the cross platforms. GFS has two servers master node and chunk servers and the HDFS has name node and data node servers.

#### 5. REFERENCES

- [1] GFS vs HDFS. (2020). Retrieved 22 August 2020, from <https://sensaran.wordpress.com/2015/11/24/gfs-vs-hdfs/>
- [2] Bonner, S., Kureshi, I., Brennan, J., & Theodoropoulos, G.(2017). Exploring the Evolution of Big Data Technologies
- [3] Jeffrey Dean., and Sanjay Ghemawat. (2004). MapReduce:Simplified Data Processing on Large Clusters. Google.
- [4] Hadoop - HDFS Overview. (2020). Retrieved from [https://www.tutorialspoint.com/hadoop/hadoop\\_hdfs\\_overview.htm](https://www.tutorialspoint.com/hadoop/hadoop_hdfs_overview.htm)
- [5] Sanjay Ghemawat., Howard Gobioff., and Shun-TakLeung. (2003). The Google File System.Google.
- [6] Hadoop wiki page documentation. Retrieved from <http://wiki.apache.org/hadoop/>
- [7] Richa Pandey. S.P. Sah. (2016). A Review on Google FileSystem. International Journal of Computer Science Trendsand Technology ( Volume 4, Issue 4).
- [8] Anup Suresh Talwalkar. (2020). HadoopT -Breaking the ScalabilityLimits of Hadoop (Doctoral thesis, Rochester Institute of TechnologyRochester, New York).
- [9] Institute of TechnologyRochester, New York). Retrived from [https://www.researchgate.net/publication/265403224\\_HadoopT\\_-Breaking\\_the\\_Scalability\\_Limits\\_of\\_Hadoop](https://www.researchgate.net/publication/265403224_HadoopT_-Breaking_the_Scalability_Limits_of_Hadoop)