# Google File System.

→ This was a system which was optimized for **concurrent**
  sequential •writes & reads. Okay Performance for random s/w.
  ↓ called secandappends.

\* Architecture:

→ Single master & multiple chunk servers. (C → Chunk)   (CS → chunk server) 64 bit immutable & unique



File
xyz.txt

Division →

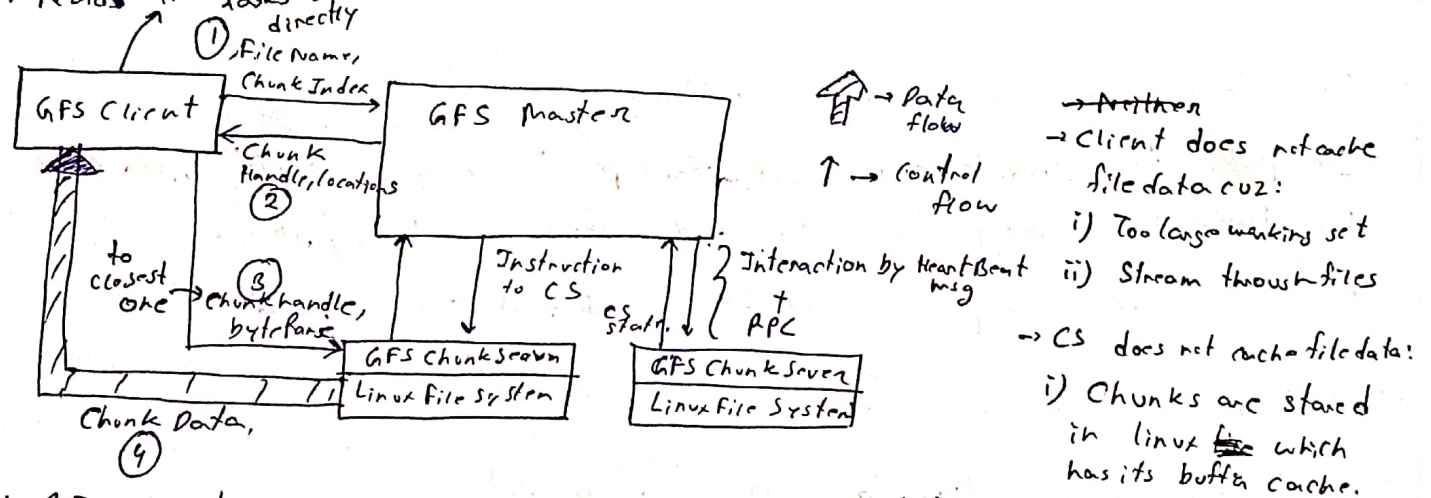| | Size | Version | Chunk handle Loc |
| C1 □ xyz-C1.txt → 64MB | V1 | f4xyzaq | CS1 |
| C2 □ xyz-C2.txt  " | V1 | q3 xyzaz | CS5 |
| C3 □ xyz-C3.txt  " | V1 | b2 xxxxx | CS2 |
| C4 □ xyz-C4.txt  " | V1 | c6xyzo8 | CS3 |

→ Each chunk is replicated ☒ 3 times (configurable) for reliability.

→ Each chunk has → i) Size = 64 MB
   ii) Version #o
   iii) Chunk Handle → 64bit unique & immutable id
   iv) a chunk location.
   v) Divided in 64KB blocks & each block has a CRC

\* Reads: Can cache the chunk locations asks CS directly



① File Name, Chunk Index →

GFS Client

② Chunk Handle, locations ←

GFS Master

③ Chunk handle, byte Range → to closest one

Chunk Data, ④

GFS Chunk Server / Linux File System

Instruction to CS

CS State ← RPC

Interaction by HeartBeat msg

GFS Chunk Server / Linux File System

⬆ → Data flow
↑ → Control flow

→ Neither
→ Client does not cache file data cuz:
  i) Too large working set
  ii) Stream through files

→ CS does not cache file data:
  i) Chunks are stored in linux fs which has its buffer cache.

∴ No caching ⟹ No cache coherence issue.

\* GFS Master:

→ The master has **multiple** non-volatile & volatile data:

• File Name → 🖐 Chunk handle mapping (Non volatile)
• Chunk handle & its version #  (Non volatile)
• Primary vote (volatile)   • Chunk location (volatile)
• Lease time  (volatile)  → Informed to Chunk Master by chunk server at startup.

⎫ All are stored in RAM for fast Response.

But File Name & Chunk handle → In persistent storage. [log + Checkpoint]
  to CS mapping + version #                              ‾‾‾‾‾‾‾‾‾‾
                                                         Replicated to remote machines.

→ Functions of master: Main → Failure Recovery
  i) Name space & Locking & Management
  ii) Replica Placement
  iii) Creation, Re-replication Rebalancing.
  iv) Garbage Collection.
  v) State Stale Replica Detection.

→ why is chunk size 64 MB?

Pros:

i) Reduces client's need to interact with master

Morechange

Large file ⟹ I have to read & write ⟹ Less Req to master.
from same file

ii) Can keep a persistent TCP connection & reuse it instead of multiple TCP conn. for smaller chunk sizes.

iii) Master keeps less meta data.

Cons:

i) Single / Small files ⟹ Less Chunks ⟹ Can become Hotspots.

⇩

Fixed by storing them with high replication factor.

→ Problem with single masta:

i) Limited by memory. 64MB size chunk needs 64B namespace.

↳ Can be fixed by more storage.

ii) No auto recovery ⟹ Solved in Raft.

→ How does master know chunk location?

→ It polls the chunk servers at startup. It keeps itself upto date by Heart Beat msg with CS and when CS leave/join.

→ This eliminates the problem of having CS & master synced
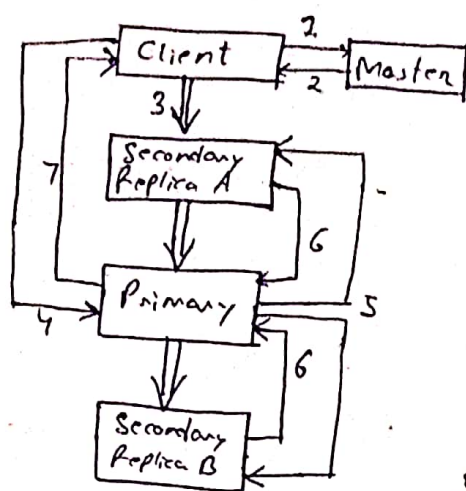
→ How does master recover from failure?

→ It has a operation log. & does checkpoints.

since

Operation log: Every item is uniquely identified by some ids. whatever operation the master is asked to do, it first logs the operation and then starts doing it. This is so cuz it it fails in between it knows what operation to do. When logs grow beyond a certain size → It checkpoints the state.

→ whenever a master restarts, it starts with state from last checkpoint and performs the log operations subsequently. The logs & checkpoint data are also stored replicated & stored in other remote machines.

* Writes → Means record appends.

→ The master leases a particular replica of a chunk for 60s to be primary. If for whatever reason, ~~most primary~~ master needs to wait for lease expiry before making someone else the primary.

1) Casks M for file's last chunk.

2 ~~8~~) If M sees chunk has no primary (or lease expired):

 a) If no chunk servers with latest Version # → Error.

 b) Pick Primary & Secondary from those with latest version #

 c) increment version # → write to disk    log on

 d) tell P and Secondaries who they are and new version #

 e) replicas write new version # to disk.

2) ~~Client sends~~ M tells C the P & S

3) Client Sends it data to all replicas & P [data flows from client to its closest then to its next closest Er: Client → Sec Rep A → P → Sec Rep B This is done to avoid network bottlenecks.]

4) C ~~sends data~~ tells P to append.    else write is done on newchunk same by all Secondaries.

~~5~~ → P checks that lease has not expired and chunk has space, Picks a offset and writes to that offset.

5) P tells @ all S to write at the same offset

6) P waits for all secondaries to reply OK if even one response is error on something. The whole 4 & 6 is retried

7) P tells C ok or error.

8) If error, it is retried multiple times.

* GFS does not guarantee that all ~~records~~ replicas are identical byte to byte but it guarantees that each ~~record~~ record is written atleast once [due to retry]

* All replicas write stuff at same offset.

* In successful record appends => The regions are defined => consistent.
* In Interleaving ~~unsuccessful~~ "     "   =>  "    "    " undefined but consistent

* Guarantees By GFS:            Mutation = ~~Write~~ Record Append.            ④

→ A file region is consistent ⇒ all replica sees same data irrespective of
→ A file region is defined ⇒ which replica it is read from.
                              If after a mutation, clients can see
* Interleav regions.           what ~~it~~ was written i·e

AAA
    ⟩⟩  AABB    instead    AAAA → AAAA
BBB            of         AAABBB


→ Concurrent successful mutations may leave interleaving ~~or~~ fragments
  . The region is undefined. but consistent.

⇒ Failed mutation is inconsistent ⇒ undefined region.

⇒ GFS may insenst padding on record duplicates in inconsistent
  region.

→ After seq of successful mutation, region is guaranteed to be defined
    i) applying mutation to all chunk replicas in same order.
                                    chunk
    ii) using chunk Version # to ensure no stale chunk replica was used

⇒ So by atomic record appends → ~~In cons~~ Undefined & Inconsistent
                                          Case is removed.


Then flow does application check for duplicate / padding regions/records?

① Applications use checksums and/or a unique Id for each record
   Valid checksum ⇒ Record valid, Unique Id ⇒ No record is read
                                                        twice

* Snapshot :
  → ~~Sta~~ Snap Shot operation makes a copy of a file or a directory tree structure
    almost instantaneously, ~~which~~

  → Used to create copies or check points.

  How is it implemented: Main idea is immeadiate copy of
                         directory and lazy copy of chunk.
  i) Master receives Snapshot request.
  ii) Revokes any outstanding request on chunk it is about
      to shapshot
      ↳ This means it anyone wants access to the file req goes through master
        which allows it to be copied.

(iii) After leases are revoked / expired, the master logs the (5) operation to disk.

(iv) Applies the log record to its in memory ~~store~~ state by duplicating the metadata for source ~~and~~ file an directory tree.

(v) The newly created snapshot files point to same chunk.

→ Now Client wants to write to chunk C after snapshot.
i) Client asks Master for primary & S.   two tree struct point to same chunk.
ii) Master sees that C has ~~ref~~ ~~cont~~ count more than 1. It defers replying to client req & instead picks & a new chunk C'.
                                                                              handle
iii) Asks chunk servers to that has a current replica of C to create a new chunk called C'. By doing this we assure that data is copied on same chunkserver as the original so no network overhead.

iv) Client replies info about C' and all writes are ~~mow~~ now done to C'.

**✳ Master operations:**

i) Name space Management and Locking
→ Multiple operation are allowed to be active
→ Locks are used for serialization.

→ GFS logically represents its ~~lookup~~ namespace as lookup table mapping full path names to meta data. [ Path compression done to save space ] only on last item
→ There is no i-node thing like Linux.
→ Each directory / file as a read / write lock.

Ex:
A directory called /home/user ~~&~~ can have a file /home/user/foo.
To create the file, we put read locks on /home → /home/user and write lock on /home/user/foo. ① The read lock on user shows that something in its children has write so it cant bedeleted
② There can be multiple writes in /home/user → /whatever
                                                                                      w.

## ii) Replica Placement.

Two purpose : ① Maximize data reliability & availability
② Maximize Network utilization.

→ Replicated across multiple racks.

↪ Trade off: reads can be faster, writes slower
  ↳ Not a problem.

→ The machines are allocated IPs such that closer machine has
  closer IP.

## iii) ß Creation, Re-balancing & Re-replication:

→ when master creates a chunk, it chooses $^\circ$ to place the empty
where
replicas. The factors are

① Place on below avg disk space utilization.

② Minimize # of recent creation on each chunk server.
  Recent creation ⇒ ~~follo~~ followed by writes there.

③ Spread across racks.

→ ~~Ref~~ Master re replicates as soon as # of replicas fall below
  user specified goal. The factors are:

① ~~b~~ How far is it from replication goal.

② Re replicate chunks for live files as opposed to recently deleted files.

③ Minimize Impact ~~on~~ of failure on running application. we
  boost priority of any chunk that is blocking user progress.

→ Rebalancing is done periodically. It examines the current
  replica distribution and moves replicas for better
  disk space & load balancing.

IV) Garbage Collection:

→ when a **file** is deleted, immediately

① Master logs the deletion immediately:

② The file is renamed to a hidden name that includes the deletion times.

③ During master's regular scan of file system, it removes any such hidden files if the have existed for 3 days (configurable)

④ Until then, file can be read under new ~~file~~ special name.

⑤ When ~~for~~ the hidden file is removed from namespace, its in-memory meta data is erased. This effectively severs its ~~protection~~ links to all its chunks.

→ In similar scan of **chunk** namespace, the master identifies orphaned chunks (unreachable from any file). and erases the meta data for those.

→ In heart beat msg regularly exchanged, with master, each chunk server reports a subset of ~~its~~ chunks it has and the master ~~appears~~ replies with the identity of all chunks that are no longer present in master's meta data. The chunk server is free to delete its replicas.

v) Stale Replica Detection:

→ Master keeps track of lastest version number & version # is increment during writes [ See ~~section~~ writes part].

→ ~~Co~~ Replicas may miss mutation if it is down.
   of chunkserver

→ During scanning, it can match the version numbers ~~to~~ of a replica and ~~lat~~ latest to and marked as stale which can be re-replicated & deleted [ Prev parts)

* Fault Tolerance:

i) High availiability
→ Fast Recovery ⇒ Fast startup incase of Discussed
master / chunkserver failure.

↘ Chunk Replication ⇒ For stale chunks, creation Discussed
of new chunks have multiple replicas.

Master Replication

→ There are replicas of master replicas with full copy of master's state.
Paper's design requires human intervention to switch to one of the replicas
after a master failure.

ii) Data Integrity: Done by check Sums! ↗ cant compare replicas since no guarante of identical byte by byte.

→ Each chunk is broken into 64 kB blocks & Each block has a check sum. They are kept in memory & logged persistently.

→ Read Path:

   i) verify the checksum of record with the existing one.

   ii) On mismatch. ⇒ Returns error & reports corruption to master.

   iii) Then the client reads from a different replica. & master
    ~~clones a~~ does re replication [discussed]

→ write Path: Check sum is CPU extensive. So optimised for
record appends.

   Append Case: When a ~~P~~read data server does not
verify for the block. It incrementally updates. the checksum
for the block. like old checksum & Check(Append) = New
                        Sum data Check-sum.
Cuz If last block was corrupt the new will be corrupt too.

→ Overwrite Case: If you overwrite, Read & verify for first &
last blocks that are touched, write the new data. Compute
new checksums. Why? → If first & last are corrupt then the
new data might cover up old corruption [ ~~or~~ Like corruption
can be hidden ] [ Paper ]

→ During idle time, ~~master~~ Chunk server can scan
and verify the checksums it corrupted ∩ ^repeated to master and re
replicated. (discussed)