

# Google File System (GFS): Paper Review

## Rohan Gore — rmg9725

Collaborators: Perplexity - ChatGPT o3

### References:

1. The Google File System (original paper)
2. Video explanation by Arpit Bhayani

## 1 Problem Statement Tackled

### 1.1 Component Failures as the Norm

GFS operates on the assumption that component failures are inevitable and frequent rather than exceptional.

Solution 1: The system implements a **heartbeat message** which acts as a way of failure detection and fault tolerance - since it does re-balancing of chunks.

Solution 2: Checksums are implemented to identify corrupted chunks. This acts as a way of self diagnose and repairing.

### 1.2 Optimization for Large Files

Traditional file systems optimize for managing large numbers of small files, but GFS targets multi-GB files containing many application objects.

Solution 1: Chunking introduced to handle and optimize the systems for large files.

Solution 2: Strategic design decision made of NOT optimizing the workflows for small files.

Solution 3: If we have important small files like executables then their replication factor will be increased in order to handle hot spots and bottlenecks.

### 1.3 Append-Heavy Workload Patterns

Most files in GFS are mutated by appending new data rather than overwriting existing data. Random writes within files are practically non-existent. Once written, files are typically read sequentially.

Solution 1: During read/write communications we use offset to identify exact actions positions.

Solution 2: Chunking helps in quickly reaching the action positions.

Solution 3: For making writes efficient there is a LRU Buffer at every chunk server which stores the updates and only writes after the write order is finalized by the primary replica.

## 2 Design Principles and Interface

- **Separation of Control and Data Path:** Metadata operations are handled by a single master; data reads/writes occur directly between clients and chunkservers to avoid bottlenecks.
- **Relaxed Consistency Model:** To simplify the system and increase performance, GFS relaxes POSIX consistency guarantees, giving 2 new operations called Snapshot and Atomic Append.
- **Efficient Metadata Management:** Metadata is stored in memory on the master for fast access with persistent operation logs for durable state.
- **Garbage Collection:** Deleted files are renamed and eventually cleaned up lazily during background scans, providing safe deletion and easy recovery from accidental removals.
- **Integrity and Recovery:** Checksumming detects data corruption on chunkservers, triggering replica reconstruction. The master coordinates re-replication and load balancing across chunkservers.

## 3 Architecture and Components

### 3.1 BLOCK 1: Master-Chunk Server

- Stores metadata (names, file size, ACI, chunk server info., chunk → server mapping) in-memory
- Heartbeat Messages that sync all the info periodically with master.
- Maintains operation log file on log servers for fault tolerance of itself and easy recovery.

### 3.2 BLOCK 2: Chunk Servers

- These are the machines where chunks are actually stored.
- The chunks are not unique to the machine as they are replicated across chunk servers - one of which is called a primary replica.

### 3.3 BLOCK 3: Application and GFS Client

- An application uses GFS at backend and has a GFS Client coded using its GFS SDK.
- When data is recieved after retrieveing/reading the internal accumulator of GFS Client packages the information in the way the application has asked for - basically working like an API.

### 3.4 ALGORITHMS Used

- **Chunking algorithm:** chunks the incoming file writes - algo. situated in master server.
- **Replication algorithm:** Distributes chunks such that it maximizes data reliability and network bandwidth usage
- **Atomic Append:** Makes file appends faster - used during writes.
- **Snapshot:** Takes a snapshot of disk at master server and uses this for check-pointing the progress at log server. Also includes the logic of updating the operation log replicas by flushing after successful checksum CRC checks.

### 3.5 WORKING 1: READ Request working

1. Client translates file name and byte offset into chunk index using fixed chunk size
2. Client requests chunk handle and replica locations from master
3. Master responds with chunk handle and locations of all replicas
4. Client caches this information and contacts closest replica directly
5. Subsequent reads of same chunk require no master interaction until cache expires

### 3.6 WORKING 2: WRITE Request working

Involves 2 main things namely, updation & communication process across chunk servers, and, leasing & mutation.

1. Master grants chunk lease to one replica (the primary) with 60-second timeout
2. Client pushes data to all replicas in pipelined fashion
3. Client sends write request to primary after all replicas acknowledge data receipt
4. Primary assigns serial numbers to mutations and applies them locally
5. Primary forwards request to secondary replicas with same serial order
6. Secondaries apply mutations and acknowledge to primary
7. Primary responds to client with operation status

The lease mechanism minimizes master overhead while ensuring consistent mutation ordering across replicas.

## 4 Related Work

- **Namespace vs. Striping:** Like AFS, GFS exposes a location-independent namespace, yet it stripes each file across many storage nodes (à la xFS and Swift) to raise aggregate throughput and fault tolerance.
- **Replication over RAID-style Parity:** Because commodity disks are inexpensive, GFS relies solely on simple, three-way replication for redundancy—easier to implement than RAID or parity codes, though it consumes more raw capacity than xFS or Swift.
- **No Client-side Caching:** Unlike AFS, xFS, Frangipani, or Intermezzo, GFS deliberately omits block caching beneath the file-system interface; its workloads either stream through large files or perform sparse random reads with negligible data reuse.
- **Centralized Metadata Service:** Whereas Frangipani, xFS, Minnesota GFS, and GPFS distribute metadata management, GFS keeps a single, lightweight master to simplify chunk placement and replication policies; high availability is achieved through shadow masters and a write-ahead log.

## 5 Conclusion

The Google File System demonstrates successful adaptation of distributed storage design to specific application requirements and operational constraints. The system's success stems from co-design with applications and willingness to challenge traditional file system assumptions. While some design decisions reflect Google's and its client's unique environment, the resulting system being ultra-reproducible is a big win. Key innovations include the separation of control and data flows, lease-based consistency management, and the recognition that relaxed consistency can significantly improve performance without compromising application requirements. These insights continue to influence modern distributed storage system design.