

A case for RAID

RAID-1: Mirrored Disks

Mirrored disks are the most expensive option. Every write to a data disk is also a write to a check disk. Since these writes happen in parallel, it should be at least half as fast.

Main drawback: duplicating all disks means doubling the cost of the disks or using only halving storage capacity.

RAID-2: Hamming Code for ECC

If all data bits in a group are read or written together, there is no impact on performance. Reads of less than the group size require reading the whole group for error checking, and writes to a portion of the group require reading the whole group, modifying the data, and writing the full group, including check information. 1 parity disk can detect 1 error, but to correct an error we need enough check disks to identify the disk with the error. For 10 data disks, we need 4 check disks. Good for supercomputers, bad for TPS.

RAID-3: Single Check Disk Per Group

Only 1 parity bit is needed to detect an error. Information on a failed disk can be reconstructed by calculating the parity of the remaining good disks and comparing bit by bit to the parity calculated for the original full group. If the check disk is the failure, read all data disks and store group parity in the replacement disk. Drawback vs RAID-2: Soft errors can be corrected without rereading a sector. Also, the extra check information for each sector is not needed, increasing disk space. Good for supercomputers but bad for TPS

RAID-4: Independent Reads/Writes

No longer spreads individual transfer information across several disks, keep each individual unit in a single disk. Allows parallelism for reads - the ability to do more than one IO per group at a time. Parity calculation is much simpler: $\text{new parity} = (\text{old data} \oplus \text{new data}) \oplus \text{old parity}$. The check disk in a group must be read and written with every small write in that group, making it the bottleneck.

RAID-5: No single check disk

Distributes data and checks across all the disks - including the check disks. Because of this, RAID 5 can support multiple individual writes per group. Best of both worlds: small read/writes perform close to the speed of a level 1 RAID and it keeps the large transfer performance and high storage capacity of RAID 3/4. RAID-5 is good for supercomputers or transaction processing.

The Google File System

GFS is a scalable distributed file system for large distributed data-intensive application. It has fault tolerance and delivers high aggregate performance to a large number of clients.

Assumptions

Component failures are common
Modest number of very large file
Reads are large streams or small & random
Large sequential writes that append data
Multiple clients concurrently append to the same file
High bandwidth is more important than low latency

New Operations

Snapshot creates a copy of a file or directory at low cost
Record append allows multiple clients to append to the same file and guarantees the atomicity of each append.

Architecture

A GFS cluster has one master and many chunkservers, and is accessed by many clients
Files are divided into 64MB chunks
Each chunk is replicated on multiple chunkservers
The master maintains all system metadata, chunk lease management, garbage collection, and chunk migration.
File data is not cached by the client or the chunkserver
Metadata is cached by the client
GFS does not have a per-directory data structure that lists all the files in that directory.
It does not support aliases for the same file (symbolic links)
GFS represents its namespace as a lookup table mapping full pathnames to metadata, stored in memory using prefix compression.

Master

Clients do not read and write files through the master, they ask the master which chunkservers it should contact, and interact with this chunkserver for subsequent operations
The master executes all namespace operations, manages chunk replicas, makes placement decisions, creates new chunks and replicas, coordinates system-wide activities to keep chunks fully replicated, to balance load, and to reclaim unused storage.

Metadata

Master stores 3 types of metadata: File/chunk namespaces, mapping from files to chunks, and locations of each chunk's replicas.
All metadata is kept in the master's memory, so master operations are fast.
Master does not keep a record of which chunkservers have a replica of a given chunk. It polls chunkservers for that information at startup. A chunkserver has the final word over what chunks it does or does not have on its own disks.
The operation log contains a historical record of critical metadata changes.
The operation log is the only persistent record of metadata and serves as a logical time line that defines the order of concurrent operations.
Files, chunks, and their versions are all uniquely identified by the logical times at which they were created.
The master recovers its state by replaying the operation log.
The master's state is checkpointed whenever the log grows too large, so it can recover by loading the latest checkpoint and replaying the log records after that.

Leases

The lease mechanism is designed to minimize management overhead at the master. Each mutation happens at all of the chunk's replicas.
Leases are used to maintain a consistent mutation order across replicas.
The global mutation order is defined first by the lease grant order chosen by the master, and then by a serial order that the machine granted the lease chooses.

Consistency model

GFS has a relaxed consistency model

File namespace mutations are atomic.

A file region is consistent if all clients will see the same data.

A file region is defined if it is consistent and all clients will see the last mutation in its entirety.

When a mutation succeeds without interference from other writers, the affected region is defined.

Concurrent successful mutations leave the region undefined but consistent.

A failed mutation makes the region inconsistent.

Mutations may be writes or record appends

A record append causes data to be appended atomically at least once even in the presence of concurrent mutations, but at an offset of GFS's choosing.

GFS may insert padding or record duplicates in between. These occupy inconsistent regions and are usually small.

After a sequence of successful mutations the file region is guaranteed to be defined and contain the data written by the last mutation

A chunk is lost irreversibly only if all of its replicas are lost before GFS can react, typically within minutes. Even in this case, it becomes unavailable rather than corrupted.

Atomic Record Appends

In a record append, the client specifies the data, and GFS appends it to the file at least once atomically at an offset of GFS's choosing and returns that offset to the client.

If a record append fails at any replica, the client retries the operation.

As a result, replicas of the same chunk may contain different data possibly including duplicates of the same record in whole or in part.

GFS does not guarantee that all replicas are identical, just that the data is written at least once as an atomic unit.

Snapshot

The snapshot operation makes a copy of a file or directory tree almost instantaneously, using standard copy-on-write techniques.

Stale Replica detection

Chunk replicas become stale if a chunkserver fails and misses mutations to the chunk while it is down.

For each chunk, the master maintains a chunk version number to distinguish up-to-date and stale replicas

Whenever the master grants a new lease on a chunk, it increases the chunk version number and informs the up-to-date replicas.

The master removes stale replicas in its garbage collection.

Before that, it considers stale replicas to not exist when it replies to client requests for chunk information.

Fault Tolerance

We cannot trust the machines, or the disks

Two strategies keep the system highly available: fast recovery and replication

Fast recovery: Both the master and the chunkserver are designed to restore their state and start in seconds no matter how they were terminated.

Chunk replication: Each chunk is replicated on multiple chunkservers on different racks

Master replication: The master state is replicated for reliability. Its operation log and checkpoints are replicated on multiple machines.

Each chunkserver independently verifies the integrity of its own copy by maintaining checksums

Concurrency Control and Recovery

Concurrency control ensures users see consistent states of the database even though their operations may be interleaved by the database

Recovery ensures that the database is not corrupted as the result of a software, system, or media failure.

Transaction

A transaction is a unit of work that must commit or abort as a single atomic unit.

If it commits, all updates it performed are made permanent and visible to other transactions.

If it aborts, all of its updates are removed from the database and the database is as if the transaction had never executed.

ACID Properties

Parity calculation is much simpler: $\text{new parity} = (\text{old data} \text{ xor } \text{new data}) \text{ xor } \text{old parity}$. The check disk in a group must be read and written with every small write in that group, making it the bottleneck.

RAID-5: No single check disk

Distributes data and checks across all the disks - including the check disks. Because of this, RAID 5 can support multiple individual writes per group.

Best of both worlds: small read/writes perform close to the speed of a level 1 RAID and it keeps the large transfer performance and high storage capacity of RAID 3/4. RAID-5 is good for supercomputers or transaction processing.

The Google File System

GFS is a scalable distributed file system for large distributed data-intensive application. It has fault tolerance and delivers high aggregate performance to a large number of clients.

Assumptions

Component failures are common

Modest number of very large file

Reads are large streams or small & random

Large sequential writes that append data

Multiple clients concurrently append to the same file

High bandwidth is more important than low latency

New Operations

Snapshot creates a copy of a file or directory at low cost

Record append allows multiple clients to append to the same file and guarantees the atomicity of each append.

Architecture

A GFS cluster has one master and many chunkservers, and is accessed by many clients

Files are divided into 64MB chunks

Each chunk is replicated on multiple chunkservers

The master maintains all system metadata, chunk lease management, garbage collection, and chunk migration.

File data is not cached by the client or the chunkserver

Metadata is cached by the client

GFS does not have a per-directory data structure that lists all the files in that directory.

It does not support aliases for the same file (symbolic links)
GFS represents its namespace as a lookup table mapping full pathnames to metadata, stored in memory using prefix compression.

Master

Clients do not read and write files through the master, they ask the master which chunkservers it should contact, and interact with this chunkserver for subsequent operations

The master executes all namespace operations, manages chunk replicas, makes placement decisions, creates new chunks and replicas, coordinates system-wide activities to keep chunks fully replicated, to balance load, and to reclaim unused storage.

Metadata

Master stores 3 types of metadata: File/chunk namespaces, mapping from files to chunks, and locations of each chunk's replicas.

All metadata is kept in the master's memory, so master operations are fast.

Master does not keep a record of which chunkservers have a replica of a given chunk. It polls chunkservers for that information at startup. A chunkserver has the final word over what chunks it does or does not have on its own disks.

The operation log contains a historical record of critical metadata changes.

The operation log is the only persistent record of metadata and serves as a logical time line that defines the order of concurrent operations.

Files, chunks, and their versions are all uniquely identified by the logical times at which they were created.

The master recovers its state by replaying the operation log.

The master's state is checkpointed whenever the log grows too large, so it can recover by loading the latest checkpoint and replaying the log records after that.

Leases

The lease mechanism is designed to minimize management overhead at the master. Each mutation happens at all of the chunk's replicas.

Two phase locking forces all transactions to be well-

formed and follow the rule: once a transaction has released a lock, it cannot obtain any additional locks. Two phase locking is sufficient but not necessary to implement serializability.

Phantom problem

One transaction reads twice, another transaction adds a tuple that satisfies the query in between the reads. The two read transactions will differ, which could not occur in a serial schedule

Isolation levels

Trades consistency for concurrency in a controlled manner.

Read uncommitted (Degree 0): No read locks obtained.

Risk is seeing updates that will be rolled back, or seeing incomplete transactions.

Read committed (Degree 1): Holds read locks on individual items for short duration.

Risk is seeing non-repeatable reads. Two reads could see two different values.

Repeatable read (Degree 2): Holds read locks on individual items for long duration.

Risk is the phantom problem.

Serializable (Degree 3): Holds read locks on predicates for long durations.

Solves the phantom problem.

Hierarchical locking

Concurrent transactions can obtain locks at different granularities

A lock at a particular level locks all items underneath it.

Lock escalation: If a transaction obtains many locks on one level, it could be given a lock at the level above.

ARIES

An implementation of recovery.

Uses WAL, with STEAL/NOFORCE buffer management policy:

Pages on stable storage can be overwritten and data pages do not need to be forced to disk to commit a

transaction.

Its REDO repeats history: It redoes updates for transactions that will eventually be undone. This allows it to employ physiological logging.

It uses page oriented REDO and logical UNDO.

First pass processes the log forward from the most recent checkpoint, determines dirty pages and active transactions.

Second pass repeats history by processing the log forward.

Third pass proceeds backwards from the end, removing the effects of all uncommitted transactions.

Uses fuzzy checkpoints that are extremely inexpensive.

Log Structured File System

Writes all modifications to disk in a log-like structure, speeding up file writing and error recovery.

Log is divided into segments, and a segment cleaner to compress live information from heavily fragmented sections.

An implementation, Sprite LFS, outperforms current unix file systems by an order of magnitude for small writes.

This system eliminates all seeks while writing.

Problems with existing file systems

Current file systems spread information around the disk in a way that causes many small accesses

Furthermore, the inode for a file is separate from its contents, as is the directory entry containing its name.

When writing small files, less than 5% of the disk's bandwidth is used for new data, the rest of the time is seeking.

The second problem is that they tend to write synchronously. The application must wait for the write to complete rather than continuing while the write happens in the background.

Log structured file systems

Fundamental idea: improve write performance buffering a sequence of file system changes in the file cache and then writing all changes to disk sequentially in a single write operation.

For workloads that contain many small files, a log struc-

tured file system converts the many small synchronous random writes of a traditional file system into large asynchronous sequential transfers that can use 100% of a disk's bandwidth. Two key issues: How to retrieve information from the log, how to manage free space on disk so that large extents of free space are always available for writing new data.

File location and reading

Outputs index structures in the log to permit random access retrievals. For each file there is an inode. Unlike unix FFS, inodes are not in a fixed position in the disk, instead, there is an inode map, which maintains the current location of each inode. The inode map is small enough to be cached, so it rarely requires disk accesses.

Segments

The goal is to maintain large free extents for writing new data.

When the log reaches the end of the disk, two things can be done: Threading (writing the log to unused portions of the disk) or copying (moving live data out of the log and writing a compressed version of it at the head of the log).

Sprite LFS uses a combination of both. The disk is divided into large areas called segments. Each segment is written sequentially, and all live data must be copied out of the segment before it can be rewritten, but the log is threaded on a segment basis.

Segment cleaning

Copying live data out of a segment is called segment cleaning. Read a number of segments into memory, identify the live data, and write the live data back out a smaller number of clean segments. After this operation is complete, the segments read are now clean.

A segment summary block identifies each piece of information written in the segment, so that the inodes can be updated when the segment is cleaned.

Sprite LFS cleans segments when the number of clean segments drops below a certain threshold, and it cleans until the number of clean segments rises above another threshold value.

Segments are divided into two types: hot and cold. Hot segments were written to very recently, and their

free space is less valuable (since they will be written to again), and cold segments have old data, and their free space is more valuable.

To calculate hot and cold segments, there is a segment usage table that contains the number of live bytes in the segment and the most recent modified time of any block in the segment.

Hot segments are cleaned more regularly than cold segments.

Crash recovery

When a system crash occurs, the last few operations may have left it in an inconsistent state.

In traditional unix file systems, the system must scan all of the metadata structures on disk to restore consistency. The cost of these scans is very high.

In LFS, the locations of the last operations are easy to determine: they're at the end of the log.

Sprite LFS uses two methods to recover: checkpoints, which define consistent states of the file system, and roll-forward, which is used to recover information written since the last checkpoint.

Checkpoints

To write a checkpoint, Sprite LFS first writes out all modified information to the log, including data, inodes, inode map, and segment usage table. Second, it writes a checkpoint region to a special fixed position on the disk.

The checkpoint region contains the addresses of all the blocks in the inode map and segment usage table, plus the current time and last segment written.

Roll-forward

Sprite LFS reads through the log entries that were written after the last checkpoint. A directory operation log entry is created for each directory operation and this is used to keep consistency between the directory entries and inodes.

PNUTS

PNUTS is a parallel and distributed database system used by Yahoo's web applications.

It provides data storage as hashed or ordered tables, with low latency for many concurrent requests, and per-record consistency guarantees.

Relaxed Consistency

No need for serializable transactions, or even transactions at all. Most users only modify one record at a time. One possible idea was eventual consistency, but that is too weak: if a user wants to remove his mother from the list of people who can access his photos before adding a photo, these updates must happen in that order.

Data and Query Model

Data is organized into tables of records with attributes. Flexible schemas: new attributes can be added, and records are not required to have values for all attributes. There are two main types of access: Point access involves a single record, range access involves a set of records. PNUTS allows hashed or ordered tables to support both types of requests.

Consistency Model

Between the extremes of serializability and eventual consistency. Per record timeline consistency: all replicas of a record apply all updates to the record in the same order. To implement, a master replica receives all updates to that record. There are several API calls with varying levels of consistency guarantees:

Read-any: returns a possibly stale version of the record

Read-critical(required-version): returns a version of the record that is at most as old as required-version

Read-latest: Returns the latest copy of the record

Write: Performs a normal write

Test-and-set-write(required-version): Performs a write to a record if and only if the version is the same as required-version. This is used for transactions that read a record and then modify it.

Architecture

Tables are horizontally split into groups of records called tablets.

Storage unit: Stores tablets, responds to get(), scan() and set()

Router: determines which tablet has a record, and which storage unit has that tablet

Tablet controller: Polled by the routers to obtain updated mappings. Determines when to move a tablet between storage units

Message Broker: A publish/subscribe system, replaces

redo log, used for replication. Provides partial ordering of published messages.

Scatter-Gather engine: Component of the router, responsible for multi-record requests.