

Distributed File Systems: Hadoop Distributed File System and Google File System

Ciprian Lucaci
ciprian.lucaci@tum.de
Technische Universität München, Germany

Daniel Straub
daniel.straub@tum.de
Technische Universität München, Germany

ABSTRACT

Distributed file systems have been a technology enabler to store and process large files which exceed the size of any drive. They also work in a distributed and fail tolerant manner. The first and most known implementations are Google File System and Hadoop Distributed File System.

Keywords

Distributed File Systems, Distributed Systems

1. INTRODUCTION

During the late 90s the Y2K issue was a popular topic in media, but together with the 2000s a bigger and more silent issue took its place. The rise of internet traffic, saving and processing every online footprint, having countless of sensors pouring out data into databases and then the rise of the smartphone to a omnipresence status meant that thousands of terabytes and petabytes of data were being generated daily. This became a problem due to multiple reasons. Firstly, there was no big enough storage medium to store all the data in a single physical location, and secondly processing such amount of data in order to extract valuable information became an even bigger challenge. Even if there were enough space to store huge amounts of data on a single drive, only seeking through such amount of data would be prohibitive in terms of access times. In this context a new IT domain was born - Big Data.

The requirements

Store very large data sets reliably High bandwidth streaming Distribute storage Distribute computation Analysis and transformation of data

The solution Moving computation where data resides Low costs - commodity machines Vertical and Horizontal scalability

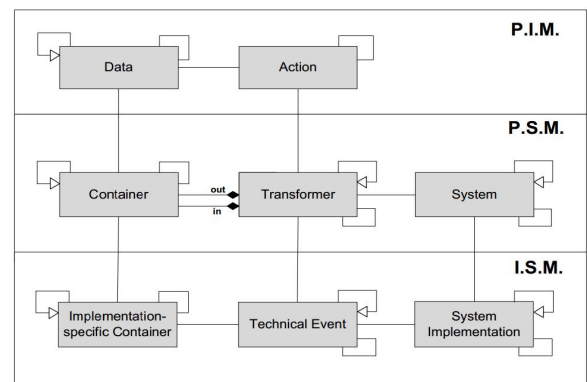


Figure 1: Domain meta-model

Google File System

Hadoop Distributed File System

2. HADOOP DISTRIBUTED FILE SYSTEM (HDFS)

.. add text cite [1]

Hadoop Ecosystem

2.1 Architecture

.. add text

2.1.1 Blocks

.. add text

2.1.2 Data Nodes

.. add text

2.1.3 Name Node

master-worker pattern namespace Inode Metadata Where is the namespace located

2.1.4 Client

.. add text

2.1.5 Secondary Node

- insert picture Is there any weakness in the architecture ..
add text Checkpoint node Backup node

2.2 Workflow

.. add text

2.2.1 Startup

2.2.2 Write

.. add text code sample

2.2.3 Read

.. add text code sample

2.3 Features

.. add text

2.3.1 Block placement policy

2.3.2 Replica management

2.3.3 Balancer

2.3.4 Block scanner

2.4 Purpose

.. add text Optimized for - large files - commodity hardware
- streaming - batch processing - multiple reads Not optimized for - big amount of small files - concurrent modification - arbitrary modification - general purpose applications Cross platform - java, thrift, rest - web access, console - opensource Companies using hdfs - linkedin, amazon, new york times, twitter, ebay, facebook, spotify, ibm, yahoo

3. GOOGLE FILE SYSTEM (GFS)

The Google File System was implemented by Google to meet Google's needs. It is widely used within their company.

3.1 Purpose

Since the concept of the GFS is the basis of HDFS the purpose of GFS is not much different to the purpose of HDFS. It was designed to store large files, meaning 100 MB and more per file. It should run on commodity hardware, therefore the system has to be fault tolerant. Optimized operations are appending to files and sequential reading of files, which are typical operations when streaming files. The goal of GFS is high data throughput. While HDFS is open source and widely used by different companies, GFS is not and is only used by Google itself.

3.2 Architecture

The GFS consists of three different types of components: the chunkservers, the master and the clients (see figure 2). While chunkservers store the data which is available in the system, the master controls the data. The clients can access available data through the master.

The GFS divides files into chunks of a fixed size. The default chunk size is 64 MB. Chunks are replicated to increase fault tolerance and load balancing (see ??). To identify chunks GFS uses a 64 bit long chunk id which is globally unique. [6]

A GFS instance usually consists of many chunkservers which store chunks and one master (see figure 2). A chunkserver knows about its stored chunks and the corresponding chunk ids. In contrast the master knows about all existing files and which chunks are part of a specific file. Furthermore

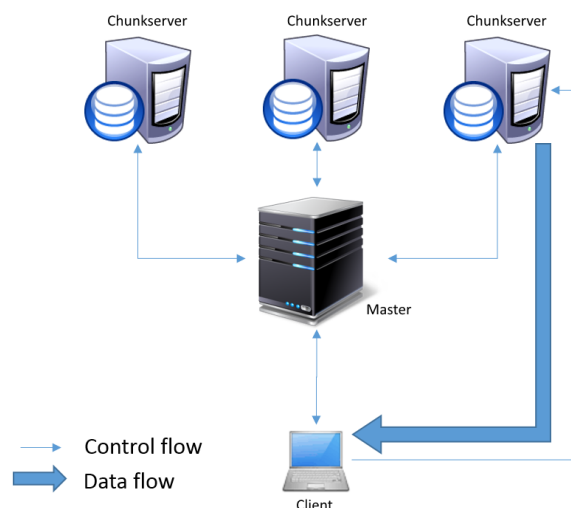


Figure 2: Google File System Architecture

the master knows where the chunks are stored. It also handles access to chunks as explained in 3.3. While the master coordinates access to files, the transmission of data is done directly between a chunkserver and a client. The master is never directly involved in a transmission of data.

3.3 Workflows

The GFS allows different operations. The most important ones are *Read*, *Write* and *Record Append*.

3.3.1 Read

Reading is the simplest of the mentioned operations and will be explained with figure 3.

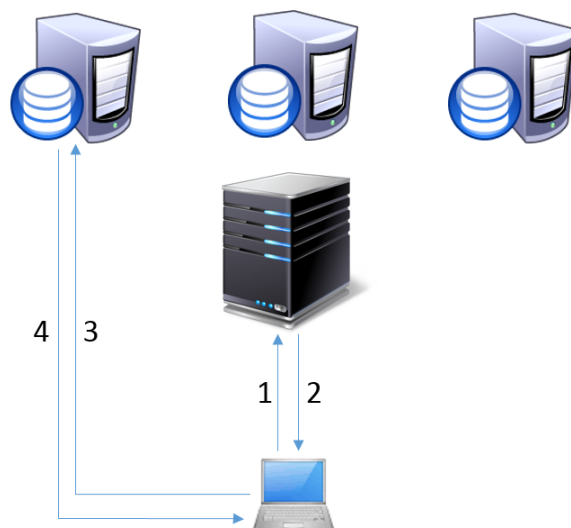


Figure 3: GFS read workflow

To read a file, the client first needs to calculate which part respectively chunk it wants to read. This can be done by the

client, because the fixed size of the chunks is known system wide. The client calculates for a specific file the chunk index which corresponds to the byte offset it wants to read (not shown in the figure). Then the client sends a request to the master, containing the name and previously calculated chunk index of the requested file (step 1 in the figure). The master responds with the unique chunk id (see 2) of the requested chunk and its positions within the system, i.e. the chunkservers containing the requested chunks (step 2). In the 3rd step the client chooses one of the chunkservers and sends a request containing the chunk id and the byte range within that chunk. Finally the chunkserver sends the requested data to the client.

3.3.2 Write

Writing to a file is a bit more complex than the *Read* operation. The *Write* operation can be defined as modifying a specific byte range within a file.

Leases

GFS uses the mechanism of leases for writing to files. This is done to ensure a consistent mutation order for all replicas of a chunk. A chunkserver may have leases for chunks, but only one chunkserver may have the lease for a specific chunk at any time. The chunkserver which has the lease for a chunk is called *primary*, the others are called *secondaries*. Leases can be requested by chunkservers and are granted by the master.

The following steps explain the execution of a *Write* operation.

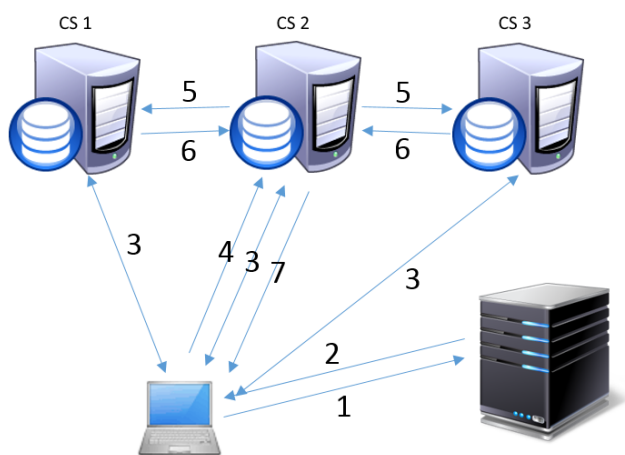


Figure 4: GFS write workflow

Assume the following scenario with reference to figure 4:

There are three chunkservers (CS), all of them containing the chunk x where x is the unique identifier of the chunk. CS2 is the primary. The master knows that there is a specific file (called file1) which is divided into different chunks with the first chunk being chunk x. It also knows that chunk x is available on the mentioned CSs CS1, CS2 and CS3. The client wants to write to the first byte of file1.

1. Like in the read operation the client's first step is to calculate the chunk index it wants to write to. With the intention to write the first byte of file x the client calculates the chunk index which is 0 (the first chunk of the file). Now the client sends a request to the master containing file name and chunk index (file1, 0).
2. The master replies with the chunk id, the primary and the secondaries (x, CS2, CS1, CS3). If at this point no CS has the lease for requested chunk, the master grants the lease to one of the chunks before answering the clients request.
3. With that information the client pushes the data which it wants to write to all replicas (the primary and the secondaries). All receiving CSs write the data to its buffer cache and respond with an acknowledgement for receiving the data.
4. After all CSs have acknowledged receiving the data, the client sends a write request to the primary. At this point the primary applies the requested write to itself.
5. Then the primary forwards the write request to all secondaries leading them to apply the changes in the same order.
6. On successful writing the data the secondaries will indicate the completion of the write with a response to the primary.
7. At last the primary returns the result of the write to the client. If there was any error, the client will retry the write from step 3 (pushing the data to CSs). If the error persists, the client will restart from step 1.

3.3.3 Record Append

With *Record Append* the GFS provides an atomic write operation. When a client uses this operation to change a file, the client transmits only the data it wants to write, not the position. The CS then decides where the data is written and returns the offset to the client.

The workflow for record append is similar to the write operation explained in 3.3.2. The differences are in the following steps:

3. The client generally pushes the data to the last chunk of the file.
4. Instead of requesting a write, the client sends a record appen request.
5. Before the primary applies the changes it checks if the data fits in the chunk (see below). If it fits the following steps are the same as for the write operation.

Something different happens, if the primary's check fails:

If a client wants to append data to a file and the last chunk does not provide enough free space the primary does not apply the changes.

0100 0100 0110 1101

Instead the chunk is padded (e.g. filled up with zeros) and a new chunk is created for the file. The primary then responds to the client with an error asking it to retry the request with the new created chunk.

0100 0100 0110 1101 0000 0000

The client will retry, now pushing the data to the new chunk and will succeed.

0010 0001 1000

Since the chunk size is limited the record append is technically limited to data with a size of 64 MB or less. GFS further restricts the record append operation to data with no more than one forth of the chunk size. This decision was made to reduce the fragmentation which is created by unsuccessful record append operations.

3.4 Features

GFS has a few features implemented to enhance performance and fault tolerance. Chunk replication and garbage collection are two of them which will be explained in the following sections.

3.4.1 Chunk Replication

As mentioned earlier, chunks in the GFS are replicated. By default there are three replicas of each chunk.

Chunk Creation: When a new chunk is created the master decides where the replicas are placed. This decision is affected by three factors. First new replicas are preferably placed on chunkservers with the most free space. This is done to fill up newly added chunkservers and to have a balanced usage of disk space. Secondly the number creations per chunkserver at any time should not be too high. This avoids flooding new chunkservers with requests. And third the replicas for each chunk should not be on the same rack of chunkservers. This enhances the fault tolerance of the system when whole racks go down e.g. due to network problems.

Chunk Re-replication: As soon as the number of replicas for a chunk is lower than specified the master will re-replicate it. When re-replicating a chunk the same factors are considered as during the chunk creation. Additionally the re-replication of chunks is prioritized: chunks which have lost more replicas are prioritized to others. Deleted files (see 3.4.2) get a lower priority. And at last files currently in use get a higher priority.

Chunk Rebalancing: The master periodically checks the

distribution of replicas. Depending on chunkserver load and disk space usage the master moves replicas between chunkservers to improve the load and disk space usage.

3.4.2 Garbage Collection

When a file is deleted in the GFS the master logs the deletion and renames the file, creating a "hidden" file. The deletion time stamp is included in the file name. The garbage collection runs periodically and searches for files which have been deleted at least three days (default value) ago. Finding such a file results in the final deletion of the file. (figure 5)

Until the deletion by the garbage collection the file can still be accessed (using the new file name) or undeleted.



Figure 5: GFS garbage collection

4. COMPARISON

.. add text

5. CONCLUSIONS AND FUTURE WORK

... not finished ...

This paragraph will end the body of this sample document. Remember that you might still have Acknowledgments or Appendices; brief samples of these follow.

6. REFERENCES

- [1] Shvachko, K. and Hairong Kuang and Radia, S. and Chansler, R.: The Hadoop Distributed File System, in *Mass Storage Systems and Technologies (MSST)*, 2010 *IEEE 26th Symposium*, pp.1-10
- [2] Hadoop: What it is and why it matters, online at http://www.sas.com/en_us/insights/big-data/hadoop.html, [accessed: May 2015]
- [3] Hadoop tutorial, online at <http://www.bigdataplanet.info/2013/10/hadoop-tutorials-part-1-what-is-hadoop.html>, [accessed: May 2015]
- [4] Thomas Kiencke: Hadoop Distributed File System, *Institute of Telematics, University of Lubeck, Germany*, online at <https://media.itm.uni-luebeck.de/teaching/ws2012/sem-sse/thomas-kiencke-hdfs-ausarbeitung.pdf>, [accessed: May 2015]
- [5] R.Vijayakumari, R.Kirankumar, K.Gangadhara Rao: Comparative analysis of Google File System and Hadoop Distributed File System, in *International Journal of Advanced Trends in Computer Science and Engineering*, Vol.3 , No.1, pp.: 553-558, 2014
- [6] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung: The Google File System, in *Proceedings of the Nineteenth ACM Symposium on Operating Systems Principles, 2003*, pp.29-43