



## International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 9, September 2016

# A Survey on Robust Sentence Based Chunking Technique for Cloud Data Deduplication using Hybrid Cloud Architecture

Vane Shital<sup>1</sup>, Pawar Sonali<sup>2</sup>, Chinke Mayuri<sup>3</sup>, Nale Komal<sup>4</sup>, Prof. Dighe M. S.<sup>5</sup>

Student, Dept. of Computer Engineering, S.C.S College of Engineering, Rahuri Factory, Ahmednagar,  
University of Pune, India<sup>1</sup>

Student, Dept. of Computer Engineering, S.C.S College of Engineering, Rahuri Factory, Ahmednagar,  
University of Pune, India<sup>2</sup>

Student, Dept. of Computer Engineering, S.C.S College of Engineering, Rahuri Factory, Ahmednagar,  
University of Pune, India<sup>3</sup>

Student, Dept. of Computer Engineering, S.C.S College of Engineering, Rahuri Factory, Ahmednagar,  
University of Pune, India<sup>4</sup>

Asst. Professor, Dept. of Computer Engineering, S.C.S College of Engineering, Rahuri Factory, Ahmednagar,  
University of Pune, India<sup>5</sup>

**ABSTRACT:** As the cloud computing services are rapidly being used in recent days for storage and other purposes, cloud deduplication is also a such service which has to be focused on. As the cloud computing has gained immense focus in last few decades, cloud storage optimal management has become must. This paper surveys various works previously done in the field of cloud deduplication and therefore, but there is still a scope of improvement in deduplication of data over cloud storages. So the survey on the papers or researches emphasizes various deduplication techniques and the ways they are different from each other for efficient deduplication. Thus as there are many ways of deduplication in clouds, an efficient technique is to be found out having minimum drawbacks and maximum outcome..

**KEYWORDS:** deduplication, POW, hybrid storage on cloud, public clouds, credentials, confidentiality.

### I. INTRODUCTION

Current period is distributed computing time. Presently a days distributed computing has extensive variety of degree in information sharing. Distributed computing is give substantial measure of virtual environment concealing the stage and working frameworks of the client. Clients utilize the assets for sharing information. Be that as it may, clients need to pay according to the utilization of assets of cloud. Presently cloud administration suppliers are putting forth cloud administrations with ease furthermore with high dependability. Client can transfer the vast sum information on cloud and shared information to a huge number of clients. A cloud supplier is offer diverse administrations, for example, framework as an administration, stage as an administration, and so forth. Clients not have to buy the assets. As the information is get transferred by the client consistently it is basic assignment to deal with this regularly expanding information on the cloud. DE duplication is best technique to make well information administration in the distributed computing. This technique is turning out to be more fascination for information DE duplication. This system is send the information over the system required little measure of information. This technique has application in information administration and organizing. Information duplication is the procedure of decreasing the extent of information Also it is the best pressure system for the information DE duplication. This technique is send the information over the system required little measure of information. This system has application in information administration and organizing. Rather than keeping excess duplicates of the same information DE duplication just keep unique duplicate and give just

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 9, September 2016

references of the first duplicate to the repetitive information. There are two techniques for the duplication check, one is document level duplication check and other is piece content level duplication check. In the document level duplication check is expel the same name record from the capacity and square level DE duplication are evacuated the copy pieces. As the information DE duplication is considering the client information there must be need of the some security system. It emerges security and protection worry of the client's delicate information. In the conventional system client need to encode his own particular information without anyone else's input so there are distinctive figure documents for each new client.

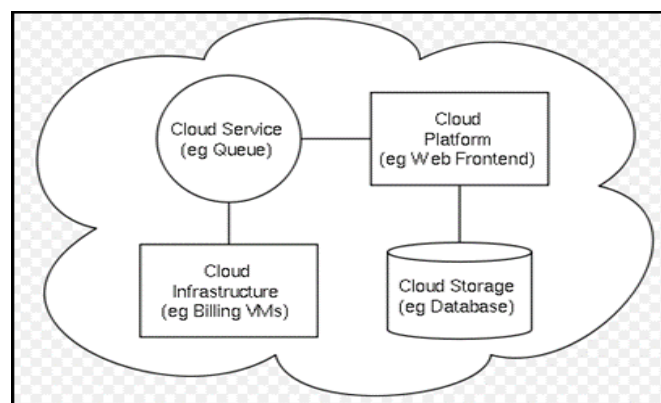


Fig 1. Cloud architecture and services

To maintain a strategic distance from the unapproved information DE duplication focalized information DE duplication is proposed to uphold the information privacy while checking the information duplication. The cloud giving numerous administrations as appeared in the above figure, for example, stage, administrations, base as an administration, and database as an administration. In this we are utilizing as a part of distributed storage as an administration. We are utilizing client accreditations to check the confirmation of the client. In the half and half cloud is available two sort of cloud such private cloud and open cloud. In private cloud store the client accreditation and client information present out in the open cloud. The half and half cloud take focal points of both open cloud and private cloud as appeared in the figure 2. open cloud and private cloud are available in the half and half cloud structural engineering When any client forward solicitation to people in general cloud to get to the information he have to present his data to the private cloud then private cloud will give a record token and client can get the entrance to the document lives on the general population cloud. We have utilized a half and half cloud construction modeling as a part of proposed. The document name is mind essential level in record information duplication and information DE duplication is checked at the square level. On the off chance that client needs to recover his information or download the information record he have to download both of the document from the cloud server this will prompts perform the operation on the same record this abuses the security of the distributed storage.

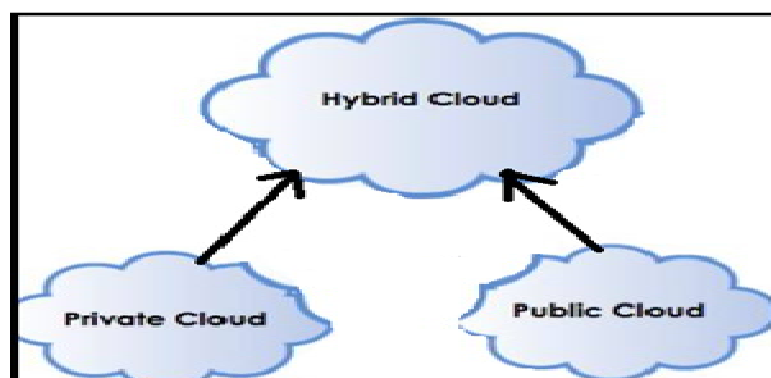


Fig 2. Hybrid Cloud Architecture.

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 9, September 2016

## II. LITERATURE SURVEY AND RELATED WORK

### A. A. DUPLESS: SERVER-AIDED ENCRYPTION FOR DEDUPLICATED STORAGE

Distributed storage administration suppliers, for example, Dropbox, Mozy, and others perform deduplication to spare space by just putting away one duplicate of every document transferred. Should customers routinely scramble their documents, be that as it may, funds are lost. Message-bolting encryption (the most unmistakable appearance of which is concurrent encryption) determines this strain. In any case it is intrinsically subject to savage power assaults that can recoup records falling into a known set. We propose a building design that acesides secure deduplicated stockpiling opposing savage power assaults, and acknowledge it in a framework called DupLESS. In DupLESS, customers encode under message-based keys acquired from a key-server by means of an absent PRF convention. It empowers customers to store scrambled information with a current administration, have the administration perform deduplication for their benefit, but then accomplishes solid privacy ensures. We demonstrate that encryption for deduplicated stockpiling can accomplish execution and space reserve funds near that of utilizing the stockpiling administration with plaintext information

### B. FAST AND SECURE LAPTOP BACKUPS WITH ENCRYPTED DE-DUPLICATION

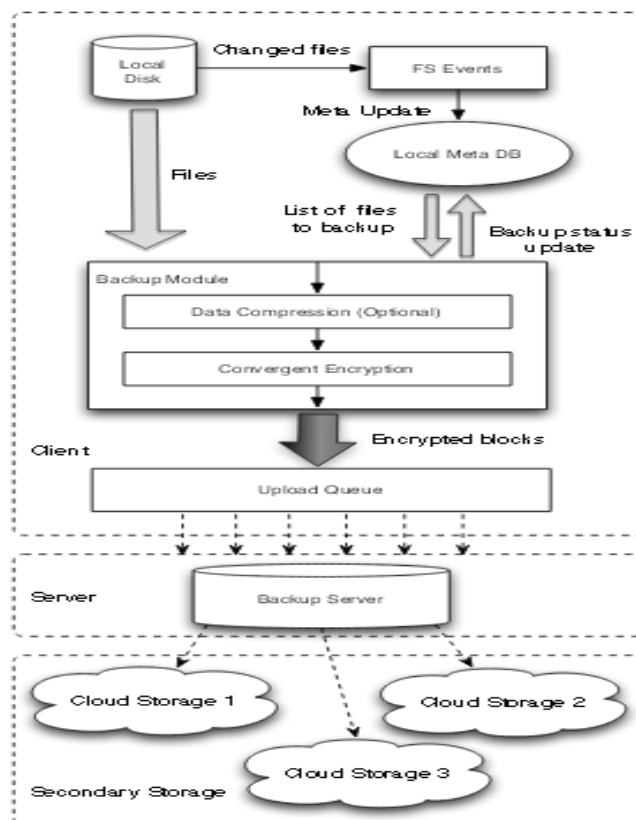


Fig 3. System diagram.

Numerous individuals now store extensive amounts of individual and corporate information on tablets or home PCs. These frequently have poor or discontinuous network, and are helpless against burglary or equipment disappointment. Ordinary reinforcement arrangements are not appropriate to this environment, and reinforcement administrations are every now and again deficient. This paper depicts a calculation which exploits the information which is basic between clients to build the pace of reinforcements, and diminish the capacity necessities. This calculation bolsters customer end per-client encryption which is essential for classified individual information. It likewise underpins a one of a kind

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 9, September 2016

element which permits prompt location of normal sub trees, dodging the need to question the reinforcement framework for each document. We portray a model usage of this calculation for Apple OS X, and present an investigation of the potential viability, utilizing genuine information acquired from an arrangement of ordinary clients. At last, we talk about the utilization of this model in conjunction with remote distributed storage, and present an investigation of the commonplace cost reserve funds

### C. SECURE DEDUPLICATION WITH EFFICIENT AND RELIABLE CONVERGENT KEY MANAGEMENT

Information deduplication is a system for taking out copy duplicates of information, and has been broadly utilized as a part of distributed storage to decrease storage room and transfer data transfer capacity. Promising as it may be, an emerging test is to perform secure deduplication in distributed storage. Albeit joined encryption has been widely received for secure deduplication, a basic issue of making focalized encryption down to earth is to productively and dependably deal with an immense number of united keys. This paper makes the first endeavor to formally address the issue of accomplishing effective and dependable key administration in secure deduplication. We first present a pattern approach in which every client holds an autonomous expert key for scrambling the focalized keys and outsourcing them to the cloud. On the other hand, such a standard key administration plan produces a tremendous number of keys with the expanding number of clients and obliges clients to dedicatedly secure the expert keys. To this end, we propose Dekey , another development in which clients don't have to deal with any keys all alone however rather safely circulate the united key shares over different servers. Security examination exhibits that Dekey is secure as far as the definitions determined in the proposed security model. As a proof of idea, we actualize Dekey utilizing the Ramp mystery sharing plan and show that Dekey brings about restricted overhead in reasonable situations.

### D. PROOFS OF OWNERSHIP IN REMOTE STORAGE SYSTEMS

Distributed storage frameworks are turning out to be progressively prominent. A promising innovation that holds their expense down is deduplication, which stores just a solitary duplicate of rehashing information. Customer side deduplication endeavors to recognize deduplication opportunities as of now at the customer and save the transmission capacity of transferring duplicates of existing documents to the server. In this work we recognize assaults that endeavor customer side deduplication, permitting an aggressor to access self-assertive size records of different clients in view of a little hash marks of these documents. All the more particularly, an aggressor who knows the hash mark of a record can persuade the capacity benefit that it possesses that document, henceforth the server lets the assailant download the whole record. (In parallel to our work, a subset of these assaults were as of late presented in the wild regarding the Dropbox record synchronization administration.) To overcome such assaults, we present the thought of verifications of-possession (PoWs), which lets a customer effectively demonstrate to a server that that the customer holds a document, as opposed to simply some short data about it. We formalize the idea of evidence of-proprietorship, under thorough security definitions, and thorough productivity prerequisites of Petabyte scale stockpiling frameworks. We then present arrangements in view of Merkle trees and particular encodings, and investigate their security. We actualized one variation of the plan. Our execution estimations show that the plan causes just a little overhead contrasted with guileless customer side deduplication.

### E. REVDEDUP: A REVERSE DEDUPLICATION STORAGE SYSTEM OPTIMIZED FOR READS TO LATEST BACKUPS

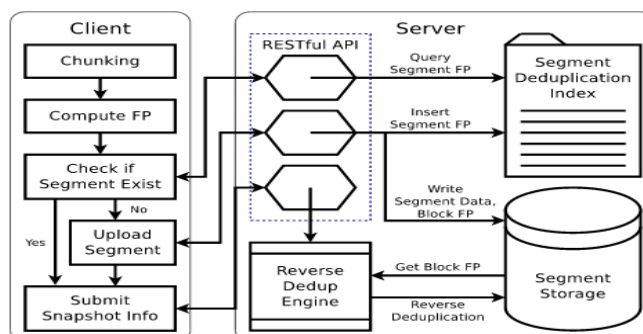


Fig 4. Reverse duplication

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 9, September 2016

Scaling up the reinforcement stockpiling for a perpetually expanding volume of virtual machine (VM) pictures is a basic issue in virtualization situations. While deduplication is known not dispose of copies for VM picture capacity, it additionally presents fracture that will corrupt read execution. We propose RevDedup, a deduplication framework that upgrades peruses to most recent VM picture reinforcements utilizing a thought called reverse deduplication. Conversely with traditional deduplication that expels copies from new information, RevDedup expels copies from old information, in this way moving discontinuity to old information while keeping the design of new information as consecutive as would be prudent. We assess our RevDedup model utilizing miniaturized scale benchmark and certifiable workloads. For a 12-week compass of certifiable VM pictures from 160 use rs, RevDedup accomplishes high deduplication productivitywith around 97% of sparing, and high reinforcement and read throughput on the request of 1GB/s. RevDedup additionally brings about little metadata overhead in reinforcement/read operations

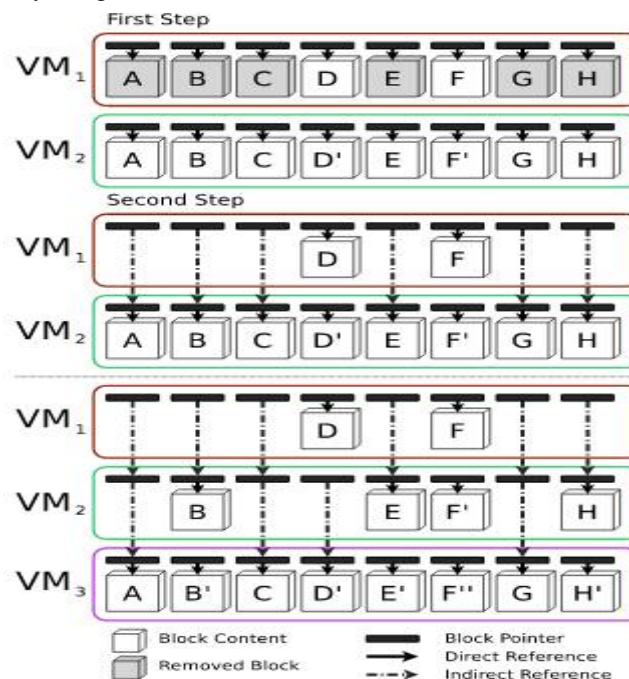


Figure 5. Reverse duplication example

## F. PRIVATE DATA DEDUPLICATION PROTOCOLS IN CLOUD STORAGE

In this paper, another idea which we call private information deduplication convention, a deduplication system for private information stockpiling is presented and formalized. Naturally, a private information deduplication convention permits a customer who holds a private information demonstrates to a server who holds a synopsis string of the information that he/she is the proprietor of that information without uncovering additional data to the server. Our idea can be seen as a supplement of the cutting edge open information deduplication conventions of Halevi et al. The security of private information deduplication conventions is formalized in the recreation based system in the connection of two-gathering calculations. A development of private deduplication conventions in view of the standard cryptographic suspicions is then introduced and examined. We demonstrate that the proposed private information deduplication convention is provably secure accepting that the basic hash capacity is crash flexible, the discrete logarithm is hard and the eradication coding calculation can deletion up to  $\alpha$ -division of the bits in the vicinity of malignant enemies in the vicinity of vindictive foes. To the best our insight this is the first deduplication convention for private information stockpiling [6].



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 9, September 2016

## III. ALGORITHMAM

### A. CHUNKING ALGORITHM

- Step 1: Upload a file
- Step 2: Read the file into buffer reader.
- Step 3: for all the sentences in the buffer reader
- Step 4: Read each line till Pull Stop is detected.
- Step 5: Consider the data till first full stop as a chunk
- Step 6: While all the sentences are read select data till every full stop as a chunk
- Step 7: Check if the currently read chunks are available in the cloud data base.
- Step 8: If yes the chunk is duplicate increase  
The count for duplicate chunk  
Else the chunk is original keep the chunk count as it is.

### B. AES ALGORITHM

1. Key Expansions  
For each round AES needs a different 128-bit block of round key also one more.
2. Initial Round  
AddRoundKey—with a block of the round key, each byte of the state is combined using bitwise xor.
3. Rounds
  - Sub Bytes—in this step each byte is replaced with another byte.
  - Shift Rows— for a certain number of steps, the state's last three rows are moved cyclically.
  - Mix Columns— on the columns of the state a mixing operation operates, in each column combining the four bytes.
  - AddRoundKey
4. Final Round (no Mix Columns)
  - Sub Bytes
  - Shift Rows
  - AddRoundKey

## IV. CONCLUSION

Here we can reason that our proposed framework information DE duplication of record is done approves way and safely. . In this we have additionally proposed new duplication check system which produce the token for the private document. The information client need to present the benefit alongside the united key as a proof of possession. We have settled more basic piece of the cloud information stockpiling which is just endured by diverse systems. Proposed routines guarantees the information duplication safely.

## REFERENCES

- [1].M. Bellare, S. Keelveedhi, and T. Ristenpart. Dupless: Serveraided encryption for deduplicated storage. In USENIX Security Symposium, 2013.
- [2].P. Anderson and L. Zhang. Fast and secure laptop backups with encrypted de-duplication. In Proc. of USENIX LISA, 2010.
- [3].J. 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.
- [4].S. Halevi, D. Harnik, B. Pinkas, and A. Shulman-Peleg. Proofs of ownership in remote storage systems. In Y. Chen, G. Danezis, and V. Shmatikov, editors, ACM Conference on Computer and Communications Security, pages 491–500. ACM, 2011.
- [5].C. Ng and P. Lee. Revdedup: A reverse deduplication storage system optimized for reads to latest backups. In Proc. of APSYS, Apr 2013.
- [6].W. K. Ng, Y. Wen, and H. Zhu. Private data deduplication protocols in cloud storage. In S. Ossowski and P. Lecca, editors, Proceedings of the 27th Annual ACM Symposium on Applied Computing, pages 441–446. ACM, 2012.