Scientific Journal of Impact Factor (SJIF): 4.14

e-ISSN (O): 2348-4470 p-ISSN (P): 2348-6406

# International Journal of Advance Engineering and Research Development

# Volume 3, Issue 11, November -2016

# Secure Distributed Deduplication With Verifiable Integrity Of Files

Chande Sayali V.<sup>1</sup>, Deore Sneha S.<sup>2</sup>, Habbli Komal P.<sup>3</sup>, Sawant Pratibha A.<sup>4</sup>, Asst. Prof. A.D. Misal<sup>5</sup>

<sup>1,2,3,4</sup>Department of Computer Engineering, P.K. Technical Campus, Chakan, Savitribai Phule, Pune. University <sup>5</sup>Assistant Professor, Department of Computer Engineering, P.K. Technical Campus, Chakan, Pune.

**Abstract** - Information may be a procedure for casting off copy duplicates of data, and has been generally utilized as a region of Cloud storage to decrease stowage and transfer transmission capability. On the opposite hand, there's one associate degree solely duplicate for each record place away in cloud despite the actual fact that such a document is possessed by an large variety of shoppers. Thus, framework enhances reposition use whereas decreasing unwavering quality. Besides, the check of security for delicate data to boot emerges after they are outsourced by shoppers to cloud. Aiming to address the higher than security challenges, this paper makes the primary endeavor to formalize the thought of spread dependable framework. We tend to propose new sent frameworks with higher unwavering quality during which the data items are spread over various cloud servers. the safety wants of data privacy and label consistency are to boot accomplished by presenting a settled mystery sharing arrange in spread reposition frameworks, instead of utilizing co-occurring coding as a region of past frameworks. Security examination shows that our frameworks are secure as way because the definitions determined within the projected security model. As a symptom of plan, we tend to execute the projected frameworks and exhibit that the caused overhead is extraordinarily restricted in wise things.

Keywords - Deduplication, Authorized duplicate check, Public auditing, shared data, Cloud computing.

## I. INTRODUCTION

Distributed storage may be a model of organized venture storage wherever info is place away in virtualized pools of capability that square measure by and enormous expedited by third gatherings. Distributed storage provides customers benefit, going from expense stinting and simplified comfort, to skillfulness opportunities and filmable administration. These extraordinary elements pull in additional shoppers to use and capability their own info to the distributed storage: as per the examination report, the degree of data in cloud is needed to accomplish forty trillion gigabytes in 2020. Despite the actual fact that distributed storage framework has been broadly speaking embraced, it neglects to oblige some very important rising wants, for instance, the capacities of examining honorableness of cloud files by cloud customers and distinctive derived files by cloud servers. We have a tendency to define each problem at a lower place. The first issue is honorableness examining. The cloud server has the capability diminish customers from the overwhelming weight of capability administration and support. The foremost distinction of distributed storage from standard in-house storage is that the knowledge is changed by means that of web associated place away in a subjective house, not in restraint of the purchasers by any stretch of the imagination that inevitably raises customer's extraordinary worries on the trait of their info.

Cloud storage provides customers with edges, starting from value saving and simplified convenience, to quality opportunities and scalable service. These nice options attract additional and additional customers to utilize and storage their personal knowledge to the cloud storage: consistent with the analysis report, the degree of information in cloud is anticipated to realize forty trillion gigabytes in 2020.

Even though cloud storage system has been wide adopted, it fails to accommodate some necessary rising wants like the skills of auditing integrity of cloud files by cloud shoppers and detective work duplicated files by cloud servers. We have a tendency to illustrate each issue below. The primary drawback is integrity auditing. The cloud server is in a position to alleviate shoppers from the serious burden of storage management and maintenance.

Cloud storage provides customers with edges, starting from value saving and simplified convenience, to quality opportunities and scalable service. These nice options attract additional and additional customers to utilize and storage their personal knowledge to the cloud storage: consistent with the analysis report, the degree of information in cloud is anticipated to realize forty trillion gigabytes in 2020.

Even though cloud storage system has been wide adopted, it fails to accommodate some necessary rising wants like the skills of auditing integrity of cloud files by cloud shoppers and detective work duplicated files by cloud servers. We have a tendency to illustrate each issue below. The primary drawback is integrity auditing. The cloud server is in a position to alleviate shoppers from the serious burden of storage management and maintenance.

## International Journal of Advance Engineering and Research Development (IJAERD) Volume 3, Issue 11, November -2016, e-ISSN: 2348 - 4470, print-ISSN: 2348-6406

## II. RELATED WORK

In this paper[4], we have a tendency to gift a mechanism to reclaim house from this incidental duplication to create it obtainable for controlled file replication. Our mechanism includes convergent coding, that permits duplicate files to amalgamated into the house of one file, notwithstanding the files are encrypted with totally different users' keys, and 2) dish, a Self-Arranging, Lossy, Associative info for aggregating file content and site info in a very localized, scalable, fault-tolerant manner. Large-scale simulation experiments show that the duplicate-file coalescing system is scalable, extremely effective, and fault-tolerant. This paper addresses the issues of characteristic and coalescing identical files within the Farsite distributed classification system, for the aim of reclaiming space for storing consumed by incidentally redundant content.

[In this paper[5],the problem of providing secure outsourced storage that each supports deduplication and resists bruteforce attacks. We have a tendency to style a system, DupLESS, that mixes a CE-type baseMLE theme with the flexibility to get message-derived keys with the assistance of a key server (KS) shared amongst a gaggle of purchasers. The purchasers move with the Kansas by a protocol for oblivious PRFs, ensuring that the Kansas will cryptographically combine on the Q.T. material to the per message keys whereas learning nothing about files hold on by purchasers. These mechanisms make sure that DupLESS provides strong security against external attacks that compromise the SS and communication channels (nothing is leaked on the far side file lengths, equality, and access patterns), which the protection of DupLESS graciously degrades within the face of comprised systems.

In this paper[6],Definitions each for privacy and for a kind of integrity that we have a tendency to decision tag consistency. Supported this foundation, we have a tendency to build each sensible and theoretical contributions. On the sensible aspect, we offer memory board security analyses of a natural family of MLE schemes that has deployed schemes. On the theoretical aspect the challenge is commonplace model solutions, and that we build connections with settled secret writing, hash functions secure on correlative inputs and also the sample-then-extract paradigm to deliver schemes underneath totally different assumptions. And for various categories of message sources. Our work shows that MLE could be a primitive of each sensible and theoretical interest.

In this paper[10], is that we will eliminate duplicate copies of storage information and limit the injury of purloined information if we have a tendency to decrease the worth of that purloined info to the assaulter. This paper makes the primary conceive to formally address the matter of achieving economical and reliable key management in secure deduplication. we have a tendency to initial introduce a baseline approach during which every user holds associate freelance key for encrypting the confluent keys and outsourcing them. However, such a baseline key management theme generates a vast variety of keys with the increasing variety of users and needs users to dedicatedly shield the master keys. Confusing the assaulter with phony info.

In this paper[2], describes version one.2 of jerasure, a library in C/C++ that supports erasure writing in storage applications. During this paper, we have a tendency to describe each the techniques and algorithms, and the interface to the code. Thus, this is a quasi-tutorial and a programmer's guide. Version 1.2 adds Blaum-Roth and Liber8tion writing to the library, provides higher examples, associate degree an example file encoder/decoder. To boot, it removes a bug from the previous write up: the pocketsize should be a multiple of sizeof (long). It ought not to be a multiple of w.

## III. SYSTEM ARCHITECTURE

## A. PROCESS BLOCK DIAGRAM

We propose Dekey, another development within which shoppers do not have to traumatize any keys on their own however rather safely applicable the simultaneous key shares over varied servers. Dekey utilizing the Ramp mystery sharing arrange and exhibit that Dekey brings regarding affected overhead in wise things we tend to propose another development referred to as Dekey, which supplies productivity and responsibleness insurances to unified key administration on each consumer and cloud warehousing sides. Another development Dekey is planned to present skilled and solid unified key administration through united key Deduplication and mystery sharing. Dekey underpins each record levels Deduplication. Security investigation exhibits that Dekey is secure as way because the definitions determined within the planned security model. Specifically, Dekey stays secure even the foe controls a collection range of key servers. We tend to execute Dekey utilizing the mystery sharing arranges that empowers the key administration to regulate to various responsibleness and classification levels. Our assessment shows that Dekey brings regarding affected overhead in typical transfer/download operations in wise cloud things. We tend to conjointly propose a 3rd party auditor for verification of files store on cloud on demand of cloud knowledge owner or user for the asking.

International Journal of Advance Engineering and Research Development (IJAERD) Volume 3, Issue 11, November -2016, e-ISSN: 2348 - 4470, print-ISSN: 2348-6406



Figure: Architecture Diagram

#### **B. ALGORITHM**

#### **AES Algorithm**

Steps:

- 1. Derive the set of round keys from the cipher key.
- 2. Initialize the state array with the block data (plaintext).
- 3. Add the initial round key to the starting state array.
- 4. Perform nine rounds of state manipulation.
- 5. Perform the tenth and final round of state manipulation.
- 6. Copy the final state array out as the encrypted data (cipher text).

### **MD5** Algorithm

Steps:

- **1. Append Padding Bits**
- 2. Append Length
- 3. Initialize MD Buffer
- 4. Process Message in 16-Word Blocks
- 5. Output

#### • Advantages of Proposed System:

The detection of masquerade activity. The confusion of the attacker and the additional costs incurred to distinguish real from bogus information, and the deterrence effect which, although hard to measure, plays a significant role in preventing masquerade activity by risk-averse attackers.

### C. MATHEMATICAL MODEL

System Description: Let S be the Whole system which consists, S=fI,P,Og Where, I-Input, P- procedure, O- Output. I-F,U F-Filesset of fF1,F2,..,FNg U- No of UsersfU1,U2,,UNg Procedure(P): P=fPOW, n ,POWB, POWF, t, i,j ,m , kg. Where,

@IJAERD-2016, All rights Reserved

POW - proof of ownership.
n - No of servers.
POWB-proof of ownership in blocks.
POWF proof of ownership in files.
t - tag.
i- Fragmentation.
j- No of server.
m-message
k- Key.
DD = It is deterministic.
NDD = If project contains large database, it is hard to determine.
NP-Complete = This project is NP Complete.
File Upload(FU):

Step 1: File level deduplication

If a file duplicate is found, the user will run the PoW protocol POWF with each S-CSP to prove the file ownership.for the j-th server with identity idj, the user first computes

 $\phi$ F;idj=TagGen'(F, idj)

and runs the PoW proof algorithm with respect to  $\phi F$ , idj. If the proof is passed, the user will be provided a pointer for the piece of file stored at j-th S-CSP. Otherwise, if no duplicate is found, the user will proceed as follows: First divides F into a set of f ragments{Bi} (where i = 1, 2, ...). For each fragment Bi, the user will perform a block-level duplicate check.

Step 2: Block Level deduplication If there is a duplicate in S-CSP, the user runs PoWBon input:

With the server to prove that he owns the block Bi. If it is passed, the server simply returns a block pointer of Bi to the user. The user then keeps the block pointer of Bi and does not need to upload Bi.

Proof of ownership (POW):

Step 1: compute and send  $\phi'$  to the verifier.

Step 2: present proof to the storage server that he owns F in an interactive way with respect to  $\phi'$  The PoW is successful if the proof is correct

 $\phi' = \phi(F)$ 

File Download(FD)-

To download a file F, the user first downloads the secret shares  $\{cij,mfj\}$  of the file from kout of n storage servers. Specifically, the user sends all the pointers for F to k out of n servers. After gathering all the shares, the user reconstructs file F, macFby using the algorithm of Recover( $\{\cdot\}$ ). Then, he verifies the correctness of these tags to check the integrity of the file stored in S-CSPs.

## **IV. CONCLUSION**

Security investigation exhibits that Dekey is secure as so much because the definitions determined within the planned security model. Specifically, Dekey stays secure even the foe controls a collection range of key servers. We have a tendency to execute Dekey utilizing the mystery sharing set up that empowers the key administration to regulate to numerous responsibleness and classification levels. Our assessment shows that Dekey brings concerning unnatural overhead in typical transfer/download operations in wise cloud things. We have a tendency to targeting the difficulty of evaluating if Associate in nursing untrusted server stores a customer's information. We have a tendency to conferred a model for demonstrable information possession (PDP), within which it's tempting to reduce the lupus erythematous piece gets to, the calculation on the server, and also the client-server correspondence. Our answers for PDP t this model: They cause an occasional (or even steady) overhead at the server and oblige a bit, consistent live of correlation.

# @IJAERD-2016, All rights Reserved

## V. REFERENCES

- [1]. Amazon, Case Studies, https://aws.amazon.com/solutions/casestudies/hash backup.
- [2].J. S. Plank, S. Simmerman, and C. D. Schuman, Jerasure: A library in C/C++ facilitating erasure coding for storage applications- Version 1.2, University of Tennessee, Tech. Rep. CS-08-627, August 2008.
- [3]. M. O. Rabin, Fingerprinting by random polynomials, Center for Research in Computing Technology, Harvard University, Tech. Rep. Tech. Report TR-CSE-03-01, 1981.
- [4]. J. R. Douceur, A. Adya, W. J. Bolosky, D. Simon, and M. Theimer, Reclaiming space from duplicate less in a server less distributed file system. In ICDCS, 2002, pp. 617624.
- [5].M. Bellare, S. Keelveedhi, and T. Ristenpart, Dupless: Serve raided encryption for deduplicated storage, in USENIX Security Symposium, 2013.
- [6]. Message-locked encryption and secure de-duplication, in EUROCRYPT, 2013, pp. 296312.
- [7]. G. R. Blakley and C. Meadows, Security of ramp schemes, in Advances in Cryptology: Proceedings of CRYPTO 84, ser. Lecture Notes in Computer Science, G. R. Blakley and D. Chaum, Eds. Springer-Verlag Berlin/Heidelberg, 1985, vol. 196, pp. 242268. 8. A. D. Santis and B. Masucci, Multiple ramp schemes, IEEE Transactions on Information Theory, vol. 45, no. 5, pp. 17201728, Jul. 1999.
- [8]. M. O. Rabin, Efficient dispersal of information for security, load balancing, and fault tolerance, Journal of the ACM, vol. 36, no. 2, pp. 335348, Apr. 1989.
- [9]. A. Shamir, How to share a secret, Commun. ACM, vol. 22, no. 11, pp. 612613, 1979.
- [10].N.O.AGRAWAL, Prof Mr. S.S.KULKARNI ,"Secure Deduplication And Data Security With Efficient And Reliable CEKM."