

# A Survey Paper on Removal of File Duplication at Block Level in a Hybrid Cloud

Pritam Dalvi <sup>[1]</sup>, Vishnu Wagh <sup>[2]</sup>, Shubham Kale <sup>[3]</sup>, Krushna Khandagale <sup>[4]</sup>

Dipalee A. Chaudhari <sup>[5]</sup>, Soniya Bastawade <sup>[6]</sup>

Student <sup>[1],[2],[3] & [4]</sup>, Assistant Professor <sup>[5] & [6]</sup>

D Y Patil College of Engineering Akurdi

Pune - India

## ABSTRACT

Data deduplication is an effective compression technique to eliminate duplicate copies of repeating data. It has been widely used in cloud storage to reduce the amount of storage space and save bandwidth. For the protection of sensitive data while supporting deduplication, the convergent encryption technique is being used to encrypt the data before outsourcing it to cloud storage. In order to make system more secure, the different privileges of users are again considered while checking duplicate content. But the problem which occurs in this approach is that even if some contents of both files are different, it stores them as two different files which leads to reduction in cloud storage space. The solution for this problem is to perform block level deduplication. In this approach the file which is to be stored on cloud storage is divided into number of different blocks based on contents and deduplication is performed on these blocks.

**Keywords:-** Deduplication, authorized duplicate check, confidentiality, hybrid cloud

## I. INTRODUCTION

Cloud computing technique which is most widely used today. In that, computing is done over the large communication network like Internet. It is an important solution for business storage in low cost. Cloud computing provide vast storage in all sector like government, enterprise, also for storing our personal data on cloud. Without background implementation details, platform user can access and share different resources on cloud. The most important problem in cloud computing is that large amount of storage space and security issues. One critical challenge of cloud storage to management of ever-increasing volume of data. To improve scalability, storage problem data de-duplication [8] is most important technique and has attracted more attention recently. It is an important technique for data compression. It simply avoid the duplicate copies of data and store single copy of data. Data de-duplication take place in either block level or file level. In file level approach duplicate files are eliminate, and in block level approach duplicate blocks of data that occur in non-identical files. De-duplication reduce the storage needs by upto 90-95 % for backup application, 68% in standard file system. Important issues in data de-duplication that security and privacy to protect the data from insider or outsider attack. For data confidentiality, encryption is used by different user for encrypt there files or data, using a secret key user perform

encryption and decryption operation. For uploading file to cloud user first generate convergent key, encryption of file then load file to the cloud. To prevent unauthorized access proof of ownership protocol is used to provide proof that the user indeed owns the same file when de-duplication found. After the proof, server provide a pointer to subsequent user for accessing same file without needing to upload same file. When user want to download file he simply download encrypted file from cloud and decrypt this file using convergent key. To make information administration versatile in distributed computing, de-duplication has been an understood system and has pulled in more consideration as of late. Information de-duplication is a particular information pressure procedure for disposing of copy duplicates of rehashing information away. The procedure is utilized to enhance stockpiling usage and can likewise be connected to network information exchanges to diminish the quantity of bytes that must be sent. Rather than keeping various information duplicates with the same substance, de-duplication wipes out excess information by keeping one and only physical duplicate and alluding other repetitive information to that duplicate. De-duplication can occur at either the document level or the square level. For file level de-duplication, it wipes out copy duplicates of the same document. De-duplication can likewise happen at the piece

level, which dispenses with copy squares of information that happen in non-indistinguishable records.

## II. CLOUD COMPUTING SYSTEM

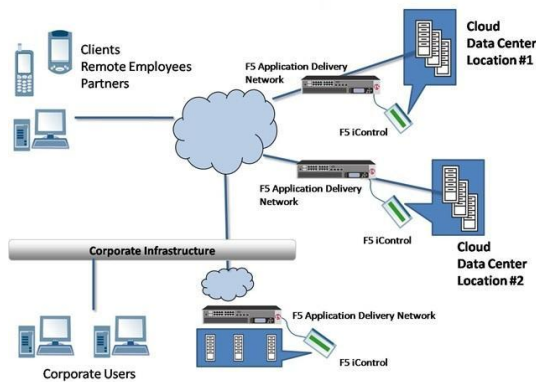


Fig. 1. Cloud Computing System

## III. TYPES OF CLOUD STORAGE

### Public Cloud:

Public clouds are made available to the public by a service provider who provide the cloud infrastructure. Most common public cloud providers are namely Google, Amazon and Microsoft own and operate their infrastructure and offer access all over the Internet. In this model, customer has no control or visibility over where the infrastructure is located. It is most important to consider that all the customers on public cloud storage share the same infrastructure with availability variances, limited configuration and security protections .Public Cloud customers generally benefit from economies of scale, because of the infrastructure costs are spread across all users which allows each and every client to operate on a low-cost, pay-as-you-go model. The main advantage of public cloud infrastructures is that public cloud storages are typically larger in scale than an in-house enterprise cloud, which in turn provides clients with seamless, on demand scalability. These clouds offer one of the greatest level of efficiency in shared resources; however, at the same time these storages are also more insecure than private cloud storages.

### Private Cloud:

Private clouds are those that are specially designed for a business. Many companies consider private clouds as a good starting point. They allow the organization to host applications, infrastructure and development environments in a cloud. Also addressing concerns regarding data security and

control that can arise in the public cloud environment. In general there are two types of private clouds: First and most common private cloud is a Private Cloud: This cloud also known as an Internal Cloud. It is hosted within an organizations data center. The benefits of a private cloud scalable, virtualized, flexible. Private cloud infrastructure are non-accepted. The type of private cloud is an externally hosted Virtual Private Cloud: It is being hosted by a third party Cloud Service Providers. This provider an exclusive private cloud environment .This takes responsibility for implementing, managing, and securing the Cloud infrastructure.

### Hybrid Cloud:

A hybrid cloud is best-of-breed. It combines the comfort level of a private cloud with the versatility and flexibility of the public cloud. Hybrid platforms use either public clouds or Hosted Virtual Private Clouds for some processes and applications. Both of these merge with on premises private clouds for high-security application environments. As with the private model, in a hybrid cloud, an organization may choose to continue to use their data center equipment and sensitive data secured on their own network. And like the public cloud, a hybrid model let an organization take advantage of a clouds accessibility, backup, scalability and disaster recovery. There are some of the limitations of the public cloud while still gaining many of the public clouds benefits.[6]

## IV. CLOUD SERVICES

Entity responsible for making services available are:

### Software as a Service(SaaS):

The SaaS is not suitable for applications that require real-time response or those for which data is not allowed to be hosted externally. The most likely candidates for SaaS are applications for which: Many competitors use the same product, such as email. Periodically there is a significant peak in demand, such as billing and payroll. There is a need for Web or mobile access, such as mobile sales management software. There is only a short-term need, such as collaborative software for a project.

### Platform as a Service(PaaS):

Platform-as-a-Service (PaaS) gives the capability to deploy consumer-created or acquired applications using programming languages and tools supported by the provider. The user does not manage or control the underlying cloud infrastructure, including network, servers, operating systems, or storage. The user has control over the deployed applications and, possibly,

over the application hosting environment configurations. Such services include session management, device integration, sandboxes, instrumentation and testing, contents management, knowledge management, and Universal Description, Discovery, and Integration (UDDI), a platform-independent Extensible Markup Language (XML)-based registry providing a mechanism to register and locate Web service applications. PaaS is not particularly useful when the application must be portable, when proprietary programming languages are used, or when the underlying hardware and software must be customized to improve the performance of the application. The major PaaS application areas are in software development where multiple developers and users collaborate and the deployment and testing services should be automated.

#### ***Infrastructure as a Service(IaaS):***

Infrastructure-as-a-Service (IaaS) is the capability to provision processing, storage, networks, and other fundamental computing resources; the consumer is able to deploy and run arbitrary software, which can include operating systems and applications. The consumer does not manage or control the underlying cloud infrastructure but has control over operating systems, storage, deployed applications, and possibly limited control of some networking components, such as host firewalls. Services offered by this delivery model include: server hosting, Web servers, storage, computing hardware, operating systems, virtual instances, load balancing, Internet access, and bandwidth provisioning. The IaaS cloud computing delivery model has a number of characteristics, such as the fact that the resources are distributed and support dynamic scaling, it is based on a utility pricing model and variable cost, and the hardware is shared among multiple users. This cloud computing model is particularly useful when the demand is volatile and a new business needs computing resources and does not want to invest in a computing infrastructure or when an organization is expanding rapidly.

## **V. RELATED WORK**

### ***Data de-duplication***

Data de-duplication is divided into many types. There is no best way to implement data de-duplication across an organization. For maximize the benefits, organizations may introduce more than one de-duplication strategy every decade. It is very simple to understand the backup and backup challenges, when selecting de-duplication as a solution. Data de-duplication has mainly divided into three forms. According to definitions, some forms of data de-duplication, such as compression, have been around for decades. Lately, single-

instance storage has introduced the removal of redundant files from storage environments and now, we have seen the introduction of sub-file de-duplication. Three types of data de-duplication are:

**Data Compression** The main use of Data compression is reducing the size of files. Data compression works within a file to identify and remove empty space that present in file as repetitive forms. This form of data de-duplication is only local to the file. Data Compression does not take into consideration of other files and data segments within those files. Data compression technique is occur many years ago, but being isolated to each single file, the benefits are limited as comparing to other forms of de-duplication. For example, data compression will not be effective in eliminating and recognizing duplicate files, but will independently compress each of the files.

**Single-Instance Storage** The main use of Single-Instance Storage is removing multiple copies of any file. Single-instance storage (SIS) environments are detect and remove redundant copies of identical files. After a file is stored in a single-instance storage system then, all those other references to same file, will refer to the original file, only single copy. Every time Single-instance storage systems compare the content of files to determine the incoming file is identical to an existing file in the storage system. Content-addressed storage is typically completed with single-instance storage functionality. In file-level de-duplication, it avoids storing files that are a duplicate of another file, many files that are considered unique or identical by single-instance storage measurement may have a more amount of redundancy within the files or between files. For example, take one small element (e.g., a new name inserted into the title slide of a file) for single-instance storage for two large files as being different and requiring them to be stored without further de-duplication.

**Sub-file De-Duplication** Sub-file de-duplication detects same or redundant data within and across files as opposed to finding identical files as in Single-Instance Storage implementations. With the help of sub-file de-duplication, redundant copies of data are detected and are eliminated even after the duplicated data exist, within different files. This form of de-duplication finds the unique data elements within an organization and detects when these elements are used or find within other files. As a result, sub-file de-duplication removes the storage of duplicate data in organization. Sub-file data de-duplication has many benefits even where files are not identical, but have data elements that are already recognized somewhere in the organization. Sub-file de-duplication implementation is mainly divided into two forms. First form is Fixed-length sub-file de-duplication uses an arbitrary fixed length of data to

search for the duplicate data within the files. If file is simple in design, fixed-length data segments miss many opportunities to produce redundant sub-file data. (Consider the case where an addition of a date or name or title is added to a document's title page the whole content of the document will shift, causes the failure of the de-duplication tool to detect same data). Second form is Variable-length implementations are usually not locked to any of arbitrary segment length. Variable-length implementations match data segment sizes to the naturally occurring duplication within files, vastly increasing the overall de-duplication ratio as compare to fixed length (In the example above, variable-length de-duplication will catch all duplicate segments in the document, no matter where the changes occur).

#### ***Data De-duplication types:***

File-level de-duplication: It is commonly known as single-instance storage, file-level data de-duplication compares a file that has to be archived or backup that has already been stored by checking all its attributes against the index. The index is updated and stored only if the file is unique, if not than only a pointer to the existing file that is stored references. Only the single instance of file is saved in the result and relevant copies are replaced by "stub" which points to the original file.

#### ***Block-level de-duplication:***

Block-level data de-duplication operates on the basis of sub-file level. As the name implies, that the file is being broken into segments blocks or chunks that will be examined for previously stored information vs redundancy. The popular approach to determine redundant data is by assigning identifier to chunk of data, by using hash algorithm for example it generates a unique ID to that particular block. The particular unique Id will be compared with the central index. In case the ID is already present, then it represents that before only the data is processed and stored before .Due to this only a pointer reference is saved in the location of previously stored data. If the ID is new and does not exist, then that block is unique. After storing the unique chunk the unique ID is updated into the Index. There is change in size of chunk as per the vendor. Some will have fixed block sizes, while some others use variable block sizes likewise few may also change the size of fixed block size for sake of confusing. Block sizes of fixed size may vary from 8KB to 64KB but the main difference with it is the smaller the chunk, than it will be likely to have opportunity to identify it as the duplicate data. If less data is stored than it obviously means greater reductions in the data that is stored. The only major issue by

using fixed size blocks is that in case if the file is modified and the de-duplication result uses the same previously inspected result than there will be chance of not identifying the same redundant data segment, as the blocks in the file would be moved or changed, than they will shift downstream from change, by offsetting the rest of comparisons.

#### ***Variable block level de-duplication:***

It Compares vary-ing sizes of data blocks that can reduce the chances of collision, stated Datalinkss Orlandini.[7]

## **VI. ENCRYPTION**

Encryption[4] in simple terms is defined as the process of obscuring information in order to make it unreadable to anyone without special knowledge, keys and/or passwords. This term is applicable to electronic signal, hard drive, message, and document. In computer context, It is the process of encoding messages, files, information etc in such a way that only authorized parties can access them. There are two terms which are commonly used in encryption those are plaintext and cipher text. The actual message which we want to send or process is known as plaintext while the message produced after implementation of encryption algorithm is known as cipher-text. Generally every algorithm produces unique key known as encryption key after implementation of algorithm. Without knowledge of encryption key it is impossible for anyone to understand the encrypted data. There are two most common approaches for encryption which are explained below:

#### ***Traditional Encryption Algorithm :***

Although it is known that data de-duplication gives more benefits, security and privacy concerns arise because the user's sensitive data is susceptible to both the outsider as well as insider attacks. So, while considering the traditional encryption techniques to secure the users sensitive data there are many issues are associated. Traditional encryption provides data confidentiality but it is not compatible with de-duplication. As in traditional encryption different users encrypt their data with their own keys. Thus, the identical data of the different users will lead to different cipher text which is making the data de-duplication almost impossible in this traditional approach.

Convergent Encryption algorithm: The convergent encryption [1] techniques are those which provide the data confidentiality to the users outsourced data stored on the public clouds. These techniques while providing the confidentiality to the data are

also compatible with the data de-duplication process. In this algorithm the encryption key is itself derived from the message. So it supports data de-duplication also, because the same file will give the same encryption key so it will generate the same cipher text irrespective of users which makes data de-duplication possible.

## **VII. CONCLUSION**

In this paper, the idea of authorized data deduplication was proposed to protect the data security by including differential authority of users in the duplicate check. In public cloud our data are securely store in encrypted format, and also in private cloud our key is store with respective file. There is no need to user remember the key. So without key anyone cannot access our file or data from public cloud.

## **REFERENCES**

- [1] Jin Li, Yan Kit Li, Xiaofeng Chen, Patrick P. C. Lee, Wenjing Lou, A Hybrid Cloud Approach for Secure Authorized Deduplication, IEEE Transactions on Parallel and Distributed Systems, 2014.
- [2] P. Anderson and L. Zhang. Fast and secure laptop back-ups with encrypted de-duplication. In Proc. of USENIX LISA, 2010.
- [3] M. Bellare, S. Keelveedhi, and T. Ristenpart. Dup-less: Serveraided encryption for deduplicated storage. In USENIX Security Symposium, 2013.
- [4] M. Bellare, S. Keelveedhi, and T. Ristenpart. Message-locked encryption and secure deduplication. In EURO-CRYPT 2013
- [5] Li, X. Chen, M. Li, J. Li, P. Lee, and W. Lou. Secure deduplication with efficient and reliable convergent key management. In IEEE Transactions on Parallel and Distributed Systems, 2013.
- [6] Neal Leavitt, "Hybrid Clouds Move to the Forefront. Published by the IEEE Computer Society, MAY 2013.
- [7] [https://www.daniweb.com/images/attachments/0/WP\\_De\\_duplication\\_US\\_Letter\\_090702.pdf](https://www.daniweb.com/images/attachments/0/WP_De_duplication_US_Letter_090702.pdf).
- [8] R. Douceur, A. Adya, W. J. Bolosky, D. Simon, and M. Theimer. Reclaiming space from duplicate files in a serverless distributed