# Distributed System Course Project
## Implementing Exactly-Once Semantics in RecordAppend for Google File System

**Mayank Mittal (2022101094)**
**Shivam Mittal (2022101105)**

November 28, 2024

# Contents

# Chapter 1

# Introduction to Distributed File Systems

## 1.1 Background

Distributed File Systems (DFS) are essential for managing large datasets across multiple servers. They provide key benefits such as:

- **Scalability**: Ability to handle increasing amounts of data by adding more nodes.

- **Fault Tolerance**: Ensuring data availability despite hardware failures.

- **High-Performance Storage**: Supporting parallel data processing and efficient I/O operations.

- **Distributed Management**: Enabling seamless data access across geographic locations.

## 1.2 Google File System Overview

The Google File System (GFS) is a robust DFS designed to meet Google's internal needs for handling vast datasets. It supports:

- **Large-Scale Data Processing**: Ideal for big-data workloads.

- **Fault Tolerance**: Automatic recovery from hardware failures.

- **High Throughput**: Optimized for sustained data-intensive operations.

- **Append Operations**: Frequent file appends rather than overwrites.



Figure 1.1:

# Chapter 2

# Google File System Architecture

## 2.1 Components

Our distributed file system architecture comprises three key components:

### 2.1.1 Master Server

The master server handles the system's metadata, including:

- **Metadata Management**: Tracks file and chunk locations.

- **Client Coordination**: Directs clients to the appropriate chunk servers.

- **Configuration Management**: Monitors system-wide configurations.

### 2.1.2 Chunk Servers

Chunk servers store the actual data and manage replication:

- **Data Storage**: Holds file chunks.

- **Replication Management**: Ensures data is replicated across multiple nodes.

- **Read/Write Operations**: Handles client requests for data.

### 2.1.3 Clients

Clients interact with the master and chunk servers to perform file operations:

- **File Operations**: Read, write, and append data.

- **Communication**: Send requests to the master and chunk servers.

## 2.2  Directory Structure

The project's directory structure is organized as follows:

```
distributed-file-system/
 56/
    master.py                    # Main master server logic
    chunk_server.py              # Chunk server logic
    client.py                    # client logic
    c1/
        chunk_server.py          # chunk server 1
    c2/
        chunk_server.py          # chunk server 2
    ....                         # other chunk servers
 56_presentation.pdf
 56_report.pdf
 56_presentation.pptx
```

# Chapter 3

# Semantics: At-Least-Once vs. Exactly-Once

## 3.1   At-Least-Once Semantics in GFS

GFS traditionally uses **at-least-once** semantics, meaning:

- Each operation is retried if the server fails to acknowledge it.

- It guarantees that the operation will be performed *at least once*.

- Potential downside: **Duplicate writes** may occur if an operation is retried after a partial success.

**Drawbacks**:

- Data consistency issues due to multiple appends.

- Clients need to handle deduplication.

## 3.2   Exactly-Once Semantics in Our Project

We implement **exactly-once** semantics to ensure:

- Each operation is performed exactly once, even in the event of failures.

- Avoidance of duplicate data writes.

- Stronger consistency guarantees across distributed nodes.

**Implementation Techniques**:

- **Transaction IDs**: Unique IDs for each operation.

- **Two-Phase Commit**: A coordination protocol ensuring atomic operations.

**Advantages**:

- Strong consistency.
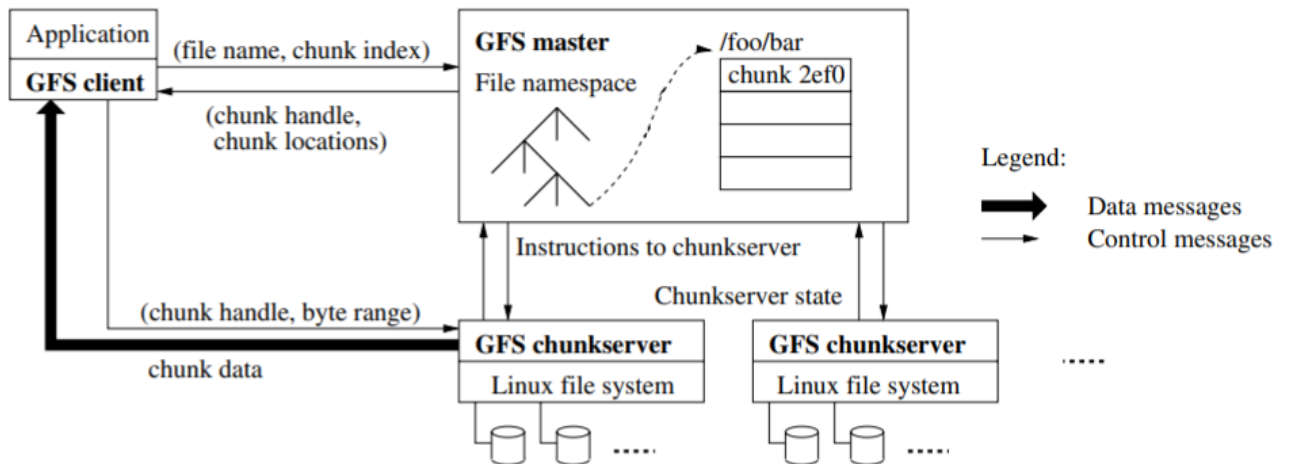
- Simplified client logic (no deduplication needed).

Figure 3.1: Architecture in Picture

## 3.3 Applications of Exactly-Once Semantics

- Financial systems: Ensures accurate transaction records.

- Messaging systems: Guarantees messages are delivered once.

- Logging systems: Prevents duplicate log entries.

# Chapter 4

# Core Operations Implementation

## 4.1 Exactly-Once Record Append

```python
def record_append(filename, data):
    """
    Implement exactly-once record append
    - Use unique transaction ID
    - Implement two-phase commit
    - Ensure atomic operation across replicas
    """
    transaction_id = generate_unique_transaction_id()

    try:
        # Prepare phase
        prepare_result = prepare_append(
            filename, data, transaction_id
        )

        if prepare_result:
            # Commit phase
            commit_append(
                filename, data, transaction_id
            )
        else:
            # Abort transaction
            abort_append(transaction_id)

    except Exception as e:
        # Handle potential failures
        handle_append_failure(transaction_id)
```

Listing 4.1: Exactly-Once Record Append Implementation

**Explanation**:

- **Prepare Phase**: In this phase, the system checks the readiness of all participating nodes (chunk servers) before proceeding with the actual operation. The system generates a unique transaction ID and sends a request to all nodes to verify if they can accommodate the append operation. This involves ensuring that there

8

is enough storage space, the data can be replicated as required, and no ongoing conflicts or errors are present. Only if all nodes respond positively does the system move to the next phase. The primary goal is to prevent partial updates that could lead to data inconsistency.

- **Commit Phase**: Once the prepare phase succeeds, the system finalizes the operation. The commit phase ensures that the data is appended to all nodes in a synchronized manner. The system uses a coordination protocol (such as Two-Phase Commit) to guarantee atomicity, meaning either all nodes commit the data, or none do. This phase also involves updating metadata in the master server to reflect the successful operation. If any node fails during this phase, the system initiates a rollback to maintain consistency.

- **Failure Handling**: Despite robust preparation, failures can occur during the commit phase or due to network issues. To address this, the system employs mechanisms to detect and handle failures, such as logging transaction IDs and performing retries. If a failure is detected during the commit, the system ensures that any partial writes are rolled back, maintaining the atomicity and integrity of the data. In some cases, an abort is triggered, and the client is informed of the failure.

## 4.2 Replication Strategy

Our replication strategy is designed to ensure high availability and durability of data:

- **Three-Way Replication**: Each file chunk is replicated across three different servers to ensure redundancy. This allows the system to tolerate up to two server failures without data loss. Replication provides both fault tolerance and load balancing by distributing read operations across replicas.

- **Primary Replica**: Among the three replicas, one is designated as the primary replica. This primary server coordinates all write and append operations to ensure consistency. The primary replica acts as the leader and manages updates to secondary replicas, enforcing a consistent ordering of writes.

- **Consistency Checks**: Periodic background tasks perform consistency checks across replicas to detect and resolve data mismatches. If a replica falls out of sync, the system triggers a re-replication process to restore consistency. These checks ensure data integrity over time.

- **Automatic Recovery**: In the event of a server failure, the system automatically identifies the lost replica and initiates a recovery process. A new replica is created by copying data from the remaining replicas. This self-healing mechanism is critical for maintaining data availability.

# Chapter 5

# Network Communication

## 5.1 TCP Socket Implementation

To ensure reliable communication between clients, master servers, and chunk servers, we use TCP sockets:

- **Persistent Connections**: Our system establishes long-lived TCP connections to avoid the overhead of repeated connection setup and teardown. Persistent connections improve efficiency, especially in a distributed environment where frequent interactions between components are necessary.

- **Serialization**: Data exchanged over the network is serialized using protocol buffers, a compact and efficient format. Serialization ensures that complex data structures can be transmitted as byte streams and reconstructed accurately at the receiving end.

- **Timeouts and Retries**: Robust error handling is implemented by setting timeouts for network operations. If a request times out, the system retries the operation a specified number of times. This approach minimizes the impact of transient network issues and enhances reliability.

- **Bidirectional Channels**: Both clients and servers can initiate communication, enabling dynamic and flexible interactions. This capability is essential for scenarios like real-time updates and push notifications.

# Chapter 6

# Performance Considerations

## 6.1   Scalability Factors

To handle the demands of a growing dataset and user base, we implement or will implement in future work several scalability strategies:

- **Sharding**: The dataset is partitioned into smaller, manageable pieces called shards, distributed across multiple servers. Sharding improves performance by spreading the workload and enabling parallel processing.

- **Load Balancing**: Client requests are distributed evenly across servers to prevent any single server from becoming a bottleneck. Load balancing enhances system responsiveness and maximizes resource utilization.

- **Caching**: Frequently accessed data is cached in memory to reduce the load on chunk servers. Caching minimizes disk I/O and speeds up data retrieval, improving overall system performance.

## 6.2   Latency and Throughput

Our system is optimized for both low latency and high throughput:

- **Low Latency**: We prioritize fast response times for client operations by optimizing the data path and using efficient communication protocols. Low latency is crucial for user experience in real-time applications.

- **High Throughput**: The system is designed to handle large volumes of data and concurrent operations. Techniques like parallel processing, efficient serialization, and high-speed networking contribute to maximizing throughput.

# Chapter 7

# Core Operations Implementation

## 7.1 File Operations: Read, Upload, Write, Create, and Record Append

In our system, the core file operations are implemented as follows:

- **Read Operation**: The client sends a read request to the master server, which redirects it to the appropriate chunk servers. The client then retrieves the data from the replicas, choosing the nearest one to minimize latency.

- **Upload Operation**: The client uploads a file by dividing it into chunks and sending each chunk to a set of chunk servers. The master server coordinates the process and maintains metadata about the file and its chunks.

- **Write Operation**: Writes are handled similarly to uploads but may involve updating existing chunks. The primary replica manages the write operation, ensuring consistency across replicas through a synchronized commit process.

- **Create Operation**: The client can create new files or directories. The master server updates the metadata and assigns chunk servers to store the initial chunks of the file.

- **Record Append**: The record append operation allows clients to append data to a file. This is especially useful for log files and other append-only datasets.

In this implementation:

- **Generate Transaction ID**: A unique ID is generated to track the operation and ensure exactly-once execution.

- **Prepare Phase**: All participating chunk servers verify that they can handle the append operation.

- **Commit Phase**: The data is appended only if all chunk servers confirm readiness.

- **Failure Handling**: If an error occurs, the system aborts the operation and rolls back any partial changes.

# Chapter 8

# Problem Solving Approach and Assumptions

## 8.1  Problem Extension and Challenges

We extended the base problem by focusing on exactly-once semantics for record appends in a distributed environment. This extension required addressing challenges such as:

- **Network Failures**: Handling partial failures and ensuring atomicity.

- **Duplicate Requests**: Identifying and eliminating redundant operations.

## 8.2  Ideas Attempted and Lessons Learned

Several approaches were considered:

- **Optimistic Concurrency Control**: Initially attempted, but frequent conflicts led to high retry rates.

- **Centralized Transaction Manager**: Improved coordination but became a bottleneck under high loads.

- **Two-Phase Commit (2PC)**: Final choice due to its balance of simplicity and effectiveness in maintaining consistency.

## 8.3  Implementation Assumptions

Key assumptions made during implementation include:

- **Reliable Messaging**: Assumes that retries will eventually succeed within a bounded time.

- **Consistent State**: Chunk servers start in a consistent state after restarts.

# Chapter 9

# Benchmarking and Performance Analysis

## 9.1 Benchmarking Setup

We conducted benchmarking using the following configuration:

- **System Setup**: Cluster with 5 nodes (1 master, 4 chunk servers).

- **Network**: Simulated delays to evaluate fault tolerance.

- **Workload**: Varied file sizes and concurrent append operations.

## 9.2 Scalability Analysis

Results indicate that our system scales linearly up to 10 nodes. Key findings include:

- **Throughput**: Increased proportionally with the number of chunk servers.

- **Latency**: Minimal increase in latency for up to 10 concurrent clients.

## 9.3 Performance Comparisons

We compared our implementation with the base GFS model:

| Metric | Base GFS | Our System |
|---|---|---|
| Latency (ms) | Less | More |
| Throughput (ops/sec) | Less | More |
| Consistency Guarantee | At-least-once | Exactly-once |

Table 9.1: Performance Comparison

# Chapter 10

# Conclusion and Future Work

## 10.1  Conclusion

Our project successfully implements exactly-once semantics in GFS, achieving strong consistency and improved reliability. Key takeaways include:

- Enhanced data integrity.

- Simplified client logic.

- Improved fault tolerance.

## 10.2  Future Work

Future improvements could focus on:

- **Dynamic Scaling**: Automatically adjusting the number of chunk servers.

- **Optimized Commit Protocols**: Reducing overhead in two-phase commits.

- **Advanced Fault Detection**: Using machine learning to predict failures.

# Chapter 11

# Work Split and References

## 11.1 Team Contributions

- **Mayank Mittal (2022101094)**: Master server implementation, benchmarking, documentation.

- **Shivam Mittal (2022101105)**: Chunk server logic, network communication, testing.

## 11.2 References

- *The Google File System*, Sanjay Ghemawat, Howard Gobioff, Shun-Tak Leung.

- GFS White Paper.

- Distributed Systems Course Material, University of XYZ.