# Storage

ARTEM SLIUSARENKO          MIHAILO CVETKOVIC          EUGEN LUCCHIARI HARTZ

`artems | mihailoc | eugenlh @kth.se`

October 25, 2024

## Contents

## List of Acronyms and Abbreviations

**GFS**  Google File System

**RDBMS**  Relational Database Management System

# 1 Introduction

Distributed computing encounters ongoing challenges in data storage management, particularly as data volumes and the demand for higher data throughput grows. To address these issues, innovative storage systems like Cassandra [Lakshman and MalikLakshman and Malik2009], Bigtable [Chang, Dean, Ghemawat, et al.C Dynamo [DeCandia, Hastorun, et al.DeCandia et al.2007], and the GFS have been developed. Each offers unique solutions to enhance scalability, fault tolerance, and performance in distributed environments. These systems make it possible to allow different kinds of applications to efficiently handle enormous amounts of data while maintaining high availability and robust data integrity which is crucial for today's data-driven decision-making processes. Additionally, by optimizing data storage and retrieval operations, these platforms enable advanced data analytics and real-time processing capabilities which empowers businesses to fully utilize their data assets. This strategic advantage is crucial in a time where timely and accurate data analysis forms the backbone of competitive business strategies and operational efficiency.

# 2 Problem Description

The main challenge addressed by the papers on Cassandra, Bigtable, Dynamo, and GFS is about efficiently managing huge amounts of data across distributed systems while ensuring high availability, performance, and scalability. Conventional storage architectures have problems handling the scale required by modern applications, such as web services and real-time analytics which manage enormous volumes of data distributed worldwide. These papers address the need for large-scale operations at companies like Google and Amazon, which require continuous data availability despite frequent component failures. Traditional databases have problems scaling horizontally in a cost-effective manner and manual partitioning of data complicates deployment. The papers propose architectural advancements that respond to these challenges by developing systems capable of handling failures without disrupting data availability. In the following sections, it will be explained what kind of solutions the different papers offer to these described challenges and it will be evaluated to which extent they have successfully implemented these innovative storage architectures.

# 3 Contribution

## 3.1 Cassandra

The Cassandra paper describes how the system addresses the challenges of managing huge amounts of data distributed across multiple servers to ensure high availability and scalability. The key focus is on overcoming the limitations of traditional databases that cannot scale horizontally in a cost-effective way and the complexities involved with manual data partitioning. Cassandra introduces a decentralised architecture that avoids single points of failure, which makes it possible to handle extensive data amounts across broad geographical locations without sacrificing performance or data availability.

## 3.2 Bigtable

The Bigtable paper explains how the system addresses the challenge of managing extensive datasets across distributed environments while ensuring high performance, availability and scalability. A key aspect is its ability to scale horizontally across thousands of servers. The approach makes it possible to manage petabytes of data that support a diverse range of Google's products such as Google Earth and Google Analytics. Bigtable introduces a non-relational data model that offers dynamic control over data layout which significantly enhances flexibility and performance compared to traditional database systems. This model supports a wide variety of data types and uses which demonstrates its adaptability and utility in handling large-scale data efficiently.

## 3.3 Dynamo

The Dynamo paper discusses Amazon's approach to achieving high availability and reliability at a massive scale. Amazon's e-commerce platform consists of hundreds of services hosted on a global infrastructure of tens of thousands of servers, across multiple data centers, where failures are frequent. These services are either stateless or stateful, with the latter relying on persistent storage. Traditionally Relational Database Management Systems (RDBMSs) are used to store application state, but they are inefficient for many services that only need to store and retrieve data by primary key. Relational databases offer unnecessary functionality, require expensive hardware, and prioritize consistency over availability. Dynamo is a key-value storage system that addresses these inefficiencies with a highly available key-value storage system, offering a simple interface, efficient resource usage, and scalability to meet growing data and request demands. To achieve this level of availability, Dynamo sacrifices consistency under certain failure scenarios, using techniques like object versioning and application-assisted conflict resolution, providing developers with a flexible interface while maintaining the scalability and reliability of the platform.

## 3.4 GFS

The [MittalMittal2015]GFS addresses the challenge of storing and managing vast amounts of data on commodity hardware, a necessity for Google's data-heavy operations. GFS was designed to handle large-scale distributed data processing tasks such as web crawling and indexing, which require massive datasets to be stored and accessed efficiently. A key contribution of GFS is its fault-tolerant design that anticipates frequent hardware failures and handles them seamlessly by replicating data across multiple servers.

Additionally, GFS introduces an efficient mechanism for handling very large files by dividing them into chunks, which are stored across multiple nodes. This approach enables parallel processing of large datasets and minimizes the impact of individual node failures. GFSs architecture also allows it to sustain high throughput, making it ideal for applications like MapReduce, which require continuous, large-scale data processing. By utilizing inexpensive hardware and focusing on availability and performance, GFS enables Google to manage its vast infrastructure efficiently, serving as the backbone for many of its core services

# 4 Solution Overview

## 4.1 Cassandra

Cassandra provides a distributed storage solution that manages structured data across numerous servers. It moves away from a traditional relational data model to a simpler, dynamic model that improves the flexibility of data layout and format. This system is designed to support high write throughput while maintaining efficient read operations, which is essential for applications like Facebook's Inbox Search feature.

The system's design integrates advanced strategies such as data partitioning using consistent hashing, extensive replication across multiple data centers, and advanced failure management. Unlike traditional systems that treat hardware failures as anomalies, Cassandra treats such failures as commonplace and typical occurrences. Therefore, the mechanisms are designed to ensure reliability and continuous availability.

This proactive approach to failure management, combined with its scalable architecture, makes Cassandra an exemplary solution for modern data-intensive applications that require resilience, agility, and operational efficiency.

## 4.2 Bigtable

Bigtable presents a solution designed to handle large-scale structured data across distributed systems. Its architecture is built around a sparse, distributed, persistent multi-dimensional sorted map, indexed by row key, column key, and timestamp, where each value in the map is an uninterpreted array of bytes. This design

allows for efficient storage and retrieval of huge amounts of data with minimal latency. This approach is ideal to support real-time applications.

Bigtable utilizes a combination of in-memory and on-disk data with optimizations for both read and write operations, including SSTable storage and a robust caching system for improved read performance. It also incorporates a compaction mechanism to maintain efficient data storage and retrieval.

This system structure enables Bigtable to meet the demands of Google's diverse applications and therefore offers a scalable, flexible, and high-performance solution for managing huge datasets.

## 4.3   Dynamo

Amazon's Dynamo database was specifically designed to meet the demanding needs of large-scale e-commerce operations by providing high availability and fault tolerance. In e-commerce, even brief outages can lead to significant financial losses and erode customer trust. Dynamo addresses these concerns by prioritizing availability over strict consistency. It is based on an eventual consistency model, which means that while data may not be immediately synchronized across all replicas, it becomes consistent over time, allowing services to remain operational during failures.

Dynamo's decentralized architecture ensures there are no single points of failure. Data is distributed across multiple nodes using consistent hashing, which balances the load and ensures that no single node becomes a bottleneck. Additionally, Dynamo replicates data across several nodes, enabling it to handle server or data center failures without disrupting service. In case of network partitions or node outages, Dynamo uses techniques like hinted handoff and anti-entropy to ensure that data remains available and can be reconciled when the system is fully restored.

Moreover, its scalability allows Amazon's services to dynamically scale up or down depending on the request load. This combination of high availability, fault tolerance, and scalability makes Dynamo an ideal solution for Amazon's complex, high-traffic e-commerce platform.

## 4.4   GFS

The GFS was designed to manage large-scale data storage needs while providing fault tolerance, high throughput, and scalability. GFS utilizes commodity hardware, meaning the system is built to expect frequent failures, yet still ensures data availability and reliability through mechanisms such as data replication. The system distributes large files across many machines, which allows efficient access and management of data, a critical feature for Google's data-intensive services such as web indexing and search.

# 5   Achievements

Each of these distributed storage systems—Cassandra, Bigtable, Dynamo, and GFS has significantly contributed to advancements in managing large-scale data. Cassandra's decentralized model has empowered companies like Facebook to handle vast amounts of structured data efficiently. Bigtable's flexible and scalable architecture supports a variety of Google's applications, from Google Maps to YouTube. Dynamo's eventual consistency model ensures uninterrupted services for Amazon, even during failures. Lastly, GFS's design, which capitalizes on inexpensive hardware, powers some of the largest data processing tasks in Google's infrastructure, providing the foundation for systems like MapReduce.

# 6   Comparative Evaluation: Strengths and Weaknesses

## 6.1   Strengths

Each system excels in different areas. Cassandra's decentralized architecture provides robust fault tolerance and high scalability, making it ideal for geographically distributed applications. Bigtable's efficient handling of large datasets enables rapid data access and processing, supporting a wide range of services. Dynamo's

emphasis on availability ensures that Amazon's platform remains operational despite failures, a critical feature for large-scale e-commerce operations. GFS stands out for its fault tolerance and ability to handle immense data processing tasks with cost-effective hardware.

## 6.2   Weaknesses

However, these systems also have weaknesses. Cassandra's focus on availability can sometimes compromise consistency, making it less suitable for applications requiring strong consistency guarantees. Bigtable's complexity in handling certain non-Google workloads may make it less adaptable outside its specific use cases. Dynamo sacrifices consistency for availability, which may introduce challenges in applications where data correctness is crucial. GFS, while efficient, is optimized for large file handling and may struggle with smaller, high-velocity data operations typically seen in modern applications like social media.

# 7   Conclusion

These distributed systems have revolutionized the way we handle large datasets in distributed environments. Each system provides unique solutions to operations involving performance, availability, and scalability. Although each system has its strengths and weaknesses, collectively they show how diverse the approaches are when handling large datasets. These innovations have not only empowered companies like Google, Facebook, and Amazon but have also set new benchmarks for the future of data storage and distributed computing.

# References

[Chang, Dean, Ghemawat, et al.Chang et al.2006] Chang, F., J. Dean, S. Ghemawat, et al. (2006). Bigtable: A distributed storage system for structured data. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI)*, pp. 205–218. USENIX Association.

[DeCandia, Hastorun, et al.DeCandia et al.2007] DeCandia, G., D. Hastorun, et al. (2007). Dynamo: Amazon's highly available key-value store. In *Proceedings of the 21st ACM Symposium on Operating Systems Principles (SOSP)*, pp. 205–220. ACM.

[Lakshman and MalikLakshman and Malik2009] Lakshman, A. and P. Malik (2009). Cassandra: A decentralized structured storage system. In *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP)*, pp. 35–50. ACM.

[MittalMittal2015] Mittal, A., J. V. . A. T. (2015). Google file system and hadoop distributed file system-an analogy. In *Innov Adv Comput Sci)*, pp. 29–43. Int J Innov.