

# Duplication Checking With Encrypted Data Storage in Cloud

<sup>1</sup>Dr.Venkatesh S, <sup>2</sup>Malarkodi M

<sup>1</sup>Assistant Professor, <sup>2</sup>Student

Department of CSE

Jeppiaar Engineering College Chennai India

**Abstract:** Cloud computing has reached a high potential leading to a strong productive phase. This means that many of the main problems with cloud computing are directed at the level of cloud computing with full commercial exploitation. However, concerns over data security still prevent many users from moving data to remote storage. Squeezing client side data in particular ensures that multiple uploads of the same content consume only network bandwidth and single upload space. The theory is that external data may require different levels of protection, depending on how popular it is: content shared by multiple users. Then present the idea of a novel that divides the data according to its popularity. Based on this idea, design a duplicate removal system that ensures semantic retention of unwanted data and provides better security and better storage and benefits of popular data bandwidth. In this way, data loss may apply to the popular data storage system. The extraction method helps to reduce storage requirements through the compatibility testing process. This allows organizations to store too much data on the same system and extend disk purchase intervals automatically. With the benefit of speed, organizations can save data on disk costs effectively. In this way, deleting data duplicates may apply to popular data, while statistically secure encryption protects unwanted content. We can use the backup system during blockchain and analyze the access system on a regular basis.

**Index Terms:** Cloud, storage, duplicate

## I. INTRODUCTION

Data extraction is a process in which the storage provider stores only one copy of a file (or part of it) for multiple users. There are four different withdrawal strategies, depending on whether duplication occurs on the client side (i.e. before uploading) or on the server side, and whether it occurs at the block level or file level. Capturing data on the client side is more advantageous than the server side as it ensures that multiple uploads of the same content consume only network bandwidth and single upload space. The aim of this project is to reduce the amount of virtual space used by removing the same data blocks and using metadata to link logical and tangible data copies. In public cloud, repository capacity of the storage platform is not displayed to the user. The reduction process helps reduce the need for storage through a process of matching similarities. This allows organizations to store as much data as possible on the same system and to extend disk purchase intervals automatically. With the benefit of speed, organizations can save data on disk costs effectively. If data segregation occurs in sources, there is no need to transfer data over the network, thus eliminating unnecessary use of network bandwidth.

## II. EXISTING SYSTEM

In the present system design a client protocol for the client release on the side called KeyD without an independent key control server using an identity-based streaming process (IBBE). Users only interact with a cloud service provider (CSP) during the data download and download process. Security analysis shows that KeyD ensures the confidentiality of data and the security of the flexible key, and effectively protects the privacy of the property at the same time. Complete and detailed performance comparisons show that the system performs better transactions between storage costs, communication and higher calculation. Convergent encryption (CE) [3], encrypting a copy of the data with a flexible key obtained by calculating the cryptographic hash of the content of the copy of the data itself and thus generating the same ciphertext text in the same transparent text, brings hope to detect duplication over time. ensuring data confidentiality. This flexible encryption feature allows CSP to duplicate encrypted data. In particular, users secretly copy copies of their data using the corresponding keys and external encrypted data to CSP. They just need to keep the conversion keys in place so they can later recover data. However, the amount of keys they convert increases according to the number of copies of the data as the copy of the data corresponds to the conversion key. Each data owner encrypts all copies of his or her data using the corresponding conversion keys and continues to encrypt these converter keys using his or her key. Both copies of encrypted data and conversion keys are stored in the cloud, and users store their key keys and metadata about exported data locally. Although encrypted data can be duplicated by the cloud, the encryption of encrypted key keys will increase according to the number of users.

## III. DRAW BACK OF EXISTING SYSTEM

- Recurrence checks only with file name and not file content.
- Storage is very expensive.
- There is no proper recovery plan during the blockade.
- Errors in copyright-based encryption

#### IV. LITERATURE SURVEY:

Hybrid Priorized Data Extraction to address the end-to-end system shared applications running on virtual machines or compounds by combining internal and post-processing processes to minimize exactly. In the on-line retrieval phase, HPDedup provides a way to save fingerprints that limit the duration of duplicates in the stream of data from various VMs or applications and prioritize the repository allocation of these streams based on scale. HPDedup also allows for a different deduplication drawback limit for streaming based on its space space to minimize disk partition. The post-processing phase removes their unsealed fingerprints due to the weak temporary location from the disks. While it offers many benefits in improving the efficiency of app sharing, the increase in the number of applications from different users sharing the same machine raises challenges in reducing basic storage.

The detailed description of the LGC is useful both in understanding its shortcomings and in the fact that this was a version that was used within the retail product for many years. We note that since the LGC came into use, there have been several publications describing other GC programs in detail and we regard the technical contributions of this paper as the details leading to the new and highly developed PGC sub-system. The PGC has undergone a transformation, which begins with a change in the device's access system and further enhanced to reduce memory usage and to avoid the need for excess data. We compare the LGC with the previous use of PGC, which has been used on client sites for a long time, and the "new GC for the new phase" (PGC +), which includes additional configurations. In short, the contributions of this paper include

Product discusses a tradeoff between small and large chunk algorithms with two levels of chunking. It transmits and stores data based on relatively large superchunks, consisting of a number of smaller pieces that form units of information for duplication. Second, it combines contributions of two novels that enable it to (i) accomplish superchunk assignments with minimal CPU and memory requirements, and (ii) balance load on cluster nodes without disrupting duplicate operations.

The Administration Interface module provides a visual interface for system administrators to configure SAR design parameters and monitor SAR operating time status. The Selected Promotion Module is responsible for monitoring and identifying episode size, sequence, and reference statistics for unique data segments processed by the Data Deduplicator module, as well as the popularity of access to separate data segments. Based on the information, SAR ultimately selects unique "hot" data components from the SSD system.

This is done by manipulating each image file such as by-stream and bypassing borders using a "continuous" function (for stitching) or Rabin fingerprints (for stitching different sizes). A chunk is simply data between two boundaries; there are vague limits at the beginning and end of each image file. Chunks are identified by their SHA1 hash, which is calculated using SHA1 over the content of the episode. We think pieces with the same chunk ID are the same; we do not make byte-by-byte comparisons to ensure that the pieces are the same

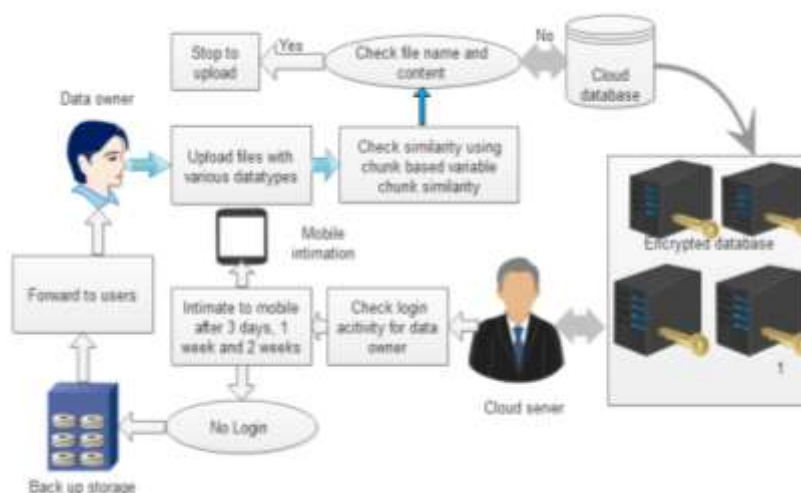
#### V. PROPOSED SYSTEM

Data compression is a process in which the storage provider stores only one copy of a file for a few of its users. There are four different compression techniques, depending on whether compression occurs on the client side (i.e. before uploading) or on the server side, and whether compression occurs at a block level or file level. Pressure is very beneficial when configured on the client side, as it also saves upload bandwidth. For these reasons, compression is an important provider of a number of popular and efficient storage services that provide cheap, remote storage to the wider community by performing pressure on the client side, thus saving both network bandwidth and storage costs. The goal of the system is to ensure the confidentiality of the data without losing the congestion benefit. Confidentiality should be guaranteed for all files, including predictable ones. The security of the entire system should not depend on the security of one part (one point of failure), and the level of security should not fall when one part is in danger. We view the server as a trusted component in terms of user authentication, access control and additional encryption.

#### VI. ADVANTAGES

- Powerful updates can be done in cloud Storage.
- File and file content analyzed.
- Save a disused storage space.
- Proper alert system is maintained by the cloud server.
- Provide advanced pressure on the file server system.

## VII. ARCHITECTURE DIAGRAM



## VIII. METHODOLOGY

On a computer, data compression is a special way to compress data to duplicate duplicate data. Synonyms are synonymous with intelligent dynamics (data) and single-sample storage (data). This process is used to improve storage usage and can be used in network data transfer to reduce the number of bytes to be sent. In the compression process, different pieces of data, or byte patterns, are identified and stored during the analysis process. As the analysis progresses, some episodes are compared to a backup copy and whenever the same thing happens, the obsolete part is replaced by a smaller reference that identifies the archive. In this module, we can scan files using the file name that contains the contents of the file. Encrypted files are fragmented. The service provider checks the pieces during file uploads. The data owner only uploads the original file to save storage space in the cloud system. We can press on text file, document file and image files.

## IX. SOFTWARE INTERFACE

PHP is a web embedded HTML language. This means that the PHP code can be embedded in the HTML of a Web page. When a PHP page is accessed, the PHP code is read or "transferred" by the server the page resides in. Outputs from PHP activities on the page are usually returned as HTML code, which can be read by the browser. Because the PHP code is converted to HTML before the page is loaded, users are not able to view the PHP code on the page. This makes PHP pages secure enough to access the site and other secure information.

Most PHP syntax is borrowed from other languages such as C, Java and Perl. However, PHP has many different features and functions. The goal of language is to allow web developers to write customizable pages quickly and easily. PHP is also excellent for creating web-driven Web sites.

Hypertext refers to files that are linked together using links, such as HTML files (Markup HyperText Language). Pre-processing to make instructions that change the output. Below is a demonstration of the differences between HTML and PHP files.

## X. RESULT



cloud owner

Welcome to Owner Register

Owner Register

Name:

Email Address:

Password:

Confirm Password:

☐ I agree to the Terms and Conditions

CLOUD SERVER

HOME OWNER DETAILS OWNER DETAILS OWNER DETAILS OWNER DETAILS OWNER DETAILS CONTACT

Owner Space Details

ID No.	Name	Phone Number	Address	Email ID	Action
1	John	9876543210	123 Main St	john.doe@example.com	<a href="#">Edit</a> <a href="#">Delete</a>

cloud owner

Owner Register

Name:

Email Address:

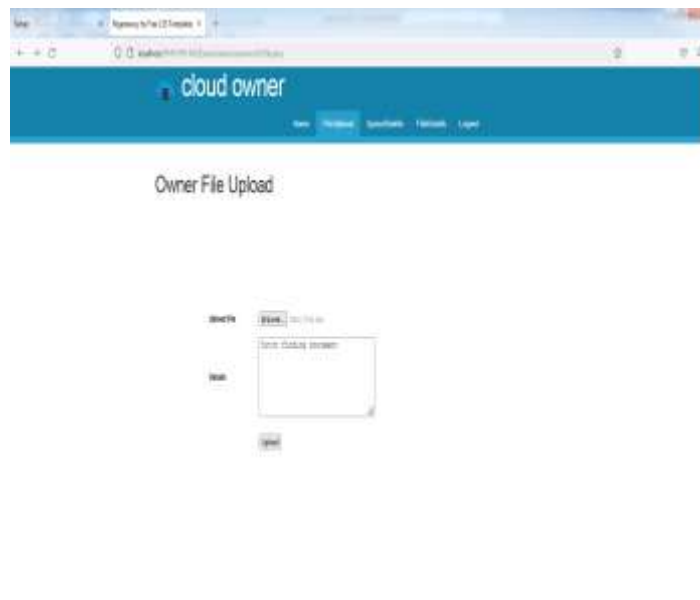
Password:

Confirm Password:

☐ I agree to the Terms and Conditions

Storage Space Details

ID No.	Name	Phone Number	Address	Email ID	Storage Space	Action
1	John	9876543210	123 Main St	john.doe@example.com	1 GB	<a href="#">Edit</a> <a href="#">Delete</a>



## XI.CONCLUSION

Security of tag and integrity has been achieved. We have developed our compression systems using a private sharing scheme and have shown that we deliver less coding / coding compared to higher network transfers for normal upload / download operations. In this work, we have identified a new privacy challenge during data access to cloud computing in order to gain access to the privacy of the archive access to the same files. Authentication was established to ensure data confidentiality and data integrity. Data anonymity is achieved as the folded values are exchanged during the transfer. User privacy is enhanced by the privacy requests of the cloud server for information about the users' accessibility requirements. Backup system to upgrade the detected system to avoid blocking and restore the number of unused spaces in the cloud system.

## REFERENCES

- [1] Wu, Huijun, Chen Wang, Yinjin Fu, Sherif Sakr, Liming Zhu, and Kai Lu. "Hpdedup: A hybrid prioritized data deduplication mechanism for primary storage in the cloud." arXiv preprint arXiv:1702.08153 (2017).
- [2] Douglass, Fred, Abhinav Duggal, Philip Shilane, Tony Wong, Shiqin Yan, and Fabiano Botelho. "The logic of physical garbage collection in deduplicating storage." In 15th {USENIX} Conference on File and Storage Technologies ({FAST} 17), pp. 29-44. 2017.
- [3] Frey, Davide, Anne-Marie Kermarrec, and Konstantinos Kloudas. "Probabilistic deduplication for cluster-based storage systems." In Proceedings of the Third ACM Symposium on Cloud Computing, p. 17. ACM, 2012.
- [4] Mao, Bo, Hong Jiang, Suzhen Wu, Yinjin Fu, and Lei Tian. "Read-performance optimization for deduplication-based storage systems in the cloud." ACM Transactions on Storage (TOS) 10, no. 2 (2014): 6.
- [5] Jin, Keren, and Ethan L. Miller. "The effectiveness of deduplication on virtual machine disk images." In Proceedings of SYSTOR 2009: The Israeli Experimental Systems Conference, p. 7. ACM, 2009.