# Survey on a Google File System (GFS)

Saleh A. Khawatreh
Dept. of Computer Engineering
Faculty of Engineering, Al-Ahliyya Amman
University
Amman-Jordan

Enas N. Al-Zubi
Dept. of Computer Engineering
Faculty of Engineering, Al-Ahliyya Amman
University
Amman, Jordan

## ABSTRACT

The Google File System (GFS), a proprietary scalable distributed file system sophisticated by Google to be used for its huge distributed data applications, primarily the search engine. It improves efficiency, reliability, scalability, transparency, security and fault tolerance despite of using cheap commodity computers and serving a large number of clients. GFS is similar with the previous distributed file systems in the goals, it divides the files into chunks of data each one is 64 MB in size and generate several copies of each data chunk. Then mounts these data chunks on several servers which could be based on global geographical locations, and that will enhance the reliability of retrieving files online and overcome any limited access to one or more servers. It will also enhance the concurrency access and control over files online due to the fact that several copies of each file do exist at the same time. GFS has successfully met Google's storage needs, for generation and data processing such as research and development issues that need huge data sets. The biggest cluster provides hundreds of TBs of storage distributed in thousands of disks on thousand machines, and it provides concurrency access by hundreds of users.

## Keywords

Destributed File Systems, Google File System (GFS) clustered storage, data storage, performance, fault tolerance.

## 1. INTRODUCTION

To meet the fast increasing applications of Google's information processing demands, the Google File System (GFS) have been designed. GFS is developed from Google effort " Big Files" which was presented by Larry Page and Sergey Brin in the early time of Google, it has the same goals of the previews systems including performance, scalability, availability, security, transparency, openness, and fault tolerance template. However, key observations of Google application workloads and new technologies have improved the GFS design and reexamined the previous ones and introduced different points in the design ways.

GFS divides the files into fixed-size chunks, each one is 64 MB, as in clusters and sectors in regular file systems, that are rarely overwritten or shrunk, but they are usually read or appended to. It is designed to work on Google's clusters, that consists of cheap commodities, so precautions must be taken against the high rates of failure proposed by the nodes itself or the subsequent data loss. When the high data throughput is needed, other design concepts are proposed and here are:

First, multiple nodes are consisted the GFS cluster, nodes in the cluster are connected as centralized style where there is one master node and a large number of chunk servers. Because of the large number of nodes which reaches to hundreds or even thousands of storage nodes that consist of cheap commodity components, which make the failure is norm not an exception, these problems occurred by application or operating system bugs, user errors, failures of connections, memory, disks, network, or power supplies. So, error detections, periodic monitoring, fault tolerance limits, and self-recovery must be taken part in the system.

Second, traditional files are huge in size, generally multi-GBs. Every file contains application objects such as web documents. Working with high increasing data sets consisting of many TBs including billions of objects, it is difficult to manage billions of KB-sized files even if the file system can support it. So, parameters and design assumptions have to be adopted.

Third, appending new data is often used rather than overwriting existing one. When written, the files are read only and only sequentially. Different types of data sharing, some constitute huge repositories which data analysis processes scan though, others may be data streams operated in continuous manner by processing applications, some may be archived information, some may be temporal results presented on one machine and to be continued on another one, however, in the same time or later. Working such as access ways on large files, appending will be the basic of performance optimizations and success guarantees, while forwarding data blocks in the client lose its appeal.

Fourth, participating the applications and the file system improves the whole system by increasing the flexibility, such as, relaxed GFS's consistency model to simplify the system without forcing an exhausted load on the applications. Also, introducing an atomic suffix process so that many users can append to a file concurrently without any synchronization among them, more details are presented later in the paper.

Many GFS clusters now are implemented for various purposes. The biggest ones have more than 1000 storage nodes, more than 300 TB of disk storage, and accessed by hundreds of clients on different machines.

## 2. DESIGN OVERVIEW

*A. Assumptions*

First, to design a file system, assumptions must be guided to offer opportunities and challenges, and these are:

The primary components that used in the system are cheap commodities which often fail. It must evolve continuous monitor and detect, tolerate, and recover from failures as a routine manner.

Storing a modest number of large files in the system, approximately few million files, each file with 100 MB or more. Multi-GB files are the basic issue and must be managed in efficient way. Small files should be supported.

Two kinds of reads are implemented in the workloads: large streaming reads and small random reads. The large streaming reads provide individual operations that read hundreds of KBs, commonly 1 MB or more. Successive processes from the same user usually read from a contiguous area of a file. A small random read often reads a few KBs at some specific offset.

Implementing well-defined semantics for multi-clients concurrently append to the same file. These files also, called

producer-consumer queues or of more ways merging. Hundreds of producers, working one in each machine, and need to append to a file. Atomicity with minimum synchronization overhead is primary. The file may be read after some time, or a consumer may be reading among the file in the same time.

The high bandwidth is more essential than low latency the applications mostly focus on processing data in bulk at a high rate, but few have stringent time to response for every read or write.

### B. Maintaining the Integrity of the Specifications

The template is used to format your paper and style the text. All margins, column widths, line spaces, and text fonts are prescribed; please do not alter them. You may note peculiarities. For example, the head margin in this template measures proportionally more than is customary. This measurement and others are deliberate, using specifications that anticipate your paper as one part of the entire proceedings, and not as an independent document. Please do not revise any of the current designations.

### C. Interface

The interface in GFS is general, files are ordered in hierarchal way and known by pathnames. Ordinary operations such as: create, open, delete, close, write, and read.

### D. Architecture

When talking about architecture, there are three types which are distributed, centralized, and hybrid. In GFS, the system is consisting from clusters, and each cluster consists of one master node and multiple chunk servers typically a cheap commodity Linux machines operating as a user-level server process, and is accessed by many users, as shown in Fig. 1. It is easy to process a chunk server and a client on that machine.

GFS divides the files into fixed-size chunks of 64 MB, every chunk has a unique global identifier chunk handle consisting of 64 bit and is assigned by the master node in the creation period, chunk servers keep chunks in local disks as Linux files, and every read or write is specified by the chunk handle.

The master stores the system metadata, and includes the namespace the translation from files to chunks, access control, and the actual locations of chunks. It also manages system-wide works such as chunk lease controlling, garbage collection of unused chunks, and chunk mixing between chunk servers. Heart Beat messages are sent from the master to the chunk servers to manage and control it and know its states.

GFS user code linked into each application presents the file system API and communicates with the master and the chunk servers, so read or write on data is occurred. Users take metadata from the master, and communicate with chunk servers for data. GFS do not implement the POSIX API and the Linux vnode layer.

Caching file data are not used in the client nor in the chunk serve. Client caches do not offer a huge improvement because the most applications stream use huge files and have working sets very large for caching.

### D. Single Master

Master manages and sophisticates chunk placement and replication orders by global information. However, it is good to minimize the master responsibilities such as reads and writes to avoid the bottleneck, clients do not read or write through the master, although, the client communicates with the master to know which chunk server it must contact. This information is

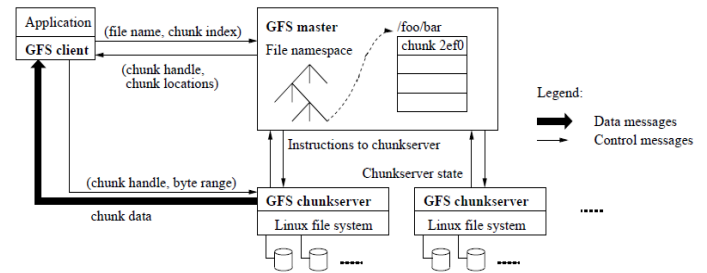cached for a time and transfers to the chunk servers directly to perform the operation.



**Fig.1 GFS architecture**.

According to Fig. 1, the simple read is as following: first, application originates the read request, then GFS client translates the request and send it to master. Master gives the chunk handle and replica location, after that, client communicates with a location and sends the request, chunk server send the data to the client, and finally client forwards the requested data to the application. In writes, first, application starts the request, and sends it to GFS client to transfer the request and sends it to the master, master responds with chunk handle and replica locations, then the client pushes the new data to all locations, this data is stored in the buffers, and when the client sends write command to primary, then primary checks serial order for data in the buffers and writes the data in that order to the chunks, then primary sends the serial order to all the secondary nodes and ask them to write, when secondary nodes finish the writes, they respond back to primary, and the primary acknowledges the client.

### E. Chunk Size

GFS choose 64 MB that bigger than ordinary file system block sizes and this is one of the design parameters. Every chunk replica is kept as a plain Linux file in a chunk server and take place as needed. Lazy space allocation is good to avoid wasting space because of the internal fragmentation.

This large chunk size has many advantages, one of them is, it decreases clients' interactions with the master because reads and writes to a chunk needs only one starting request to the master to ask for chunk location. And this decrease is important for this workload because applications usually read and write huge files in sequential way. Also, small random reads, the client cache the chunk locations for a multi-TB working methods. Another advantage is that, the client is better performing multi-operations in a specific chunk, and reduce network overhead by using TCP connection to the chunk server in a limited period of time. Finally, it decreases the size of the metadata in the master, and that provides more spaces in the memory as it is shown in the next section.

Although, there are several advantages for this large size chunk, there are some disadvantages, such as when the file is small, it needs a small number of chunks may be only one, and when many clients ask for this junk, a hot spot is happened.

GFS solves this problem by keeping executable chunks that have a higher replication factor, and by operating the batch queue system stagger application start times. A likely solution is by making clients read data from other clients in some cases.

### F. Metadata

Three types of metadata are stored in the memory of the master and these are: the file and chunk namespaces, the translation

from the files to chunks, and the locations of chunk's replicas. The two types are stored permanently through logging mutations to an operation log kept in the local disk on the master and replicated on other machines. The advantage of using the log is to update the master state without risking inconsistencies of the master crash. Chunk locations are not stored persistently, although, it requests its chunk in all chunk servers at the setup or when a new chunk server is added to the cluster.

Since metadata is stored in memory, expert operations are quick. Besides, it is simple and effective for the expert to occasionally look over its whole state out of sight. This intermittent examining is utilized to actualize chunk garbage accumulation, re-replication in the vicinity of chunk server failures, what's more, chunk movement to adjust load and disk space utilization crosswise over chunk servers.

One potential sympathy toward this memory-just approach is that the number of chunks and thus the limit of the entire framework is restricted by the amount of memory the expert has. This is not a genuine constraint. The master maintains under 64 bytes of metadata for each 64 MB chunk. Most chunks are full on the grounds that most records contain numerous chunks, just the remainder of which may be somewhat filled. Likewise, the file namespace information ordinarily requires less than 64 bytes for each file because it stores document names utilizing prefix pressure.

The master does not store a constant record of which chunk servers have a reproduction of a given chunk. It essentially surveys chunk servers for that data at startup. The master can stay up with the latest from that point that it controls all chunk placement and screens chunk server status with customary Pulse messages named "Heart Beat".

The authors, at first endeavored to keep chunk location data diligently at the master, yet they concluded that it was much easier to demand the information from chunk servers at startup, what's more, intermittently from that point. This dispensed with the issue of keeping the master and chunk servers in a state of harmony as chunk servers join and leave the cluster, change names, come up short, restart, and so on. In a cluster with hundreds of servers, these occasions happen very frequently.

Another approach to comprehend this design decision is to figure it out that a chunk server has the last word over what chunks it does or does not have on its own disks. There is no point in attempting to keep up a perspective view of this data

on the master due to the fact that mistakes on a chunk server may bring about chunks to vanish suddenly (e.g., a disk might have crashed also, be crippled) or an administrator might rename a chunk server.

The operation log contains a verifiable record of basic metadata changes. It is integral to GFS. Not just is it the just tenacious record of metadata, yet it additionally serves as a time stamp of the sequence of operations. Files and chucks and in addition their variants are all interestingly and unceasingly recognized by the sensible times at which they were made.

Due to the critical factor of the operation log, it must be stored reliably and not allow changes to be visible to clients till metadata changes are be continuous. On the other hand, effectively loose the all file system or modern client operations although the chunks are alive. Thus, replicate it to many remote machines and serve client's processes only after flushing the symmetrical log record to disk both locally and remotely. The master impresses many log records together before flushing

thereby reducing the impact of flushing and replication on all the system.

Through replaying the operation log, the master returns its file system state. Log must be kept small to decrease the startup time. To be sure that the log does not increase, the master must check its state so it will restart by loading the new checkpoint from the local disk and only replaying a specific range of records. This checkpoint is in the form of B-tree and could mapped directly to memory and used in namespace search without additional parsing. That can improve the speed of restart and the availability.

Checkpoints take some time, so the master's internal state is operated in a manner that a latest checkpoint could be made without extra delay. The master switches to a new log file and makes the latest checkpoint in an individual request. These checkpoints contain the whole mutations before this switch. It could be made in a less than minute or for a cluster that has a million files. When finished, it is stored to disk in local and remote.

### G. Consistency Model

GFS has a consistency model which develops distributed applications and keeps the simple and efficient implementations. In this part, GFS guarantees are discussed and how GFS keeps these guarantees.

File namespace changes (e.g., file creation) are atomic. They are made by the master: namespace locking guarantees atomicity and correctness, the master's operation log provides a global order of these operations.

The type of the change is responsible of the state of a file region after data change. Never the less, it succeeds or fails, or there are continuous changes. Table 1 show the result.

A file location is reliable if all clients will continuously see the same information, regardless to which replicas they read from. A location is characterized after a file information change in the event that it is reliable and customers will see what the change writes in its whole. At the point when a change succeeds without impedance from simultaneous access, the influenced location is characterized (and by suggestion reliable): all clients will dependably see what the transformation has composed. Simultaneous fruitful transformations leave the locale location however steady: all clients see the same information; however, it may not reflect what any one change has composed. Regularly, it comprises of blended parts from different transformations. A fizzled transformation makes the location conflicting (subsequently likewise unclear): different clients might see diverse information at distinctive times. The authors depict underneath how the applications can recognize characterized locales from indistinct locations.

GFS applications may adjust the relaxed consistency model with simple methods needed already for some purposes: depending on attaches rather than overwrites, check pointing, and writing self-validating, self-identifying records.

In practice, all the applications turn into files by affixing rather than overwriting. In one ordinary use, a writer produces

a file from the start to the end. It basically changes the name of the file to a permanent name after writing all the data, or checkpoints in every specific period of time how much has been successfully written. Checkpoints may also have application-level checksums. Readers verify and process only the file region up to the last checkpoint, which is known to be in the defined state. Regardless of consistency and concurrency issues, this approach has served us well. Appending is far more efficient and

more resilient to application failures than random writes. Check pointing allows writers to restart incrementally and keeps readers from processing successfully written file data that is still incomplete from the application's perspective.

**Table 1. File Region State After Mutation**

| | Write | Record Appen |
|---|---|---|
| Serial success | *defined* | *defined* interspersed w |
| Concurrent successes | *consistent* but *undefined* | *inconsistent* |
| Failure | *inconsistent* | |

## 3. SYSTEM INTERACTIONS

The system was designed to decrease the master's responsibilities. Here are how client, master, and chunk servers communicate with each other to perform data changes, snapshot, and atomic record append.

### E. Leases and Mutation Order

A mutation is a process which changes the contents or metadata of a chunk like a write or an append process. Every mutation is implemented at all the chunk' replicas. Leases are used to maintain a continuous mutation order across replicas. The master awards a chunk lease to a replica, that called the primary. The primary chooses a serial order for all mutations to chunks. All replicas do the same as that order applying mutations. So, the global mutation order is created at the start through the lease award order that picked by the master, and with a lease by the serial identifiers assigned by the primary.

The lease technique is created to eliminate management overhead at the master. A lease has an primary timeout of 60 seconds. Although, thus, the chunk is being mutated, the primary can ask and typically get domains from the master indefinitely. These domains requests and grants are piggybacked on the Heart Beat messages in regular way

exchanged between the master and all chunk servers. The master can sometimes try to drag a lease before it endears (e.g., when the master wants to weaken mutations on a file that was renamed). Although, if the master loses interactions with a primary, it can safely allow a new lease to another replica after the old lease dies. Figure 2, shows this process through the control flow of a write by these steps:
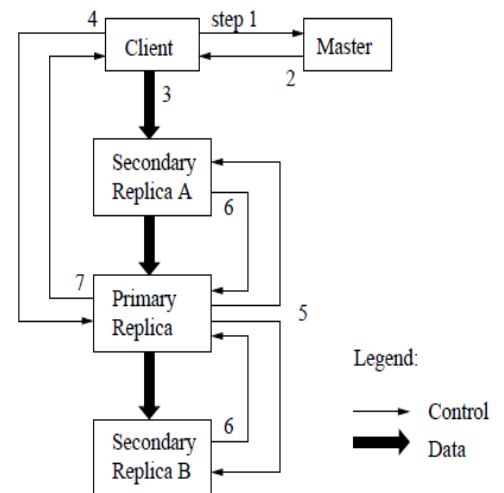
1. The client requests the master that chunkserver keeps the present lease for the chunk and the locations of the all replicas. If no one has a lease, the masterpicks one to a replica it chooses.
2. The master responds with the identity of the primary and all locations of the secondary replicas. The client stores these data to next mutations. It needs to communicates with the master next time when only the primary can not be reached or replies which it no longer has a lease.
3. The client caches the data to the replicas. A client may do so in different order. Every chunk server will keep this data in an internal LRU buffer cache till the data is in use or dropped out. Through separating the data sent from the control flow, so, performance will be improved by ordering the important data flow according to the network regardless of which chunk server is the primary.

4. When all replicas receive the acknowledgment, the client requests write to the primary. The request knows the data previous cached to all replicas. The primary specifies s serial numbers to the all mutations it receives, may be from multiple clients, that supplies the primary serialization. It supplies the change to its own local state in sequence.
5. The primary sends the write requests to all secondary replicas. Every secondary replica supplies changes in the same sequential number assigned by the primary.
6. The secondary nodes reply to the primary signalizing that they have finished the process.
7. The primary sends back to the client. Any errors met at any of replicas are reported to the client. If there are errors, the write may have passed at the primary and a subset of the secondary replicas. (If it had not passed at the primary, it would not have been created a serial number and send it.) The client request is assumed to have dropped, and the modified location is left in not continuous state. The client code treats these errors through trying the failed mutation again. It will make a few tries at steps (3) to (7) before falling back to a retry from the starting of the write.



**Fig. 2. Write Control and Data Flow**

### F. Data Flow

The flow of data is different from the flow of control to use the network maximum efficiency. But control flows from client to primary and then to secondary nodes, data is cached along choose sequence of chunk servers in a pipelined model. The goals are to utilize fully every machine's network bandwidth, against network bottleneck and huge-latency links, and decrease the latency to cache by all the data. To utilize every machine's network bandwidth, the data is cached along a series of chunk servers not as distributed in some other topology (e.g., tree). So, every machine's full outbound bandwidth is used to transfer data quickly as can rather than divided between multiple recipients.

### G. Atomic Record Appends

GFS applies an atomic approach process called record append. In an ordinary write, the client chooses a specific offset in that data is to be written. Continuous writes to the same locations are not serialize: the location may end up consisting of data fragments from many clients. In a record append, although, the client specifies only the data. GFS appends it to the file once at least atomically (i.e., as one continuous serialize order of bytes) at an offset of GFS's choosing and returns that offset to the client. This is the same as to writing to a file opened in O APPEND mode in UNIX without the strain conditions when multiple writers do so concurrently. Record append is mostly used by distributed applications in that multi clients get the same file concurrently.

## 4. MASTER OPERATION

The master operates all namespace processes, and controls chunk replicas throughout the system: it chooses placement decisions, makes the latest chunks and hence replicas, and coordinates many system-wide activities to keep chunks replicated, to poises load across all the chunk servers, and to reused the unused storage. Now these topics are discussed below:

### A. Namespace Management and Locking

Several masters processing require a long time: for example, a snapshot process has to change chunk server leases on all chunks covered by the snapshot. Delay another processes to be active and use locks over locations of the namespace to ensure the best serialization. Not as many traditional file systems, GFS does not use a per-directory data structure which lists all the files in which directory. Nor does it provide aliases for the same file or directory (i.e., hard or nominal links in UNIX terms). GFS performs its namespace as a lookup table mapping full pathnames to metadata. With prefix compressions, this table prepared in memory. Each part in the namespace tree (either a free file name or a standard directory name) has a read-write lock. Every master operation acquires a number of locks before it runs.

### B. Replica Placement

A GFS cluster is almost distributed at deferent levels from each other's. It almost has thousands of chunk servers manipulated across multi machines. These chunk servers may be accessed from thousands of clients from the same or different machines. Communication among two machines on different formats may access one or more network switches. In addition, bandwidth into a rack may be less than the normal bandwidth of all machines with the rack. Multi-level distribution presents a unique challenge to distribute data for, reliability, availability, and scalability.

## 5. FAULT TOLERANCE AND DIAGNOSIS

One of the hardest challenges in implementing the system is operating with frequent component crashes, The quality and quantity of nodes together make these problems more than the exception: this leads to cannot completely trust the machines, and cannot trust the disks too. Component crashes can result in an unavailable system or, worse, corrupted data.

### A. High Availability

Keeping the system highly available with two simple effective methods: fast recovery, replication and master replication.

### B. Data Integrity

Every chunk server makes check summing to detect corruption of received data. As that a GFS cluster usually has hundreds of disks on hundreds of machines, it often organizes disk failures that cause data corruption or loss on both the read and write operations. They could recover from corruption using another chunk replica, but, it could be not practical to detect corruption by comparing replicas across chunk servers. However, divergent replicas may be legal: the semantics of GFS mutations, in particular atomic record append, does not guarantee identical replicas. So, every chunk server must verify in independence the integrity of its own replica by maintaining checksums.

### C. Diagnostic Tools

Extra and detailed diagnostic logging has helped immeasurably in such problems: isolation, debugging, and performance evaluations, while using only a little cost. Without logs, it is difficult to understand transient, non-repeatable communications between machines. GFS servers create diagnostic logs which record many significant rolls (for example chunk servers going up and down) and all RPC requests and replies. Those diagnostic logs could be free to delete without affecting the correctness of the system. Although, trying to keep these logs close as far as space permits.

## 6. MEASUREMENTS

Many micro-benchmarks are used to measure the bottlenecks inherent in the GFS systems and architectures, and also many measurements from real clusters that used in Google.

### A. Micro-benchmarks

Measuring performance on a GFS cluster consist of one master, two master replicas, 16 chunk servers, and16 clients. This configuration was set up to be easy to use. Ordinary clusters have thousands of chunk servers and hundreds of clients. All the machines are configured with dual 1.4 GHz processors, 2 GB of memory, two 80 GB 5400 rpm disks, and a 100 Mbps full-duplex Ethernet connection to an HP 2524 switch. All 19 GFS server machines are interacted with one switch, and all 16 client machines to the other. The two switches are connected with a 1 Gbps link.

### 1. Reads

Clients read at the same time from the file system. Every client reads a in a random way selected 4 MB location from a 320 GB file set. This is repeated 256 times so every client finishes reading 1 GB of data. The chunk servers taken together have only 32 GB of memory, thus improving at most a 10% hit rate in the Linux buffer cache. These results would be close to actual cache results, and Fib.3 (a) shows the results.

### 2. Writes

Clients write at the same time to N distinct files. Every client writes 1 GB of data to another file in a sequence of 1 MB writes. The write rate and its theoretical ranges are presented in Figure 3(b). The limit at 67 MB/s because of the need to write every byte to 3 of the 16 chunk servers, every one with a 12.5 MB/s input connection. The write rate for one client is 6.3 MB/s, half of the limit. The idea for that is that networks acknowledgments. It does not interact well with the pipelining model so, using for storing data to chunk replicas. Delays in propagating data from a replica to another decrease the overall write rate. Aggregate write rate is 35MB/s for 16 clients (or2.2 MB/s per client), is half the theoretical limit. As in the state of reads, it is better that multiple clients write simultaneously to the same chunk server as the number of clients increases. In addition, collision is better for 16writers than for 16 readers because of the three different replicas from every write involves.

### 3. Record Appends

Showing record appends performance in Fig. 3(c). Clients append simultaneously to a single file. Performance is decreased

by the network bandwidth of the chunk servers that keep the end chunk of the file, independently of the number of clients. It starts at 6.0 MB/s for one client and falls to 4.8 MB/s for 16 clients, mostly because of the congestion and variances in network transfer rates caused by the clients.

### B. Real World Clusters

Examining two clusters using in Google that are representative of many others like that. Cluster A is used usually for research and improvement by hundred engineers. An ordinary task is identified by a user and processing to several hours. It reads by many MBs to a few TBs of data, exchanges or analyzes data, and sends the results back to the cluster. Cluster B is used for production data processing. Table 2 Shows that, the tasks last longer and continuously generate and process multi-TB data sets with only occasional human intervention. In both cases, a single "task" consists of various processes on many machines read and write a number of files at the same time. Table 3. Shows the performance metric for them, and Fig, 3. Shows the throughput
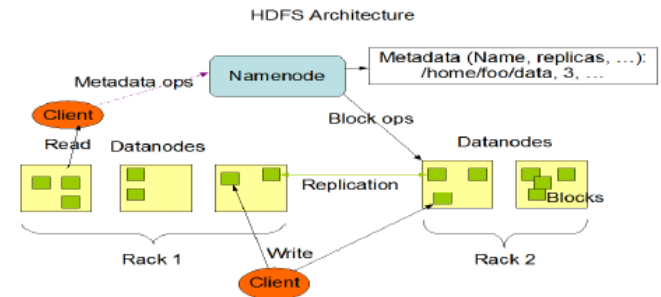
**Table 2. Characteristics of two GFS clusters**

| Cluster | A | B |
|---|---|---|
| Chunkservers | 342 | 227 |
| Available disk space | 72 TB | 180 TB |
| Used disk space | 55 TB | 155 TB |
| Number of Files | 735 k | 737 k |
| Number of Dead files | 22 k | 232 k |
| Number of Chunks | 992 k | 1550 k |
| Metadata at chunkservers | 13 GB | 21 GB |
| Metadata at master | 48 MB | 60 MB |

**Table 3: Performance Metrics for Two GFS Clusters**

| Cluster | A | B |
|---|---|---|
| Read rate (last minute) | 583 MB/s | 380 MB/s |
| Read rate (last hour) | 562 MB/s | 384 MB/s |
| Read rate (since restart) | 589 MB/s | 49 MB/s |
| Write rate (last minute) | 1 MB/s | 101 MB/s |
| Write rate (last hour) | 2 MB/s | 117 MB/s |
| Write rate (since restart) | 25 MB/s | 13 MB/s |
| Master ops (last minute) | 325 Ops/s | 533 Ops/s |
| Master ops (last hour) | 381 Ops/s | 518 Ops/s |
| Master ops (since restart) | 202 Ops/s | 347 Ops/s |

## 7. COMPARISON BETWEEN GOOGLE FILE SYSTEM AND HADOOP DISTRIBUTED FILE SYSTEM

GFS is discussed above, and now this is small introduction about Hadoop system (HDFS). As shown in Fig. 4. The architecture of Hadoop system. And Table 4. Presented the comparison,



**Fig. 4. The HDFS architecture.**

## 8. CONCLUSIONS

The Google File System employs the important essential role in supporting large-scale data operating workloads on commodity components. While some design choices are specific to the specific setting, various of them may apply to data processing operations of the same magnitude and cost issues.

Starting by repeating ordinary file system theories in case of the current and anticipated application workloads and technological mechanism. Our observations led to obvious different views in the design areas. Treating component crashes as the normal rather than the exception, optimize for large files that are mostly appended to and then read, both extend and relax the ordinary file system interface to improve the overall system. This system provides fault tolerance by constant monitoring, replicating critical data, fast and automatic recovery. Chunk replication tolerate chunk server failures. The frequency of these failures improves a novel online repair mechanism which in regular and transparency fixes the failures and compensates for lost replicas as soon as can. Additionally, using the check summing to detect data corruption at the disk or IDE subsystem level, that becomes all common give the quantity of disks in the system. GFS met the storage needs in successful manner and is mostly used in Google as the storage platform for research and development as well as production data operations. It is an important method which enables user to continue to innovate and attack errors on the all of the whole web.

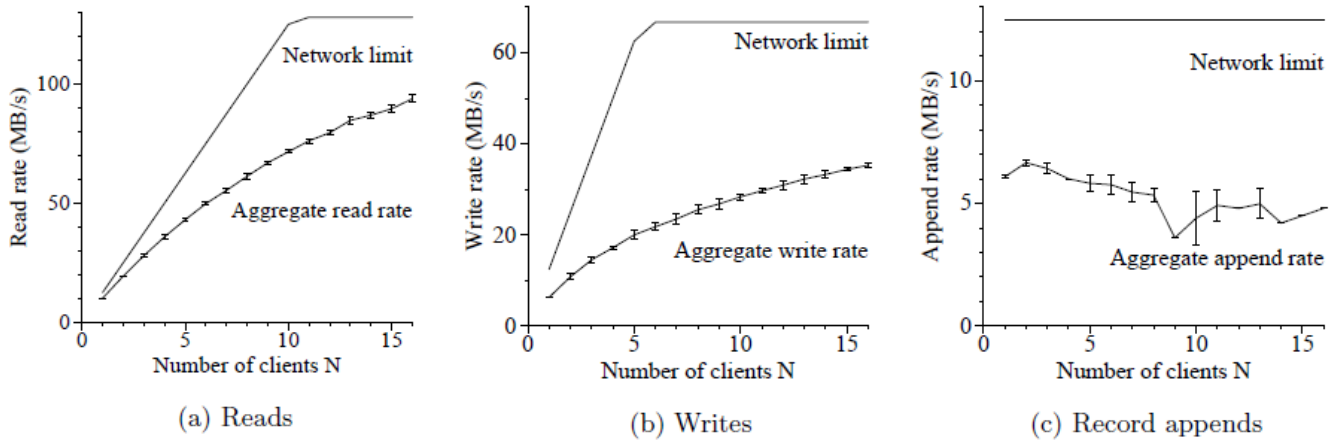(a) Reads  (b) Writes  (c) Record appends

Fig. 3: Aggregate Throughputs. Top curves show theoretical limits imposed by our network topology. Bottom curves show measured throughputs. They have error bars that show 95% confidence intervals, which are illegible in some cases because of low variance in measurements.

**Table 4. Some of the Comparative Analysis of GFS and HDFS**

| Properties | GFS | HDFS |
|---|---|---|
| Design Goals | The basic goal of GFS is to support large files<br>· made based on the proposals that:<br>TB data sets will be distributed across thousands of disks attached to Commodity compute nodes.<br>· Used for intensive computing data.<br>· reliably, till when failures occur within chunk servers, master, or network partitions.<br>· GFS is created more for batch processing better than interactive use by users. | The main goals of HDFS is to support large files.<br>· made based on the proposals that:<br>TB data sets will be distributed across thousands of disks attached to commodity compute nodes.<br>· Used for intensive computing data.<br>· reliably, till when failures occur within chunk servers, master, or network partitions<br>· HDFS is created more for batch processing rather than interactive use by users. |
| Processes | Master and chunk server | Name node and Data node |
| File Management | Files are supported hierarchically in directory and defined by path names.<br>· GFS is exclusively for Google application. | HDFS supported an ordinary hierarchical file organization<br>· HDFS also supports third-party file systems such as Cloud Store and Amazon Simple Storage Service. |
| Scalability | Cluster based architecture<br>· The file system consists of hundreds or even thousands of storage machines built from inexpensive commodity parts.<br>· The largest cluster has over 1000 storage nodes, over 300 TB of disk storage, and are heavily accessed by hundreds of clients on distinct machines on a continuous basis. | · Cluster based architecture<br>· Hadoop currently runs on clusters with thousands of nodes.<br>· E.g. Face book has 2 major clusters:<br>- A 1100-machine cluster with 8800 cores and about 12PB raw storage.<br>- A 300-machine cluster with 2400 cores and about 3PB raw storage.<br>- Each (commodity) node has 8 cores and 12 TB of storage.<br>EBay uses 532 nodes cluster (8*532 cores, 5.3PB)<br>· Yahoo uses more than 100,000 CPUs in >40,000 computers running Hadoop<br>- biggest cluster: 4500 nodes (2*4cpu boxes w 4*1TB disk & 16GB RAM) |

| | | [10] · K.Talattinis et.al concluded in their work<br>that Hadoop is efficient when operating in a full distributed system, although, in order to achieve optimal results and obtain advantage of Hadoop scalability, it is necessary to use large clusters of computers. |
|---|---|---|
| Protection | Google has their own file system called GFS. With GFS, files are divided and stored in multiple pieces on many machines. | The HDFS supposed a permission method for files and directories which shares<br>data of the POSIX model.<br>· File or directory have individual permissions |

## 9. REFERENCES

[1] Ghemawat, Sanjay, Howard Gobioff, and Shun-Tak Leung. "The Google file system." *ACM SIGOPS operating systems review*. Vol. 37. No. 5. ACM, 2003.

[2] Lu, Lanyue, et al. "A study of Linux file system evolution." *ACM Transactions on Storage (TOS)* 10.1 (2014): 3.

[3] 5- Yang, Jade. "From Google File System to Omega: a Decade of Advancement in Big Data Management at Google." Big Data Computing Service and Applications (BigDataService), 2015 IEEE First International Conference on. IEEE, 2015.K. Elissa, "Title of paper if known," unpublished.

[4] .''www.wekipidia.com''//last seen: Saturday, 19-12-2015.

[5] Vijayakumari, R., R. Kirankumar, and K. Gangadhara Rao. "Comparative analysis of Google File System and Hadoop Distributed File System."ICETS-International Journal of Advanced Trends in Computer Science and Engineering 3.1 (2014): 553-558.