

Novel Approach for Secure Data deduplication System

¹ARADHYULA VENKATA RAMU, ²D.RAVI KIRAN

¹M.Tech (CSE), Priyadrshini Institute of Technology & Management

²Professor (Dept.of CSE), Priyadrshini Institute of Technology & Management

Abstract:- Data deduplication is a strategy for evaluating duplicate copies of information and has been broadly utilized as a part of cloud storage to abatement storage space and transfer data transmission. Then again, there is stand out duplicate for every record put away in cloud regardless of the fact that such a document is claimed by countless. In like manner, deduplication framework progress storage use while diminishing unwavering quality. Moreover, the challenge of protection for delicate information likewise happens when they are outsourced by clients to cloud. Wanting to address the above security test, this paper builds the first push to commend the thought of scattered dependable deduplication framework. This paper prescribes another disseminated deduplication frameworks with upper reliability in which the information pieces are circulated from corner to cornering various cloud servers. The wellbeing needs of information security and label strength are likewise performed by presenting a deterministic mystery sharing plan in appropriated storage frameworks, rather than utilizing convergent encryption as a part of past deduplication frameworks.

Keywords – Cloud Storage, Convergent Key, Deduplication, secret sharing, reliability



1. INTRODUCTION

By the unpredictable improvement of computerized information, deduplication methods are extensively connected with to reinforcement information and diminish system and storage straightforwardness by notification and destroy repetition among information. As an option of keeping up numerous information copies with the same substance, deduplication diminishing repetitive information by keeping up just single duplicate and alluding other excess information to that duplicate. Deduplication has internal much fixation from both scholastic world and industry since it can truly recoup storage usage and keep storage space, especially for the applications with high deduplication proportion, for example, archival storage frameworks. Various deduplication frameworks have been anticipated in light of different deduplication plan, for example, customer side or server-side deduplication, document level or piece level deduplications. Specially, with the coming of

cloud storage, information deduplication strategy develop to be more stunning and fundamental for the administration of regularly expanding amount of information in cloud storage administrations which rouses Try and club to outsource information storage to outsider cloud suppliers, In the event that we consider a percentage of the illustrations as confirmations:

[i] Today's cloud storage administrations, for example, Google Drive, Drop box have been relating deduplication to spare the system data transmission and the storage cost with customer side deduplication.

Two sorts of deduplication as far as the size :(a) piece level deduplication, which figure out and take out redundancies among information blocks.(b)file-level deduplication, which decide redundancies between diverse documents and destroy these redundancies to abatement capacity requests, and The

record can be isolated into lesser settled size. Utilizing altered size pieces abbreviate the figuring of square bound-emerge, even as utilizing variable-size squares .[ii]Despite the way that deduplication strategy can gather the storage space for the cloud storage administration suppliers, it diminishes the consistency of the framework. Information consistency is truly an exceptionally crucial issue in a deduplication storage framework on the grounds that there is one and only duplicate for every record amasses in the server pooled by every one of the proprietors. On the off chance that such a pooled record was lost; an unnecessarily expansive measure of information gets to be inaccessible due to the inaccessibility of the considerable number of documents that share this document. In the event that the estimation of a document were figured regarding the measure of record information that would be lost if there should arise an occurrence of behind a solitary lump, then the amount of client information lost when a document in the storage framework is ruined develops with the quantity of the solidarity of the piece. In this manner, how to affirmation of high information consistency in deduplication framework is an essential issue. The vast majority of the first deduplication plan has just been measured in a solitary server area. then again, as bunches of deduplication frameworks and cloud storage frameworks are arranged by clients and capacity for higher constancy, especially in archival storage frameworks where information are key and ought to be pruned over long time point. This includes the deduplication storage frameworks give unwavering quality practically identical to other high-accessible frameworks.

Besides, the test for information protection likewise emerges as more delicate information is being outsourced by clients to cloud. Encryption instruments have for the most part been used to secure the privacy before outsourcing information into cloud. Most business storage administration supplier is hesitant to apply encryption over the information on the grounds that it makes deduplication outlandish. The reason is that the conventional encryption components, including open key encryption and symmetric key encryption, require diverse clients to scramble their information with their own particular keys. Therefore, indistinguishable information copies of diverse clients will prompt distinctive ciphertext. To take care of the issues of classification and deduplication, the thought

of convergent encryption has been star postured and generally embraced to uphold information secrecy while acknowledging deduplication. In any case, these frameworks accomplished privacy of outsourced information at the expense of diminished mistake versatility. In this manner, how to ensure both secrecy and unwavering quality while accomplishing deduplication in a cloud storage framework is still a test.

2. OUR CONTRIBUTIONS:

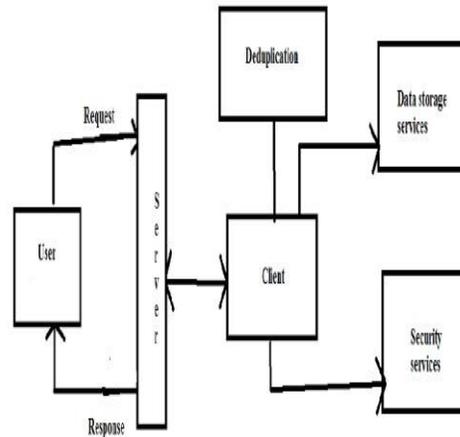
In this paper, we demonstrate to plan secure deduplication frameworks with higher unwavering quality in cloud processing. We present the circulated cloud storage servers into deduplication frameworks to give better adaptation to internal failure. To advance secure information classification, the secret sharing procedure is used, which is likewise perfect with the appropriated storage frameworks. In more subtle elements, a document is first part and encoded into pieces by utilizing the strategy of secret sharing, rather than encryption systems. These shares will be appropriated over numerous free storage servers. Moreover, to bolster deduplication, a short cryptographic hash estimation of the substance will likewise be registered and sent to every storage server as the unique mark of the piece put away at every server. Just the information proprietor who first transfers the information is required to process and convey such secret shares while every single after client who possess the same information duplicated don't have to register and store these shares anymore. To recoup information copies, clients must get to a base number of storage servers through validation and acquire the secret shares to recreate the information. At the end of the day, the secret shares of information may be open by the approved clients who claim the relating information duplicate.

Another recognizing highlight of our proposition is that information trustworthiness, including label consistency, can be accomplished. The conventional deduplication techniques can't be straightforwardly broadened and connected in dispersed and multi-server frameworks. To clarify further, if the same short esteem is put away at an alternate cloud storage server to bolster a duplicate check by utilizing a conventional deduplication technique, it can't avoid the plot assault propelled by numerous servers. At the end of the day,

any of the servers can get shares of the information put away at alternate servers with the same short esteem as verification of proprietorship. Besides, the label consistency, which was initially formalized by [5] to keep the duplicate/ciphertext substitution assault, is considered in our convention. In more subtle elements, it keeps a client from transferring a malevolently produced ciphertext such that its tag is the same with another genuinely created ciphertext. To accomplish this, a deterministic secret sharing system has been formalized and used. As far as anyone is concerned, no current work on secure deduplication can legitimately address the dependability and label consistency issue in disseminated storage frameworks. This paper makes the accompanying commitments.

Four new secure deduplication frameworks are ace postured to furnish proficient deduplication with high unwavering quality for document level and piece level deduplication, individually. The secret part strategy, instead of conventional encryption strategies, is used to ensure information privacy. In particular, information is part into sections by utilizing secure secret sharing schemes and put away at distinctive servers. Our proposed developments support both record level and square level deduplication.

Security examination exhibits that the proposed deduplication frameworks are secure as far as the definitions indicated in the proposed security model. In more points of interest, privacy, unwavering quality and honesty can be accomplished in our proposed framework. Two sorts of agreement assaults are considered in our answers. These are the plot assault on the information and the intrigue assault against servers. Specifically, the information stays secure regardless of the possibility that the enemy controls a set number of storage servers. We execute our deduplication frameworks utilizing the Incline secret sharing scheme that empowers high re-obligation and classification levels. Our assessment results exhibit that the new proposed developments are proficient and the redundancies are enhanced and practically identical with the other storage framework supporting the same level of unwavering quality.



System architecture

3. THE DISTRIBUTED DEDUPLICATION SYSTEMS:

The distributed deduplication systems future aim is to reliably store data in the cloud while achieving privacy and consistency. Its main objective is to allow deduplication and distributed storage of the data diagonally multiple storage servers. As an alternative encrypting the data to keep the privacy of the data, new structures put on the top-secret intense technique to split data into shards. These shards will then be distributed transversely in multiple storage servers.

3.1 The File-level Distributed Deduplication System

To maintain efficient duplicate check, tags for each file will be calculated and are directed to S-CSPs. To avoid a conspiracy attack hurled by the S-CSPs, the tags deposited at different storage servers are computationally autonomous and different. The details of the structure as follows.

System setup. In our structure, the number of Storage servers S-CSPs is expected to be i with identities denoted by id_1, id_2, \dots, id_n , correspondingly. Describe the security parameter as 1 and set a secret sharing scheme $SS = (Share, Recover)$, and a tag generation algorithm $TagGen$. The file storage system for the storage server is set to be $\#File Upload$. To upload a file F , the user relates with S-CSPs to achieve the deduplication. More exactly, the user firstly calculates and sends the file tag $\phi_F = TagGen(F)$ to S-CSPs for the file duplicate check. When a duplicate is found, the user calculates and sends $\phi_F; id_j = TagGen'(F, id_j)$ to the j -

th server with identity id_j via the secure channel for $1 \leq j \leq n$. The motive for presenting an index j is to avoid the server from receiving the shares of other S-CSPs for the same file or block, which will be described in detail in the security analysis. If $X_{F;id_j}$ equals the metadata stored with X_F , the user will be provided a pointer for the shard stored at server id_j . Else, if no duplicate is found, the user will continue as follows. He runs the secret sharing algorithm SS over F to get $\{c_j\} = \text{Share}(F)$, where c_j is the j -th shard of F . He also calculates $X_{F;id_j} = \text{TagGen}'(F, id_j)$, which helps as the tag for the j -th S-CSP. As a final point, the user uploads the set of values $\{\varphi_F, c_j, X_{F;id_j}\}$ to the S-CSP with identity id_j via a secure channel. The S-CSP stores these values and returns a pointer back to the user for local storage.

File Download. To download a file F , the user first downloads the secret shares $\{c_j\}$ of the file from k out of n storage servers. Exactly, the user sends the pointer of F to k out of n S-CSPs. After meeting enough shares, the user reconstructs file F by using the algorithm of $\text{Recover}(\{c_j\})$. This method provides fault tolerance and lets the user to remain available even if any limited subsets of storage servers fail.

3.2. The Block-level Distributed Deduplication System

We demonstrate how to attain the fine-grained block-level distributed deduplication. In a block-level deduplication system, the user also needs to firstly achieve the file-level deduplication before uploading his file. If no duplicate is found, the user splits this file into blocks and does block-level deduplication. The system arrangement is the same as the file-level deduplication system; excluding the block size parameter will be defined in addition. Following, the details of the algorithms of File Upload and File Download are mentioned.

File Upload. To upload a file F , the user first achieves the file-level deduplication by sending φ_F to the storage servers. If a duplicate is found, the user will achieve the file-level deduplication, else, if no duplicate is found, the user achieves the block-level deduplication as follows.

Initially divides F into a set of fragments $\{A_i\}$ (where $i = 1, 2, \dots$). For each fragment A_i , the user will achieve a block-level duplicate check by computing $X_{B_i} = \text{TagGen}(A_i)$, where the data handling and duplicate check of block-level deduplication is the same as that

of file-level deduplication if the file F is substituted with block B_i . Upon getting block tags $\{X_{B_i}\}$, the server with identity id_j computes a block signal vector R_{B_i} for each i . i) If $R_{B_i} = 1$, the user additionally computes and sends $X_{B_i;j} = \text{TagGen}'(B_i, j)$ to the S-CSP with identity id_j . If it also equals the matching tag stored, S-CSP sends a block pointer of B_i to the user. At that time, the user keeps the block pointer of B_i and does not need to upload B_i .

ii) If $R_{B_i} = 0$, the user runs the secret sharing algorithm SS over B_i and gets $\{c_{ij}\} = \text{Share}(B_i)$, where c_{ij} is the j -th secret share of B_i . The user also computes $X_{B_i;j}$ for $1 \leq j \leq n$ and uploads the set of values $\{X_F, X_{F;id_j}, c_{ij}, X_{B_i;j}\}$ to the server id_j through a secure channel. The S-CSP returns the consistent pointers back to the user.

File Download. To download a file $F = \{A_i\}$, the user first downloads the secret shares $\{c_{ij}\}$ of all the blocks A_i in F from k out of n S-CSPs. Exactly, the user sends all the pointers for A_i to k out of n servers. Subsequently gathering all the shares, the user recreates all the fragments A_i using the algorithm of $\text{Recover}(\{c_{ij}\})$ and gets the file $F = \{A_i\}$.

4. ALGORITHMS USED:

Here we discuss about Secret Sharing Scheme. Let us have a look on two algorithms in a secret sharing scheme, which are Share and Recover. The secret is separated and shared by using Share. With enough shares, the secret can be pull out and improved with the algorithm of Recover. Here, the Ramp secret sharing scheme (RSSS) [7], [8] is assumed to secretly split a secret into shards. Definitely, the (i, j, p) -RSSS (where $n_i > j > p \geq 0$) produces n shares from a secret so that (i) the secret can be improved from any j or more shares, and (ii) No evidence about the secret can be assumed from any p or less shares. Two algorithms, Share and Recover, are defined in the (L, j, p) -RSSS. Share splits a secret S into $(j - p)$ pieces of equal size, generates p random pieces of the same size, and translates the j pieces using a non-systematic j of- i removal code into i shares of the same size; Improve takes any j out of i shares as inputs and then outputs the original secret S . We can say that when $p = 0$, the $(i, j, 0)$ -RSSS turn into the (i, j) Rabin's Information Dispersal Algorithm (IDA) [9]. When $p = j - 1$, the $(L, j, j - 1)$ -RSSS becomes the (i, j) Shamir's Secret Sharing Scheme (SSSS) [10].

Tag Generation Algorithm. In our structures

below, two kinds of tag generation algorithms are defined, that is, TagGen and TagGen'. TagGen is the tag generation algorithm that records the original data copy C and outputs a tag T (C). This tag will be produced by the user and practical to achieve the duplicate check with the server. Alternative tag generation algorithm TagGen' precedes as input a file C and an index j and outputs a tag. This tag, generated by users, is used for the proof of ownership for C.

Message authentication code. A message authentication code (MAC) is a tiny piece of data used to authenticate a message and to make available integrity and validity assurances on the message. Here the message verification code is applied to attain the reliability of the contract out stored files. It can be simply made with a keyed i.e cryptographic hash function, which takes input as a secret key and an arbitrary-length file that supplies to be authenticated, and outputs a MAC. Individual users with the same key making the MAC can confirm the exactness of the MAC value and notice whether the file has been changed or not.

4.1. Advantages of Proposed work:

- Unique feature of the proposal is that data integrity, as well as tag consistency, can be achieved.
- For our knowledge, no current work on safe deduplication can appropriately address the reliability and tag consistency problem in distributed storage systems.
- The proposed constructions maintain both file-level and block-level deduplication.

Security analysis determines that the proposed deduplication systems are safe in terms of the definitions stated in the proposed security model. If we want to elaborate we can also say that confidentiality, reliability and integrity can be achieved in the proposed system. Two kinds of collusion attacks are measured in our solutions. These are the collusion attack on the data and the collusion attack against servers. In specific, the data remains secure even if the opponent controls a limited number of storage servers. The implementation of deduplication systems using the Ramp secret sharing scheme allows high reliability and confidentiality levels. The evaluation results prove that the proposed constructions are efficient and the redundancies are optimized and similar with the other storage system supporting the same level of dependability.

5. CONCLUSION:

The proposed data deduplication frameworks are to expand the consistency of information however achieving the protection of the client's outsourced information without an encryption apparatus. The security of label consistency and respectability were accomplished. The execution of deduplication frameworks utilizing the Incline secret sharing scheme here gives the exhibit that it obtains little encoding/interpreting overhead contrasted with the system transmission overhead in standard download/transfer operations.

REFERENCES

- [1] J. Gantz and D. Reinsel, "The digital universe in 2020: Bigdigital shadows and biggest growth in the fareast," <http://www.emc.com/collateral/analyst-reports/idcthe-digital-universe-in-2020.pdf>, Dec 2012.
- [2] M. O. Rabin, "Fingerprinting by random polynomials," Center for Research in Computing Technology, Harvard University, Tech.Rep. Tech. Report TR-CSE-03-01, 1981.
- [3] J. R. Douceur, A. Adya, W. J. Bolosky, D. Simon, and M. Theimer, "Reclaiming space from duplicate files in a serverless distributed file system." in ICDCS, 2002, pp. 617-624.
- [4] M. Bellare, S. Keelveedhi, and T. Ristenpart, "Dupless: Server-aided encryption for deduplicated storage," in USENIX Security Symposium, 2013.
- [5] "Message-locked encryption and secure deduplication," in EUROCRYPT, 2013, pp. 296-312.
- [6] 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. 242-268.
- [7] A. D. Santis and B. Masucci, "Multiple ramp schemes," *IEEE Transactions on Information Theory*, vol. 45, no. 5, pp. 1720-1728, 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. 335-348, Apr. 1989.
- [9] A. Shamir, "How to share a secret," *Commun. ACM*, vol. 22, no. 11, pp. 612-613, 1979.
- [10] 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, 2014, pp. vol.25(6), pp. 1615–1625.

[11] S. Halevi, D. Harnik, B. Pinkas, and A. Shulman-Peleg, “Proofs of ownership in remote storage systems.” in ACM Conference on Computer and Communications Security, Y. Chen, G. Danezis, and V. Shmatikov, Eds. ACM, 2011, pp. 491–500.

[12] 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. [13] J. S. Plank and L. Xu, “Optimizing Cauchy Reed-solomon Codes for fault-tolerant network storage applications,” in NCA-06: 5th IEEE International Symposium on Network Computing Applications, Cambridge, MA, July 2006.

[14] C. Liu, Y. Gu, L. Sun, B. Yan, and D. Wang, “Radmad: High reliability provision for large-scale deduplication archival storage systems,” in Proceedings of the 23rd international conference on Supercomputing, pp. 370–379.