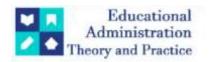
## **Educational Administration: Theory and Practice**

2024, 30(11), 406-410 ISSN: 2148-2403 https://kuey.net/

**Research Article** 



# Method Proposed for Avoidance of Data Replication in Distributed Storage System

Prof. Swati D. Ghule1\*, Dr. Anup Girdhar2

<sup>1\*</sup>Assistant Professor, P.E. S. Modern College of Engineering, Pune – 5, Email: swati.ghule@moderncoe.edu.in <sup>2</sup>Guide TMV, Email: anupgirdhar@gmail.com

**Citation:** Prof. Swati D. Ghule et al. (2024), Method Proposed for Avoidance of Data Replication in Distributed Storage System, *Educational Administration: Theory and Practice*, 30(11) 406-410
Doi: 10.53555/kuey.v30i11.8422

#### **ARTICLE INFO**

### **ABSTRACT**

Data replication is the process of copying and maintaining files in multiple locations to ensure consistency, reliability, and accessibility. Data replication in the cloud offers many benefits, but it also comes with certain drawbacks like increase storage cost. The process of replication can consume significant resources, including CPU, memory, and network bandwidth, which might impact the performance of primary systems or applications. Also Replicating data across different regions or cloud providers may expose it to additional security risks if proper encryption and access controls are not implemented. As the volume of data grows, the replication process might become more complex and resource-intensive, requiring additional infrastructure or optimization to handle the increased load [1].

The study consider example of GFS and HDFS which are two most used distributed file systems for dealing with huge clusters where big data lives [3]. Avoiding data replication can be advantageous in terms of cost efficiency, resource utilization, and system simplicity. Efficient resource utilization in cloud storage involves optimizing how storage resources are allocated, used, and managed to achieve cost savings, performance improvements, and operational efficiency. The study suggest HoneyBee algorithm for resource optimization. HoneyBee is a cloud-native storage system designed to optimize data placement and resource utilization in distributed storage environments.

Keywords: GFS, HDFS, HoneyBee, Distributed Storage, Resource Utilization.

#### **Introduction:**

Data replication is widely used in various contexts, including database management systems, cloud services, and distributed file systems, to ensure data integrity, improve performance, and support disaster recovery efforts.

Data replication in the cloud offers many benefits, but it also comes with certain drawbacks. Here are some common challenges and drawbacks associated with cloud data replication:

- 1. Cost:
- Storage Costs: Replicating data across multiple locations or regions can increase storage costs, as additional storage space needs additional cost.
- o **Bandwidth Costs**: Transferring data between locations can incur significant bandwidth charges, especially if the replication involves large volumes of data or frequent updates [1].
- 2. Complexity:
- Management Overhead: Managing and configuring replication policies, schedules, and conflicts can be complex and may require specialized tools or expertise.
- **Data Consistency**: Ensuring data consistency across replicas, especially in asynchronous replication scenarios, can be challenging. Data may become inconsistent or outdated if replication lags [2].
- 3. Security Concerns:
- **Data Exposure**: Replicating data across different regions or cloud providers may expose it to additional security risks if proper encryption and access controls are not implemented.
- Compliance Issues: Different regions may have different regulatory requirements for data storage and protection. Ensuring compliance across replicated data can be challenging.

- 4. Complex Failures:
- Replication Failures: Replication processes themselves can fail due to network issues, configuration
  errors, or software bugs, potentially leading to data inconsistencies or loss.
- **Conflict Resolution**: In scenarios where data can be written to multiple replicas, resolving conflicts (e.g., simultaneous updates) can be complex and may require manual intervention.
- 5. Performance Impact:
- **Resource Utilization**: The process of replication can consume significant resources, including CPU, memory, and network bandwidth, which might impact the performance of the primary systems or applications [1].
- o **Data Synchronization**: Keeping data synchronized in real-time or near real-time can be resource-intensive and might affect the performance of the primary system [2].
- 6. Scalability Issues:
- **Scaling Challenges**: As the volume of data grows, the replication process might become more complex and resource-intensive, requiring additional infrastructure or optimization to handle the increased load.
- 7. Backup and Recovery Complexity:
- **Recovery Time**: In case of a failure, recovering data from replicated sources might involve complex processes, especially if multiple replicas are involved.

Addressing these drawbacks typically involves careful planning, including optimizing replication strategies, leveraging cloud service features, and implementing robust monitoring and management practices. This study mainly focus on resource optimization.

The GFS has two replicas: Primary replicas and secondary replicas. A primary replica is the data chunk that a chunk server sends to a client. Secondary replicas serve as backups on other chunk servers. User can specify the number of replicas to be maintained. The HDFS has an automatic replication rack based system [3].

## **Need of Study:**

Avoiding data replication in distributed systems can be desirable in certain scenarios due to the following reasons:

- 1. Cost Efficiency
- 2. Resource Utilization
- 3. Consistency and Complexity
- 4. Performance Considerations
- 5. Data Freshness
- 6. Simplicity

In GFS and HDFS also need to take care of the following challenges which shows relation with above mentioned list

- 1. Increased Storage Cost
- 2. Write Latency
- 3. Network Overhead
- 4. Complexity in management
- 5. Potential for imbalanced data [6]

Here in this study the main focus has been given on Resource optimization.

#### **Resource Optimization**

- Compute and Storage Efficiency: Maintaining multiple replicas of data requires additional resources
  for storage and possibly for maintaining consistency. Avoiding replication can free up these resources for
  other tasks or reduce the overhead on the system.
- **Power Consumption**: Less replication means fewer active disks and servers, which can result in lower power consumption and a smaller environmental footprint.

With the efficient resource optimization, it will be helpful in resolving the issues like increase storage and bandwidth cost, scalability of data, complexity in recovering the data.

### **Problem Statement:**

Efficient resource utilization in cloud storage involves optimizing how storage resources are allocated, used, and managed to achieve cost savings, performance improvements, and operational efficiency.

Efficient resource utilization in cloud storage offers several significant benefits:

1. Cost Efficiency: Optimizing storage usage helps minimize costs by reducing the

1. **Cost Efficiency**: Optimizing storage usage helps minimize costs by reducing the need for excess storage capacity and lowering expenses associated with data transfer and retrieval.

- 2. **Scalability**: Efficient utilization allows businesses to scale storage needs seamlessly, enabling them to easily accommodate growing data volumes without over-provisioning resources.
- **3. Performance Optimization**: Proper resource management can enhance access speeds and reduce latency, leading to improved application performance and user experiences.
- **4. Enhanced Data Management:** Efficient resource utilization promotes better organization of data, making it easier to manage, retrieve, and analyze information.
- **5. Reduced Waste**: Minimizing unused or underutilized storage resources helps decrease waste, aligning with sustainable practices and reducing the environmental impact.
- **6. Improved Security**: Effective utilization often includes better data governance and access controls, leading to enhanced security measures for sensitive information.
- **7. Faster Backups and Restores**: Efficiently managed storage can streamline backup processes and reduce restore times, ensuring data integrity and availability.
- **8. Data Redundancy and Recovery**: Optimized resource use often incorporates better redundancy strategies, improving data recovery options in case of loss or corruption.
- 9. Increased Flexibility: Organizations can more easily adapt their storage resources based on changing needs, whether that involves increasing capacity during peak times or downsizing when demand decreases.
- 10. Better Insights and Analytics: Efficient storage utilization can facilitate better analytics capabilities, allowing organizations to gain insights from their data more quickly and effectively.

During data replication both GFS and HDFS depends on a single master node which at times proves to be a failure point [5].

## **Hypothesis:**

HoneyBee is a cloud-native storage system designed to optimize data placement and resource utilization in distributed storage environments. It uses a different approach from traditional data replication strategies to achieve efficiency.

Here's a comparison of how HoneyBee's approach contrasts with data replication strategies:

Feature	HoneyBee	Data Replication
	Optimizing data placement and resource utilization	Ensuring data availability through multiple copies
	Achieved through advanced placement and erasure coding	Achieved through multiple replicas of data.
	Lower, as it uses erasure coding and optimized placement	Higher, due to maintaining multiple replicas
	Maximized by balancing load and adapting to node conditions	May be less efficient due to fixed replication policies
Scalability		Scaling requires managing additional replicas
Complexity	Can be complex to manage data placement and erasure coding	Simpler to understand, but can be resource-intensive

#### **Research Methodology:**

The Honeybee optimization algorithm on a cloud storage system to optimize the resources and minimize the transaction response and execution time.

The HoneyBee algorithms is as follows

Input: List\_statehost , List\_stateVM

Output: VM(j)

// return the number (i) of host which have minimum processing time

// i is the number of specified host

1: i ← -1

2: minPT← Integer.MAX\_VALUE

3: For each host(i) in List statehost

4: If host(i) available then

5: If (PThost(i) < minPT) then

6: minPT = PThost(i)

7:  $i \leftarrow$  number of the current host

8: End if

9: End if

```
10: End for
// return the number (j) of VM which has minimum count of requests
// j is the number of specified VM

11: j ← -1

12: mincount← Integer.MAX_VALUE

13: For each VM(j) in List_VMhost(i)

14: If VM available then

15: If (Count_REQVM(j) < mincount) then

16: mincount = Count_REQVM(j)

17: j ← number of the current VM

18: End if
```

Honeybee is useful to improve the efficiency of data storage in distributed systems like GFS (Google File System) and HDFS (Hadoop Distributed File System) by minimizing data replication.

#### **Result and Discussion:**

Benefits of Using HoneyBee for Resource optimization.

- 1. **Cost Efficiency**: By reducing the need for multiple replicas and optimizing data placement, HoneyBee can lead to significant cost savings in terms of storage and operational expenses.
- **2. Improved Performance**: HoneyBee's approach to balancing load and optimizing data access can result in better performance compared to traditional replication methods, where performance might be affected by the overhead of managing multiple replicas [2].
- **3.** Enhanced Fault Tolerance: Using erasure coding and intelligent data placement, HoneyBee can achieve fault tolerance with less redundancy, thereby reducing storage overhead while maintaining reliability [2].
- **4. Dynamic Adaptation**: HoneyBee's ability to dynamically adapt to changes in the storage environment allows for more flexible and efficient resource utilization compared to static replication strategies [7].

By implementing these strategies, Honeybee enhances the efficiency of storage systems like GFS and HDFS, leading to reduced storage costs, improved performance, and optimized resource usage.

## **Findings:**

- 1. Data Placement Optimization: HoneyBee intelligently places data across different storage nodes based on factors like access patterns, node availability, and resource utilization. It minimizes the number of replicas needed by efficiently using available resources.
- 2. Resource-Aware: It takes into account the current state and performance of storage nodes to make decisions about where data should be placed. This helps in balancing the load and ensuring efficient use of resources.
- **3. Elasticity**: HoneyBee is designed to scale with the addition of new nodes or changes in resource availability, adapting dynamically to the storage environment.
- **4. Fault Tolerance without Over-Replication**: Instead of relying solely on replication for fault tolerance, HoneyBee uses techniques like erasure coding, which can provide similar levels of redundancy with less storage overhead [2].

In GFS and HDFS following features can be achieved using HoneyBee:

- 1. Smart Replication Strategy
- 2. Dynamic Data Placement
- 3. Load Balancing
- 4. Data locality awareness
- 5. Fault Tolerance [4]

#### **Conclusion:**

HoneyBee offers an alternative approach to traditional data replication by focusing on optimizing data placement and resource utilization, often through techniques like erasure coding. It can be particularly useful in scenarios where cost efficiency and resource management are critical, and where traditional replication strategies might be too resource-intensive or inflexible. By leveraging HoneyBee's capabilities, organizations can achieve efficient storage management while still ensuring fault tolerance and performance.

## **Future Scope for Further Research:**

Efficient resource utilization in cloud storage can be achieved through a combination of data tiering, compression, deduplication, automated lifecycle management, right-sizing, object storage, caching, backup management, monitoring, encryption, and hybrid/multi-cloud strategies. Each method contributes to cost

savings, improved performance, and better overall resource management. Implementing these strategies requires a thoughtful approach, considering the specific needs and constraints of the organization and the characteristics of the data being managed.

#### References

- [1] Rambabu D, Govardhan A, "Survey on data replication in cloud systems", Article in Web Intelligence January 2024, 83–109, DOI: 10.3233/WEB-230087
- [2] Mokadem et. al., "A review on data replication strategies in cloud systems", International Journal of Grid and Utility Computing, 2022, 13 (4), pp.347-362. 10.1504/IJGUC.2022.125135. Hal-03828293
- [3] Nader Gemayel, "Analyzing Google File System and Hadoop Distributed File System", Research Journal of Information Technology · March 2016, DOI: 10.3923/rjit.2016.66.74
- [4] Sanjay Ghemawat, "The Google File System", SOSP'03, October 19–22, 2003, Bolton Landing, New York, USA. ACM 1-58113-757-5/03/0010
- USA. ACM 1-58113-757-5/03/0010
  [5] Zahid Ullah et. al., "Analytical Study on Performance, Challenges and Future Considerations of Google File System", International Journal of Computer and Communication Engineering, Vol. 3, No. 4, July 2014, DOI: 10.7763/IJCCE.2014.V3.336
- [6] M. A. Ahmed et al., "Dynamic Replication Policy on HDFS Based on Machine Learning Clustering", IEEE Access, VOLUME 11, 2023, pp. 18551-18559.
- [7] Mohammad Khalilzadeh, "A Honey Bee Swarm Optimization Algorithm for Minimizing the Total Costs of Resources in MRCPSP", Indian Journal of Science and Technology, Vol 8(11), DOI: 10.17485/ijst/2015/v8i11/71405, June 2015