# A Survey on Reduce the Storage Space in Multi Cloud using DeDuplication and Improving Security

[1] Shweta Joshi,[2] Prof.Shradhdha Bhalodiya,[3] Dr.Vipul Vekariya

[1]PG Scholar,Computer Engineering,Noble Group of Institution,Junagadh,India

[2]Assistant Professor, Computer Engineering,Noble Group of Institution,Junagadh,India

[3]Principal,Noble Group of Institution,Junagadh,India

*Abstract:*   With the increasing use of cloud storage platform for storage and process, there's a growing demand of some mechanism or methodology which can offer the ability of eliminating redundant knowledge and thereby achieving higher area and information measure necessities of storage services. In alternative words, it ought to offer a capability of a much better potency once applied across multiple users. During this context, Cloud Service suppliers usually use Deduplication, that stores solely one copy of every file or block by eliminating redundant knowledge. However providing a secure Deduplication is AN uphill task. During this regard, a trial is formed to gift a survey on the various aspects of Deduplication.

*IndexTerms* – **Deduplication,security**

## I. INTRODUCTION

With the increasing quality and cost-effectiveness of cloud storage systems, several firms and organizations have migrated or attempt to migrate knowledge from their non-public knowledge centers to the cloud. However, entirely reckoning on a specific cloud storage supplier incorporates a variety of probably serious issues. First, it will cause the supposed trafficker lock-in downside for the shoppers [1,2], which ends in prohibitively high price for purchasers to change from one supplier to a different. Second, it will cause service disruptions that successively can result in SLA violation, because of cloud outages, leading to penalties, financial or alternative forms, for the service suppliers. Examples embrace a series of high-profile cloud outages within the year of 2013 for cloud suppliers, like Amazon, Microsoft and Google [3], from a 5-min failure that costed [*fr1] 1,000,000 greenbacks to a week-long disruption that costed AN immeasurable quantity of brand name harm.

From Jan to March 2014, DropBox has full-fledged doubly service outages [3]. additional seriously, Nirvanix filed for Chapter eleven bankruptcy protection on Gregorian calendar month one, 2013 [4]. the corporate gave customers 2 weeks' notice to retrieve their knowledge. Some users had petabytes of knowledge with single copy hold on in Nirvanix. Third, entirely reckoning on a specific cloud storage supplier also can end in attainable augmented service prices and knowledge security problems, like the information escape downside [5]. Therefore, victimization multiple freelance cloud suppliers, referred to as Cloud-of-Clouds, is an efficient thanks to offer higher availableness for the cloud storage systems. in a very Cloud-of-Clouds storage system, the information redundancy is introduced to judiciously distribute the information among the clouds. Thus, the redundant knowledge distribution theme is critically necessary for the storage availableness, performance, price and area potency. many systems are planned for Cloud-of-Clouds. RACS [1] uses the erasure cryptography to mitigate the seller lock-in downside encountered by a user once shift the cloud vendors. It transparently stripes the information across multiple cloud storage suppliers with RAID-like techniques. HAIL [6] provides integrity and availableness guarantees for the hold on knowledge. It permits a group of servers persuade a shopper that a hold on file is unbroken and recoverable by the approaches adopted from the cryptological and distributed-systems communities. NCCloud [7] achieves the costeffective repair for a permanent single-cloud supplier failure to enhance the supply of cloud storage services. it's designed supported network-coding-based storage schemes referred to as create codes with a stress on the storage repair, excluding the failing cloud in repair.

In a base paper they need planned a system for deduplication checking however however they need instructed that System can became additional powerful if it mix for block level further as file level deduplication check. this can overcome the disadvantage that if whole file is same then additionally here it check the duplication chunk by chunk. it's avoided if we tend to initial perform the file level deduplication and if file is exclusive then solely opt for the block level deduplication check. additionally checked {for another|for an additional|for one additional} cryptography algorithmic program to supply the more security

## II. DATA DEDUPLICATION

Data deduplication refers to methodologies that store solely one copy of redundant knowledge and thereby offer one copy. By eliminating redundant knowledge each space and network bandwidth[2].With relation to service suppliers, it offers secondary price savings in power and cooling that is achieved by reducing the amount of spindles[3]. It ensures that just one copy of knowledge is hold on within the datacenter. thus it clearly decreases the dimensions of datacenter. therefore it primarily implies that the amount of the replicated copies of knowledge that were typically duplicated on the cloud server may be controlled and managed simply to shrink the physical space for storing. The recent statistics has recognized that deduplication is that the most influential storage technology and is expected to supply seventy fifth of all backups within the next few years.The potency of any knowledge Deduplication application may be effectively measured by the one.

Dedupe quantitative relation wherever Dedupe Ratio= Size of Actual Data/Size of knowledge once Deduplication two. turnout (Megabytes of knowledge

Deduplicated per sec).Following square measure the parameters that govern the dedupe quantitative relation and throughput-

1. Nature of knowledge to be deduplicated.
2. wherever is that the Deduplication applied?-either on supply device or target device
3. If knowledge Deduplication is inline or a post process application
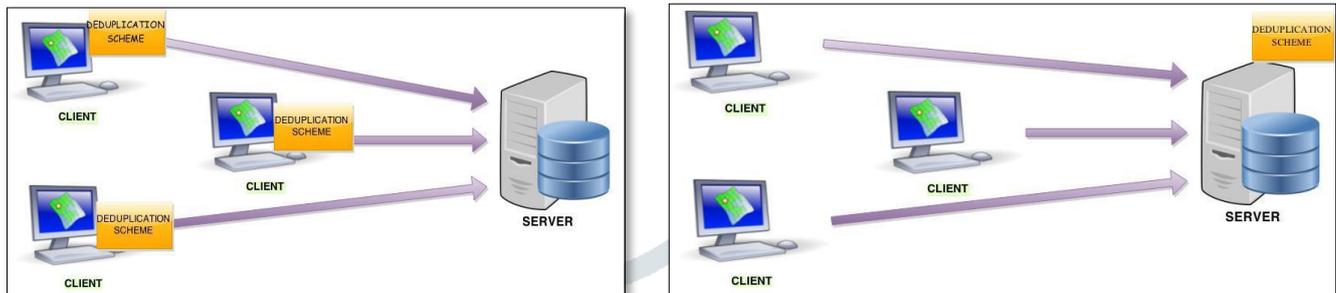4. Implementation of knowledge Deduplication[4]



Fig 1 Source and Target Deduplication

Some advantages of deduplications square measure

1) Storage-based knowledge deduplication ready to scale back the number of storage needed for a given set of files. it's simpler wherever several copies of terribly similar or perhaps identical knowledge square measure hold on on one disk.
2) To scale back the amount of bytes that has got to be transferred between endpoints ,Network knowledge deduplication is employed which may scale back the number of information measure needed.
3) Virtual servers and virtual desktops like deduplication as a result of files for every virtual machine to be fused into one space for storing that permits nominally separate system
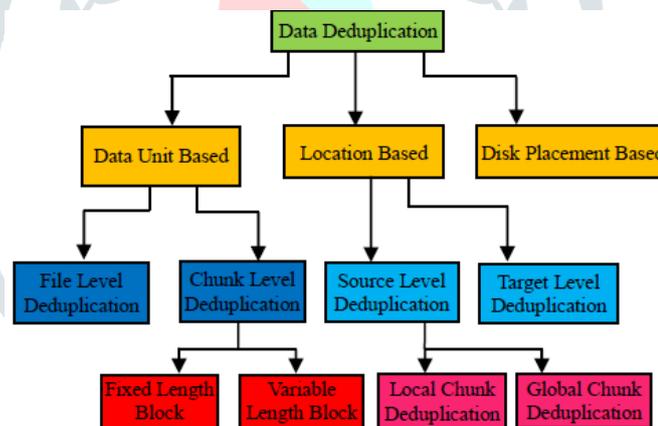


Fig 3 Classification of Deduplication Process

Deduplication could occur "in-line", as knowledge is flowing, or "postprocess" when it's been written.

(1) Post-process deduplication

In post-process deduplication, knowledge is 1st hold on on the device and so method at a later time can analyze the info searching for duplication. The profit is that before storing the info there's no have to be compelled to look ahead to the hash calculations and operation to be completed , thereby guaranteeing that store performance isn't degraded.

(2) In-line deduplication

Alternatively, as knowledge enters the target device deduplication hash calculations may be tired period of time .Only a relevancy the prevailing block is hold on, instead of the total new block if the storage system identifies a block that it's already hold on .The advantage of in-line deduplication over post-process deduplication is that since duplicate knowledge isn't hold on it needs solely less storage . On the negative facet, it's oftentimes argued as a result of hash calculations and lookups take ciao and knowledge body process may be slower, thereby reducing the backup output of the device

## III. LITERATURE REVIEW

. This paper addresses [1] this difficult open issue by, on one hand, process and implementing access policies supported knowledge attributes, and, on the opposite hand, permitting {the knowledge|the info|the information} owner to delegate most of the computation tasks concerned in fine-grained knowledge access management to untrusted cloud servers while not revealing the

underlying data contents. we tend to accomplish this goal by exploiting and unambiguously combining techniques of attribute-based coding (ABE), proxy coding, and lazy coding. Our projected theme additionally has salient properties of user access privilege confidentiality and user secret key responsibility. intensive analysis shows that our projected theme is very economical and demonstrably secures beneath existing security models.

Secure source that records possession and method history of knowledge objects is important to the success of knowledge forensics in cloud computing, however it's still a difficult issue nowadays. [2] during this paper, to tackle this unknown space in cloud computing, we tend to projected a brand new secure source theme supported the additive pairing techniques. because the essential bread and butter of knowledge forensics and post investigation in cloud computing, the projected theme is characterised by providing the data confidentiality on sensitive documents hold on in cloud, anonymous authentication on user access, and source trailing on controversial documents. With the demonstrable security techniques, we tend to formally demonstrate the projected theme is secure within the customary model.

With the character of low maintenance, cloud computing provides a cost-effective and economical resolution for sharing cluster resource among cloud users. sadly, sharing knowledge Associate in Nursing exceedingly|in a very} multi-owner manner whereas conserving knowledge and identity privacy from an untrusted cloud remains a difficult issue, thanks to the frequent amendment of the membership. [3] during this paper, we tend to propose a secure multi-owner knowledge sharing theme, named Mona, for dynamic teams within the cloud. By investment cluster signature and dynamic broadcast coding techniques, any cloud user will anonymously share knowledge with others. Meanwhile, the storage overhead and coding computation price of our theme ar freelance with the amount of revoked users. additionally, we tend to analyze the safety of our theme with rigorous proofs, and demonstrate the potency of our theme in experiments
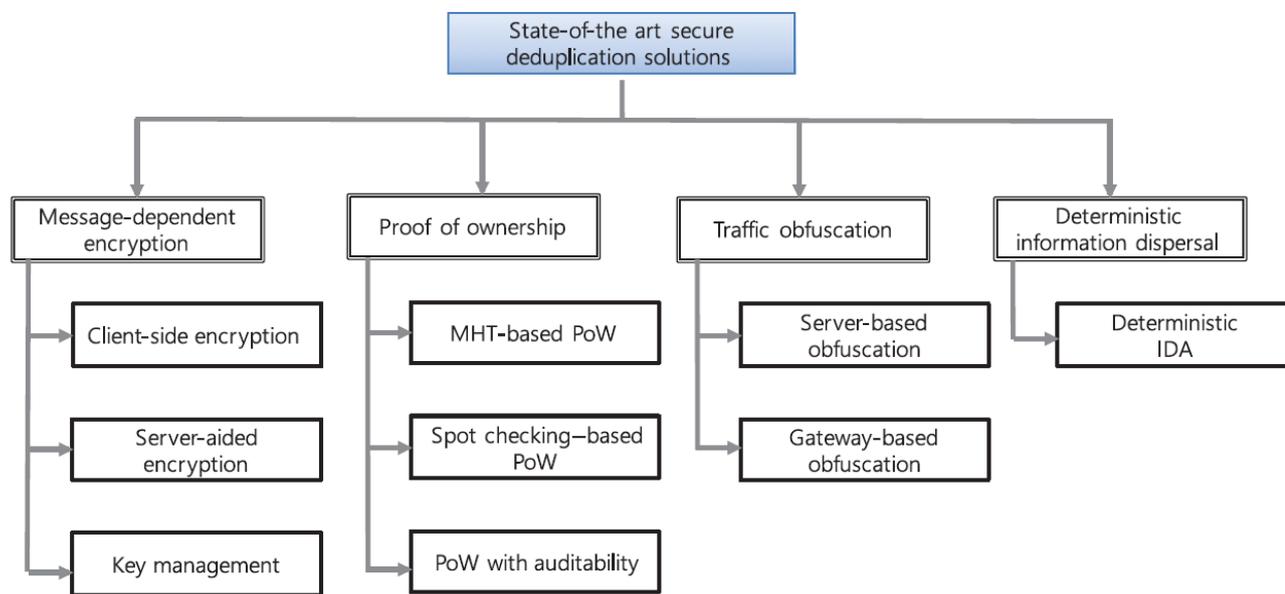


Fig. 4 state-of-the-art secure deduplication solutions

It is fascinating to store knowledge on knowledge storage servers like mail servers and file servers in encrypted type to cut back security and privacy risks. however this typically implies that one should sacrifice practicality for security. as an example, if a consumer desires to retrieve solely documents containing bound words, it absolutely was not antecedently best-known the way to let the info storage server perform the search and answer the question, while not loss of knowledge confidentiality. we tend to describe our science schemes for the matter of looking out on encrypted knowledge and supply proofs of security for the ensuing crypto systems. [4] Our techniques have variety of crucial benefits. they're demonstrably secure: they supply demonstrable secrecy for coding, within the sense that the untrusted server cannot learn something regarding the plaintext once solely given the ciphertext; they supply question isolation for searches, that means that the untrusted server cannot learn something a lot of regarding the plaintext than the search result; they supply controlled looking out, so the untrusted server cannot hunt for Associate in Nursing discretional word while not the user's authorization; they additionally support hidden queries, so the user could raise the untrusted server to go looking for a secret word while not revealing the word to the server. The algorithms given ar easy, quick (for a document of length n, the coding and search algorithms solely want O(n) stream cipher and block cipher operations), and introduce nearly no house and communication overhead, and thus ar sensible to use nowadays.

Searchable interchangeable coding (SSE) permits a celebration to source the storage of its knowledge to a different party (a server) in an exceedingly personal manner, whereas maintaining the power to by selection search over it.[5] This drawback has been the main target of active analysis in recent years. during this paper we tend to show 2 solutions to compass point that at the same time fancy the subsequent properties:

Both solutions ar a lot of economical than all previous constant-round schemes. above all, the work performed by the server per came back document is constant as hostile linear within the size of the info.

Both solutions fancy stronger security guarantees than previous constant-round schemes. In fact, we tend to suggests delicate however serious issues with previous notions of security for compass point, and show the way to style constructions that avoid these pitfalls. Further, our second resolution additionally achieves what we tend to decision adaptive compass point security,

wherever queries to the server may be chosen adaptively (by the adversary) throughout the execution of the search; this notion is each necessary in follow and has not been antecedently thought of.

Surprisingly, despite being safer and a lot of economical, our compass point schemes ar remarkably easy. we tend to take into account the simplicity of each solutions as a very important step towards the readying of compass point technologies. As an extra contribution, we tend to additionally take into account multi-user compass point. All previous work on compass point studied the setting wherever solely the owner of the info is capable of submitting search queries. we tend to take into account the natural extension wherever Associate in Nursing discretional cluster of parties aside from the owner will submit search queries. we tend to formally outline compass point within the multi-user setting, Associate in Nursing gift an economical construction that achieves higher performance than merely exploitation access management mechanisms

| Operations and vendors | Amazon S3 [14] | Windows Azure [15] | Aliyun [16] | RackSpace [17] |
|---|---|---|---|---|
| Storage (per GB/month) | $0.033 | $0.157 | $0.029 | $0.13 |
| Data in (per GB) | Free | Free | Free | Free |
| Data out to internet (per GB) | $0.201 | Free | $0.123 | Free |
| Put, copy, post, and list (per 10k transactions) | $0.047 | Free | $0.0016 | Free |
| Get and others (per 10k transactions) | $0.0037 | Free | $0.0016 | Free |

Most standard searchable coding schemes suffer from 2 disadvantages. First, looking out the hold on documents takes time linear within the size of the info, and/or uses significant arithmetic operations.[6] second, the prevailing schemes don't consider adaptive attackers; a search-query can reveal data even about documents hold on within the future. If they are doing take into account this, it's at a big price to the performance of updates. during this paper we tend to propose a completely {unique|a unique} interchangeable searchable coding theme that provides looking out at constant time within the variety of unique keywords hold on on the server. we tend to gift 2 variants of the fundamental theme that disagree within the potency of search and storage. we tend to show however every theme may be utilized in a private health record system.
.

Table 1 Comparison of existing deduplication technique

| Techniques | Metadata Processing | Chunking Method | Chunk Granularity | Dedupe Scalability |
|---|---|---|---|---|
| Deduplication using Byte-index Chunking method | Index-matrix table | Fixed size | Chunk | Small scale storage |
| Probabilistic Deduplication for cluster based storage | Bitmap Vector | Content Based | Super-chunk | Cluster based storage |
| Scalable Deduplication using data routing technique | Similarity Index | Variable Length | Super chunk | Cluster based storage |
| Deduplicaton using multi-layer metadata | Tree map Global and local metadata | Variable length | Chunk | Small Scale cloud based storage |
| Application aware deduplication for cloud backup | Local and global metadata | All types of chunking methods | Chunk | Large scale cloud based storage |

Table 2 Literature survey of different papers

| Paper Title | Method | Publication | Summary |
|---|---|---|---|
| Secure Storage as a Service in Multi-Cloud Environment | AES-256,RSA | Ad-hoc, Mobile, and Wireless Networks. ADHOC-NOW 2017. Lecture Notes in Computer Science, vol 10517. Springer, Cham | This paper is work on to create secure storage in multi-cloud environment architecture which handles confidentiality and integrity issues. |
| A Multi-Cloud Approach for Secure Data Storage on Smart Device | | 2016 Sixth International Conference on Digital Information and Communication Technology and its Applications (DICTAP) | This paper gives a secure data storage for mobile cloud computing using multiple cloud based system |
| MECCAS: Collaborative Storage Algorithm Based on Alternating Direction Method of Multipliers on Mobile Edge Cloud | hyper-graph model theory, Bayes-theorem-based heuristic | 2017 IEEE 1st International Conference on Edge Computing | This paper work on alternating direction-Method -of-multipliers-based collaborative storage algorithm named MECCAS (Mobile Edge Cloud Collaborative Storage).This algorithm compares with other three algorithm ADM, RDM, ERASURE |

| | method for workload distribution. | | |
|---|---|---|---|
| "Secure scheme on mobile multi cloud computing based on homomorphic encryption," | a multi-factor authentication method, an authentication scheme using MDA | 2016 International Conference on Engineering MIS (ICEMIS), | This paper gives a security architecture that generate homomorphic signature to ensure the integrity |
| "A mobile cloud middleware for data storage and integrity," | asymmetric encryption techniques | 2015 International Conference on Cloud Technologies and Applications (CloudTech) | This paper e present a new mobile cloud middleware, which aims to provide the mobile clients with an extendible storage area and data integrity service. |
| "A Cost Efficient Multi-cloud Data Hosting Using Heuristic Data Placement Algorithm" | CHARM method | 2018 IEEE | The system will become more powerful if it combine for block level as well as file level deduplication check. This will overcome the drawback that if whole file is same then also here it check the duplication chunk by chunk. It is avoided if we first perform the file level deduplication and if file is unique then only go for the block level deduplication check. Also checked for another encryption algorithm to provide the more security. |
| "An Efficient and Secure Deduplication Scheme Based on Rabin Fingerprinting in Cloud Storage" | Rabin fingerprinting in cloud storage | 2017 IEEE International Conference on Computational Science and Engineering (CSE) and on Embedded and Ubiquitous Computing (EUC) | blocks file by using the Rabin fingerprinting, which supports various changes in the file. Specially, the proposed scheme introduces the trusted third-party server to add the secret information for convergent key randomization. It prevents the offline brute-force dictionary attacks. The scheme allows cloud storage server to perform block-level deduplication before blocks are encrypted by users. This can significantly avoid invalid encryption and hence improve the efficiency of deduplication. |
| "Implementing Deduplication Technique for RDF Files with Enhanced Security using Multi Cloud servers" | Privileges of users and duplicate check. | International Conference on Energy, Communication, Data Analytics and Soft Computing (ICECDS-2017) | In this project we realize a prototype of our considered authorized duplicate check scheme and conduct test bed experiments on our prototype. From this project we show that our sanctioned duplicate check scheme acquire negligible overhead balance to convergent encryption and network relocate. |

## IV. PERFORMANCE PARAMETERS

Following are vital metrics for achieving deduplication efficiency:

4.1 De-duplication potency

The unitization algorithmic rule ought to take less time to get chunks. Variable size unitization is additional economical compare to fastened sized unitization. De-duplication outturn the speed at that the chunks are with success generating. a perfect unitization algorithmic rule ought to have a high outturn, which suggests it ought to generate the chunks at quicker rate.

4.2 procedure overhead

To have higher deduplication potency, the computation overhead ideally ought to be low. once the info stream is processed within the deduplication engine, each computer memory unit must check to search out chunk boundaries. This cause monumental use of central processing unit resources which cause high procedure overhead [3].

4.3 Chunk size variance

It indicates variance within the chunk size. The chunk size shouldn't be too little or overlarge. This variance in chunk size affects the deduplication potency that results in low unitisation outturn. The smaller the variance in chunk size, the higher the

deduplication potency [3].

### 4.4 Low-entropy string

In some strings it's tough to search out its chunk boundaries. These strings will have same characters repetitively, that makes it tough to search out chunk boundaries. The unitization algorithmic rule ought to ready to conclude the low entropy strings [3].

## V. CONCLUSION

In this paper we've conferred a review on knowledge deduplication challenges in unitization method, quantifiability, throughput, information process, parallelizing dedupe method and deploying knowledge deduplication on cluster so as to realize sensible deduplication performance. For the comparison of unitization algorithms and deduplication magnitude relation with them, initial we tend to were ready to divide knowledge into fastened and variable size chunks exploitation unitization algorithms with most potency and fewer time. Second, with the utilization of correct hashing algorithmic rule and deduplication method we tend to were with success ready to comment that variable-size unitisation provides higher deduplication magnitude relation as compared to fixed-size unitization. unitization time needed with switch divisor algorithmic rule was diminished more or less to twenty fifth. we tend to conjointly propose to pose the dedupe method to realize higher outturn

## REFERENCES

[[1] Andre Brinkmann, Sascha Effert,"Snapshots and Continuous Data Replication in Cluster Storage Environments",Fourth International Workshop on Storage Network Architecture and Parallel I/O, IEEE,2008.

[2] Q. Liu, Y. Fu, G. Ni, R. Hou,"Hadoop Based Scalable Cluster Deduplication for Big Data" ,2016 IEEE 36th International Conference on Distributed Computing Systems Workshops.

[3] N Kumar, R. Rawat, and S. C. Jain,"Bucket Based Data Deduplication Technique",5th International Conference on Reliability, Infocom Technologies and Optimization (ICRITO) (Trends and Future Directions), Noida, 2016, pp. 267-271.[4] Z. Sun, J. Shen, and J. Yong,"A novel approach to data deduplication over the engineering-oriented cloud systems",Integrated Computer-Aided Engineering, vol. 20, no. 1, pp. 45–57, 2013.

[5] A. Venish and K. Siva Sankar,"Study of Chunking Algorithm in Data Deduplication",'Springer India 2016. R. Vikraman and A. S," A Study on Various Data Deduplication Systems",International Journal of Computer Applications, vol. 94, no.4, pp. 35-40, 2014.

[6] A. Venish and K. Siva Sankar,"Study of Chunking Algorithm in Data Deduplication",'Springer India 2016. R. Vikraman and A. S," A Study on Various Data Deduplication Systems",International Journal of Computer Applications, vol. 94, no.4, pp. 35-40, 2014.

[7] George Crump Optimization(2011, September 30),Which Primary Storage is beast? [Online] Available: http://www.storage-switzerland.com

[8] E. Manogar and S. Abirami,"A Study on Data Deduplication Techniques for Optimized Storage",2014 Sixth International Conference on Advanced Computing(lCoAC), IEEE 2014, pp. 161-166.

[9] R-S Chang, C-S Liao, K-Z Fan, and C-M Wu," Dynamic Deduplication Decision in a Hadoop Distributed File System"International Journal of Distributed Sensor Networks, pp. 1-14, April 2014

[10] Min Xu, Yunfeng Zhu, Patrick P. C. Lee, Yinlong Xu," Even Data Placement for Load Balance in Reliable Distributed Deduplication Storage Systems" In Proc. of IEEE International Symposium on Quality of Service (IWQoS), pp. 349-358, 2015.

[11] Deepu S ,Bhaskar ,Shylaja,"PERFORMANCE COMPARISON OF DEDUPLICATION TECHNIQUES FOR STORAGE IN CLOUD COMPUTING ENVIRONMENT",Asian Journal of Computer Science And Information Technology 4 : 5 (2014) 42 - 46.

[12] Amanpreet Kaur,Sonia Sharma,"An Efficient Framework and Techniques of Data Deduplication in Cloud Computing",IJCST Vol. 8,April - June 2017.

[13] Shengmei Luo, Guangyan Zhang, Chengwen Wu," Boafft: Distributed Deduplication for Big Data Storage in the Cloud",IEEE TRANSACTIONS ON CLOUD COMPUTING, VOL. 4, NO. X, XXXXX 2016.

[14] Deepavali Bhagwat,Kave Eshghi,Darrell D. E. Long,"Extreme Binning: Scalable, Parallel Deduplication for Chunk- based File Backup",in Proc. IEEE Int. Symp. Modell. Anal. Simulation Comput. Telecommun. Syst., 2009, pp. 1–9.

[15] C. Liu, Y. Lu, C. Shi, et al.,"ADMAD: Application-driven metadata aware deduplication archival storage System",in Proc. 5th IEEE Int. Workshop Storage Netw. Archit. Parallel I/Os, 2008, pp. 29–35.7.