

DATA DEDUPLICATION IN CLOUD

CLOUD COMPUTING AND VIRTUALIZATION -(ITE3007)

REVIEW 2 PROJECT REPORT

Submitted by

ARIHANT JAIN

20BIT0006

SAHIL NANDAL

20BIT0010

RAJ RATNESH SINGH

20BIT0028

Submitted to

Dr Siva Rama Krishnan S



**SCHOOL OF INFORMATION TECHNOLOGY AND
ENGINEERING**

VIT UNIVERSITY: VELLORE 632 014

JANUARY, 2023

Contents

1. Abstract.....	3
2. Introduction.....	4
3. Literature Survey.....	5
4. Overview and Planning.....	15
4.1.Proposed System Overview.....	15
4.2.Challenges.....	15
4.3.Assumptions.....	15
4.4.Architecture Specifications.....	15
4.5.Hardware Requirements.....	16
4.6.Software Requirements.....	16
5. System Design.....	17
5.1.High Level Design.....	17
5.1.1. Sender Side.....	17
5.1.2. Receiver Side.....	17
5.2.Low Level Design.....	18
6. Implementation.....	19
6.1. Admin.....	19
6.2. Server.....	20
6.3. User.....	21
7. Summary.....	23
8. References.....	24

1. Abstract

Cloud storage is one of the cloud computing services that provides consumers with virtualized storage on demand. However, as the amount of data in the cloud grows, customers expect to be able to use on-demand cloud services at any time, while providers must maintain system availability and handle a significant amount of data. As a result, removing data redundancy has become one of the most important approaches in the cloud. Providers require a method to drastically reduce data quantities in order to save money while running large-scale systems. One of the reasons for the increase in the cost of cloud hosting services can be the bandwidth pricing on the usage-based whereas in the pay as you go service model helps in decreasing the bandwidth usage cost. One of the points of attention is the bandwidth cost which is a serious topic of concern. Research on Traffic Redundancy Elimination (TRE) [31] technologies have been done to decrease the bandwidth cost while data transfer. TRE technologies eliminate the transmission of duplicate information. TRE has garnered a lot of attention in the cloud environment due to its excellent rate of reduction in bandwidth cost. One of the ways for reducing data redundancy is to implement sender-based TRE or receiver-based TRE but it cannot capture all types of traffic redundancy as traffic can be either short-term (time span of seconds or minutes) or long-term (time span can hours or days). Moreover, the receiver-based TRE solution is inefficient when the data changes occur which creates a need for a new method which is more efficient and the solution proposed for that is Collaborative end-to-end (CoRE) [32].

2. Introduction

The advent of the 20th century brought an immense increase in digital computation and, along with it, the start of an information era. This has further led to the development of research to ensure easily accessible data, which led to the introduction of cloud computing as a basic building block for the same. Cloud computing is a synthesis of numerous computing domains that allows for lower capital expenditures and greater flexibility in resource supply, scaling up, and scaling down. This is where cloud storage plays a major role in providing any kind of service. Because of benefits and characteristics such as multi-tenancy, improved server utilization, energy efficiency, and elasticity obtained from on-demand utility computing services, cloud computing use has increased tremendously. Data deduplication is a type of data compression technique that eliminates multiple copies of repeated data. It is also known as intelligent compression or single instance storage. This approach is used to enhance storage utilization and can also be used to reduce the amount of bytes that must be delivered during network data transfers. During the deduplication process, unique blocks of data, or byte patterns, are discovered and stored. Deduplication techniques use data similarity to locate duplicate data and save storage space. Data deduplication is so closely related to incremental backup, which transfers just the data that has changed since the previous backup.

3. Literature Survey

Manju Bala et al. [1] In the given paper, the authors suggested deduplication as the strategy used to make efficient use of the space available for storing data in the cloud. The authors mention data de-duplication as the effective technique used to reduce the space occupied by reducing duplicates while maintaining data integrity. The paper presents various data deduplication techniques in reference to chunking (file-level chunking and block-level chunking), location (source and target), and time (inline and post-process). The choice of deduplication methodology depends on the needs of the user. The paper compares the different data deduplication strategies prevalent today. The paper concludes by mentioning the future work of the paper, which includes exploring different security techniques to secure data during deduplication and to execute deduplication in a more secure manner.

Ashok Kumar S et al. [2] The paper discussed the role of storage capacity optimization in the efficient use of cloud computing. The paper makes use of data deduplication as a strategy for efficient storage utilization. Redundant data is recognized by making use of hash functions which identify unique sequences (which are compared to the other sequences already calculated using the function). Thus, this allows the authors to reference the previously stored data rather than replicate it again. SHA3 calculation is used in the paper for hashing. When compared to the SHA-3 algorithm, MD5 takes up a lot of spare space. The MD5 hash capacity implementation is quite limited in terms of memory usage, time differentiation, and SHA-3 computation. SHA-3 aids in recovering significant plate space while also enhancing capabilities. SHA-3 is the most astounding in recognising the richness of cloud records space. The work done in the paper can be used to reduce excess space taken up by the same data resources by a large percentage. The security of the data is also taken care of while evaluating the use of the hashing function in data deduplication. Thus, storage space is utilized effectively, hence reducing the cost of using the services.

Junbeom Hur et al. [3] The proposed paper makes use of the Message-Locked Encryption technique to tackle the problem of performing secure data deduplication on encrypted data. This problem arises because of the difference in cipher texts obtained while using different algorithms and keys (here, the plain text/actual test remains the same whereas the cipher text differs). The paper also proposed many block-level deduplication schemes which are used to provide fine-grained storage by dividing data into blocks, wherein, even if some data changes in a particular file, only that data's index is changed in the index table. Thus, it ensures efficient use of storage. The paper also introduces the readers to the various attacks taking place on Message-Locked Encryption schemes when the message set becomes predictable. The performance of the proposed model is calculated and compared with UMLE. As a result, this study provides a novel and safe deduplication technique that guarantees successful updating of data and resistance to brute-force assaults, even when the messages obtained are predictable. The paper also assesses and indicates that increasing block size reduces bandwidth and time utilization. The system also ensures the acceptability and usefulness of the proposed model in real-world cloud scenarios.

Mauro Conti et al. [4] The given paper introduces the readers to the creation of side channels, which aids in providing the details of the file existence status to the attackers as a result of client-side data deduplication. The paper mentions that client-side data duplication is useful for reducing bandwidth and storage costs at the client-side but introduces vulnerability in the system, leading to intelligent attacks when the vulnerability is exploited by the attackers. The paper mentions the weakness of the contingency plans available before the paper was published and proposes a RANdom REsponse model (RARE) for achieving greater privacy and security of data, thereby reducing the number of vulnerabilities paving the way for attacks by hackers. The methodology of the RARE model involves the user sending two data chunks at once. The authors mention that the cloud that receives the deduplication request responds with a randomized deduplication response that has been carefully designed to preserve the deduplication advantage while minimizing privacy leakage. The analytical results of the given study corroborate the privacy guarantee, and the results show that both the deduplication benefit and the privacy of the RARE model can be preserved.

Jie X et al. [5] The given paper includes the role of cloud computing and the benefits provided by using cloud computing. The paper by the author discusses in length the cloud deployment models (software as a service, platform as a service, infrastructure as a service), types of clouds, and their benefits. The usage of cloud (private, public, hybrid, community) according to the need is also discussed in detail. It is emphasized that the cloud provides a platform for implementing big data analysis methodologies by different firms to enhance the utilization of data and the generation of useful information from the given unstructured datasets. The method of data deduplication is used in the paper, which aims at reducing the data storage by reducing duplication. The paper proposes a model which includes load balancer, deduplicators, cloud storage and redundancy manager for managing duplication in data and reducing redundancies. The given paper makes use of CloudSim and HDFS simulators, which work as simulating environments and are Java-based toolkits. Hash functions are used for deduplication and various sized files are used for conducting the experiments. The paper concludes by evaluating the percent reduction in size of the files and the time taken for performing deduplication. When the number of upload files is increased to a thousand, five and ten deduplicators can still aid to minimize processing time, but their effectiveness is reduced to 91.40% and 95.58%, respectively. However, time savings diminish dramatically when the number of upload files is increased to ten thousands, as five and ten duplicators may minimize 60.10% and 79.71% of processing time, respectively.

Zhenhua Liu et al. [6] The proposed paper uses a randomized client-side deduplication technique to tackle the problem of the security of data. The paper discussed the challenges proposed by the ever-increasing amount of data and the impact of the same on the storage capacity of the cloud environment, which therefore, makes it necessary to ensure that techniques like data deduplication are employed in order to reduce the amount of redundant data uploaded on the cloud. But, adoption of hashing methods can introduce some other attacks which result in threat to the data. Thus, the paper proposed methodology to tackle the proposed problem. The paper discussed the preliminaries, which include discrete logarithm problem, the lifted ElGamal encryption, hash collisions and static KEK tree. The problem formulation made

in the paper included cloud storage system, the security model and the adversarial model and the security requirements. The proposed scheme of the paper included the main idea and the construction (system setup, data upload, ownership management, data encryption and data download). The security analysis employed by the paper included security proof, privacy, integrity, forward secrecy and backward secrecy, CAA and BFA resistance, and subsequently the paper discussed the comparisons and efficiency of the proposed model.

Im-Yeong Lee et al. [7] The paper discussed the importance of data deduplication strategy in improving the utilization of cloud storage and enhancing the amount of data that can be uploaded to the cloud. The paper makes use of comparisons to tackle the problem of security while making use of data deduplication as a strategy to reduce redundant data to improve data storage in the cloud. The paper classified data deduplication according to server-side, client-side and appliance side usages. Deduplication was classified according to different levels in the paper which included file-level deduplication, block-level deduplication (fixed length and variable-length). The paper also discussed the secure data deduplication methods and its importance in maintaining data security while reducing redundancies in data. Various attacks are also discussed in the paper, which included dictionary attack, poison attack, ownership forgery attack and thereafter various technologies vital for secure data deduplication (eg convergent encryption, oprf, mle data encryption, proof of ownership, preparing for poison attacks, proof of ownership and bloom filter) were discussed in depth. Technologies such as DupLESS, CloudDedup, PerfecDedup and Hur et al.'s scheme were also mentioned in the paper. The authors of this study explored ways for doing data deduplication as well as techniques for performing secure data deduplication. These technologies were not totally changed, but certain modified security technologies were adapted to the necessity for data deduplication. While these technologies have achieved their original goals, it is difficult to expect high levels of completeness while fulfilling both computation and traffic efficiency due to the advent of more security threats and the adoption of other security technologies to counter them. However, the recent desire for privacy and secure technologies is projected to enhance the demand and supply of safe data deduplication technologies in the future, leading to more advanced technology research.

Alexandra Shulman-Peleg et al. [8] The paper discussed the use of data deduplication as a strategy to reduce data duplication as a technique to reduce data redundancy while uploading user data to the cloud. This helps in saving cloud storage, hence improving space utilization. The research showed how deduplication can be used as a side channel to divulge information about the contents of other users' files. In another situation, as the paper discussed, deduplication can be used as a covert channel through which malicious software communicates with its command and control center, regardless of any settings, especially firewall settings on the attacking the system employed/used by the user. Cloud storage providers are unlikely to abandon cross-user deduplication due to the significant savings it provides. As a result, we suggest simple strategies for enabling cross-user deduplication while significantly minimizing the danger of data leakage.

Weiwei Chen et al. [9] The paper discussed in length the data deduplication techniques employed to tackle the problem of redundant data storage in the cloud. The authors discussed the different distributed data deduplication methodologies. The paper introduced the readers to the difference in application redundancy analysis paving way to the super-chunk resemblance analysis. The AppDEDUPE design (including the design principles: throughput, capacity and scalability) were discussed in depth. The overview of the AppDEDUPE included the clients (data partitioning, chunk fingerprinting, two-tiered data routing), director, (fire recipe management, application-aware routing decision) and the dedupe storage nodes (application-aware similarity index lookup, chunk fingerprint caching, parallel container management. The overview of the entire system included the clients, dedupe storage nodes and the director). The complete system introduced in the paper was evaluated based on the evaluation metrics, which included the deduplication efficiency, normalized deduplication ratio, normalized effective deduplication ratio, number of fingerprint index lookup messages, usage of RAM for intra-node data deduplication and data skew for distributed storage, which made use of MaxLoad, MinLoad and MeanLoad values for calculation of the performance. As a result, the authors conclude by stating that real-world trace-driven evaluation clearly revealed AppDedupe's significant advantages over current distributed deduplication algorithms for big clusters.

Robert H. Deng et al. [10] The authors developed a Bloom filter-based location selection approach and a safe data deduplication scheme in this research. The paper discussed the use of re-encryption methodology for efficient data deduplication. The preliminaries discussed in the paper include the convergent All-or-nothing transform, ciphertext-policy attribute-based encryption, and the bloom filter. The analysis of the reed scheme included the review and the sub-reserved attack. The models and security goals of the paper included the system model, threat model, and three folded goals, which ensure integrity of data. The scheme introduced in the paper included a bloom filter-based location selection mechanism and package generation. The analysis of the technique introduced in the paper is done by comparing different data deduplication techniques. The authors further demonstrated that the proposed approach achieves the desired security objectives and provided comprehensive simulation testing. The results of the experiments demonstrated that the technique was effective at re-encryption.

Zheng, Qinghua et al. [11] The paper discussed the storage management of heterogeneous data with data deduplication in a cloud computing context. The authors proposed a model for providing data deduplication along with simultaneous access control across various service providers of the cloud. The system and security model of the proposed paper included key generation centers, cloud service providers, authorized parties, data holders, proxy re-encryptions (PRE) for converting the cipher text to one that can be easily decrypted at a proxy and attribute-based encryption. The system design included fundamental algorithms for system setup, data encryption and decryption, symmetric key management, partial key control (based on ABE), partial key control (based on PRE), the deduplication schemes, etc. The paper also evaluated the efficiency of the proposed model and in addition, compares it with pre-existing models. As stated in the study, as future work, the authors will strengthen user privacy and improve the performance of our approach in preparation for real deployment. Furthermore, the

authors will perform game theory analysis to further demonstrate the rationality and security of the suggested scheme.

Huafei Zhu et al. [12] The paper introduced the notion of data deduplication for secure data storage. The authors demonstrated that the proposed private data deduplication protocol is safe if the underlying hash function is collision-resilient, the discrete logarithm is hard, and the erasure coding method can erase up to ϵ -a fraction of the bits in the presence of malevolent adversaries. The authors also provided correctness and sound proof of the proposed system in the paper. In this research, the authors introduced a novel concept known as private data deduplication protocols in the context of two-party computations. A possible outcome of private data deduplication procedures has been developed and assessed in the paper.

Vakalapudi Bhavya Sri Kanthi et al. [13] The paper presented a system of secure data deduplication using advanced encryption cryptography standards. The paper presents a theoretical analysis of the AES system. The work proposed in the paper included user registrations, file uploading, encryption, checking for duplicates, and the final data dump diagram of the proposed model. The execution and assessment of encryption strategies resulted in the development of a more powerful encryption algorithm, whereas the deduplication strategy allowed the authors to save space in the cloud, cutting prices. Only clients who have been granted authorization will be able to encrypt and decrypt data using the keys generated during encryption. The paper experimented with and addressed the shortcomings of data encryption systems and performance.

Jin Li et al. [14] The paper presented data deduplication as the strategy for identifying and reducing redundancy in data. The paper discussed the method of validating image deduplication while storing data in the cloud. The authors' methodology in the study includes calculating the hash value of the image to be submitted by the user as its fingerprint. Following that, the fingerprint is transferred to cloud servers for validation and duplication detection. The methodology also involved the addition of a response mechanism by the storage and verification servers (no deduplication) and the subsequent transfer of data by the user to the servers. However, if the fingerprint is always discovered, the user will not upload the data for deduplication. This led to the conclusion that the server validation process was inefficient.

Huiba Li et al. [15] The paper discussed data deduplication in a cloud computing system and in addition, defines cloud computing as a paradigm shift and data deduplication as a way of increasing storage efficiency in the cloud in addition to increasing the system overload. The paper introduced the users to the various challenges and proposed the solutions possible for those challenges. The preliminary study done in the paper on data-deduplication included live data deduplication in an open-source cloud framework, extreme binning, data deduplication backup mechanism, SAM framework (used for backing-up cloud, which reduces overload on the server-side by using excess CPU power and storage). The report examined various recent studies on applying data deduplication techniques to cloud systems and highlighted the inadequacies of previous work. Furthermore, the writers provided many viable solutions. Because most earlier deduplication work focused on centralized backup systems, the authors

believed that their approach would pave the way for efficient deduplication for cloud computing environments.

Farzad Farnoud et al. [16] This paper provides an information-theoretic analysis of the performance of deduplication algorithms on data streams when repeated data segments are not necessarily exact copies. This paper introduces a source model in which probabilistic substitutions are considered and both the fixed-length scheme and the variable length scheme deduplication algorithms are studied. The fixed-length deduplication algorithm is shown to be unsuitable for the proposed source model as it does not take into account the edit probability while the conventional variable-length deduplication algorithm show that as source entropy becomes smaller, the size of the compressed string vanishes relative to the length of the uncompressed string, leading to high compression ratios. This paper takes the source model introduced by Niesen (the first person to perform an information-theoretic analysis of deduplication algorithms) as a reference which is incompatible with the current analysis.

Srilatha Puli et al. [17] This paper focus on the attacks from malicious clients that are grounded on the manipulation of data identifiers and attacks based on backup time and network traffic observation. This paper defines intra-user deduplication, inter-user deduplication, client-side deduplication, server-side deduplication while describing the advantages, disadvantages, need and evolution of deduplication in today's explosive growth of data. The deduplication scheme provided in the paper is simple and robust and is efficient in terms of storage space and bandwidth savings for both clients and cloud service provider, the data deduplication is considered at a file level granularity but the solution can be extended to the block level. Two-phase deduplication that combines both intra- and inter-user deduplication techniques by introducing deduplication proxies between the clients and the storage server is the approach taken.

Y. Toaf et al. [18] This paper deals with the management of systems that are often compressed by means of deduplication techniques that partition the input text into chunks and store recurring chunks only once. It describes the design choices made during the development of an approximate hash function, serving as the basic tool of the new suggested deduplication system and report on extensive tests performed on a variety of large input files. The idea underlying a deduplication system is to locate repeated data and store only its first occurrence. The chunk size may indeed have a major impact on the performance: if it is too small, the number of different chunks may be so large as to jeopardize the whole approach, because the data structure D might not fit into RAM, so the system might not be scalable. On the other hand, if the chunk size is chosen too large, the probability of getting identical chunks decreases: many instances of chunks might exist, that could have been deduplicated had the chunk size been chosen smaller, but which, for the larger chunk size, have to be kept. A possible solution is to look for similar rather than identical chunks. If such a similar chunk is located, only the difference is recorded, which is generally much smaller than a full chunk. The idea of the current work is to implement the required similarity by what we call an approximate hash scheme. In this paper the first concern was to verify that the proposed approximate hash indeed spreads its values evenly. Once this has been confirmed, check that this uniformity does not come at the price of

sensitivity, as it would for a standard hashing scheme. Thus checked the impact of the signature scheme in some artificial perturbation and clustering tests. Finally, bringing examples of applying the whole deduplication process in comparison with an identity based approach.

Mandeep Singh Devgan et al. [19] This paper provides the concepts, methods and the schemes that can make the cloud services secure and reduce the incidence of data duplication. In this paper the proposed scheme works for deduplication of data with arithmetic key validity operations that reduce the overhead and increase the complexity of the keys so that it is hard to break the keys. This paper describe that without key management and key validation(components of the user management) process deduplication would remain unsecure process and file would always remain under multiple thread including integrity loss and breach of privacy and data duplication can be taken care of either by minimizing the number of writes for saving I/O bandwidth or denormalization. This paper shows that at each level of duplication process (file, block, chunk, zone) there is a needs of keys to be arithmetically valid and there ownership also need to be proved for proper working of any secure (Source, Target, Semantic ,Local, Hardware etc.). Deduplication system. From this study, it can also be concluded that there is no absolute or perfect solution of deduplication.

Zhiqiang Yao et al. [20] This paper propose a novel secure role re-encryption system (SRRS), which is based on convergent encryption and the role re-encryption algorithm to prevent the privacy data leakage in cloud and claims that it also achieves the authorized deduplication and satisfies the dynamic privilege updating and revoking. This paper discusses three major issues first, the adversary may utilize the relative attack methods to intercept the user's privacy information second, in the data deduplication based on the traditional convergent encryption the unauthorized users can obtain the user's information only by supplying the hash value of the file third, the privilege of the authorized user is dynamic and flexible, it is difficult to guarantee the access permission of authorized user and achieve the key updating and revoking management when performing data deduplication in response to which SRRS is introduced. In this proposed system, firstly exploited the convergent encryption algorithm to prevent privacy data leakage and used the role re-encryption algorithm to achieve authorized deduplication efficiently. Specifically, created a role authorized tree to manage the user's roles and the corresponding role keys, and introduced the management center to reduce the computation cost and management overhead of the client, and implement the dynamic updating of the authorized user's privilege.

B. Haritha et al. [21] This paper defines cloud computing, cloud advantages, types of cloud while listing data deduplication advantages and uses. This paper mention two strategies to deduplicate excess information - Inline and post-processing deduplication. This paper evaluates deduplication at three levels: File-level data deduplication strategy, Block level data deduplication strategy, Byte level data deduplication technology.

B. Shanthini et al. [22] This paper examines the three-level cross-space design and propose an effective and protection safeguarding huge information deduplication in distributed storage which is referred to as EPCDD which accomplishes both protection safeguarding and

information accessibility and opposes beast power assaults and beats existing contending plans, as far as calculation, correspondence, and capacity overheads. In this paper the process is mentioned to begin by registering a user in database who then uploads the data from system to cloud which is pre-processed and deduplicated and once uploaded to cloud then provides proof of storage that the data is stored in cloud.

Shufen Niu et al. [23] This paper proposes a novel verifiable attribute-based keyword search over encrypted data supporting data deduplication to address the problems of data integrity detection, data deduplication and reasonable access authorization. This paper introduces a scheme to achieve effective access authorization and data confidentiality using attribute-based encryption. Searchable encryption is used to address the problem of the limited use of encrypted data. This paper details PRELIMINARIES AND SECURITY DEFINITIONS, SYSTEM AND SECURITY MODEL and the proposed system contains algorithms for system initialization, key generation, data encryption, trapdoor generation, search, result verification, transform and decryption. The results of the experiment reveal that the scheme has a low bandwidth communication, storage and computability.

Chen Chen et al. [24] This paper details a data routing strategy based on distributed Bloom Filter where to improve system throughput superchunk is used as the basic unit of data routing. The optimal node is selected as the routing node by matching the BloomFilter, and the storage capacity of the node and maintained in the memory of the storage node. In this paper the specific parameters of all kinds of routing strategies are obtained through experiments, and the routing strategies proposed are tested. The paper forwards with the cluster deduplication system as it has the advantage on improving storage space utilization and decreasing the reduplication rate, so the overall deduplication rate would decrease with the number of nodes in the cluster increases while the routing algorithm need to ensure the load balance of the system while keeping the deduplication rate. One of the challenges in the paper is how to send data to the storage nodes.

Rong Hao et al. [25] The author proposes a cloud storage auditing scheme with deduplication supporting strong privacy protection that ensures that the privacy of the user's file would not be disclosed to the cloud and other parties when this user's file is predictable or from a small space. This paper provides a new method for duplicate check by generating a file index and use new strategy to generate a key for file encryption. The paper elaborates the issues like, when the file is predictable or from a small space, Convergent Encryption cannot resist brute-force dictionary attacks, in which the malicious cloud can recover the entire file with a number of guesses which leads us to a secure auditing and data deduplication scheme where a key server is introduced to help user generate the convergent key. But comes up with an issue where the key server is able to guess or derive the file's content from the file's hash value sent from the user by launching the brute-force dictionary attacks making it unable to fully prevent brute-force dictionary attacks. This paper focuses on solving this issue by introducing new strategies like generating file index with the help of Agency Server(AS) and the key for file encryption is generated with the file and the file label and the file label is kept by the user secretly and in

order to improve the storage efficiency, the users, who own the same file, are able to generate the same ciphertext and the same authenticators.

Lucas Kencel et al. [26] This paper includes an encryption scheme that guarantees semantic security for unpopular data and provides weaker security and better storage and bandwidth benefits for popular data to increase the data deduplication efficiency. The paper proves the scheme is secure under the Symmetric External Decisional Diffie-Hellman Assumption in random oracle model.

M. V. P. Chandra Sekhara Rao et al. [27] In this paper, the authors have leveraged the advantages of perfect hash functions and made use of a probabilistic data structure to ensure the ownership of the data items. The paper defines data in two ways as files uploaded/owned by several users and files uploaded/owned by very few users. Files belonging to the former group will be benefited strongly by deduplication as they are popular and may not be sensitive from the confidentiality perspective. Files owned by very few users may contain very sensitive information and require stronger security, but will not benefit a lot by deduplication which bring out a major challenge in secure design of the scheme - the secrecy in verifying the popularity of the data files. The proposed scheme consists of CSP to hold data and will verify the ownership of the data before which the file will be distributed in blocks and Key server which will maintain the list of all unpopular data items and their identifiers. As per this the user will have a bit of freedom on deciding the popularity of the data which is being uploaded. The popularity and the ownership of the data will then be checked from the input of the other users. The proposed mechanism dealt with the conflict of interest between storage optimization and security of the data in the cloud storage, while focus on providing varying levels of security for the data items based on how popular the data is.

Parth Shah et al. [28] This paper provides an overview of different secure deduplication techniques in cloud storage like Provable Data Possession (PDP), Proof of Retrievability (POR), secure keyword search, DupLESS, Proof of Storage with Deduplication (PoSD), Dekey, Message-Locked Encryption, Attribute Based Encryption (ABE) and Identity Based Encryption (IBE). This paper do a case study and mentions the security issues and threats analytically that the users are facing while also stating the various vulnerabilities that can be exploited. The paper concludes that there are certain issues which remain to be resolved, one of them is storage and management overhead caused by cryptographic keys and also the issue of breach of client's key. Moreover, in case of single user and multiple users (cross user) it's more difficult to the implement edit and deletion operations with secure data deduplication solution and efficiency can be increased by the combination of file-level and block-level deduplication.

Zhiqiang Bai[et al[29] This paper introduce a novel server-side deduplication scheme for encrypted data in a hybrid cloud architecture, where a public cloud manages the storage and a private cloud manages the role of the data owner to perform deduplication and dynamic ownership management and to reduce the communication overhead an initial uploader is used for mechanism check to ensure only the first uploader needs to perform encryption, and adopt

an access control technique that verifies the validity of the data users before they download data. In the paper some of the issues are mentioned such as, to protect privacy, many users encrypt data before uploading it to the cloud storage where the encryption key is randomly generated, the same data encrypted with different keys will produce different ciphertext, which will hinder deduplication where the proposed solution is to revoke the cloud users from the valid ownership list. The proposed deduplication scheme works in the following way – the key is generated for the data and the private cloud gets the key, duplication check is performed by the public cloud, if the test is negative then the data is to be stored before which the private cloud encrypts the data, if the test is positive then the private cloud further clarifies the ownership of data by challenging hash code and then performs deduplication.

Mohammad R. Khosravi et al. [30] This paper develops a new method using Convergent and Modified Elliptic Curve Cryptography (MECC) algorithms over the cloud and fog environment to construct secure deduplication systems focusing on two goals - the redundancy of data needs to be reduced to its minimum and to develop a robust encryption approach. In the methodology in the paper when the user uploads a file the file is being tested if it is duplication if the test is negative the file is encrypted and if the test is positive then the metadata of the already saved file is sent to the user whereas, in case of downloading a file the access to the file is checked. The proposed methodology is analyzed in four ways, i.e., a) when a new user tries to upload a new file, b) when the same user tries to upload the same file c) when different users try to upload the same file to the cloud server and d) when the users try to download the file. The methodology make use of convergent encryption and convergent decryption. The implemented deduplication methodology is deployed in the JAVA programming.

4. Overview And Planning

4.1. Proposed System Overview

A system will be designed that takes data from the transmitter and forwards it to the receiver. The entering TCP data stream is broken into chunks at the receiver. The chunks are linked in sequence to form a chain and are then stored in a local chunk storage. Each arriving chunk is compared to the chunk storage by the receiver. Once a matched chunk is found on a chain, it retrieves a number of consecutive chunks down the chain as projected chunks in future incoming data. The retrieved chunks' signatures and expected offsets in the incoming data stream are provided to the sender in a PRED message as a prediction for the sender's upcoming outgoing data.

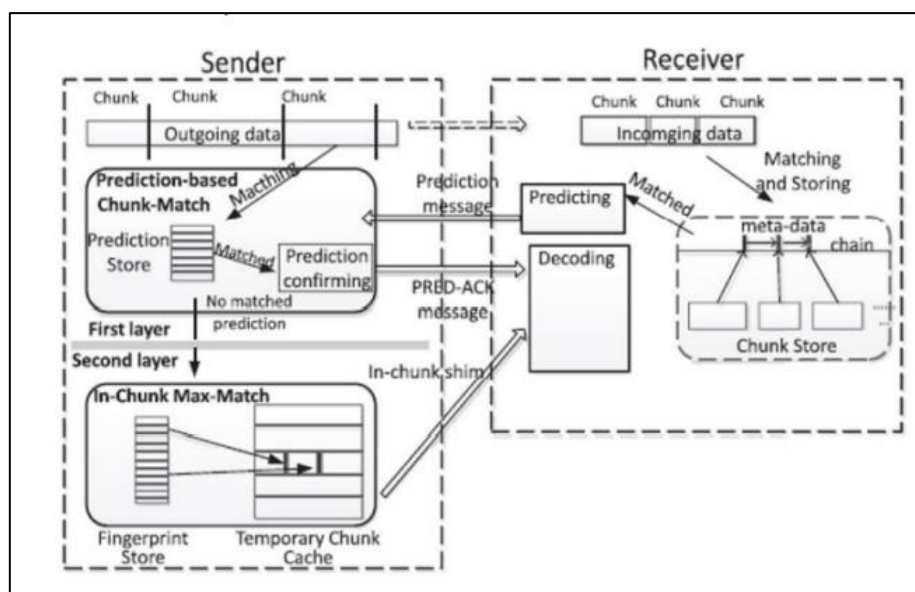
4.2. Challenges

- Faced issues while generating the hash values.
- Calculating the fingerprints for sub window lengths in the second layer.

4.3. Assumptions

Our present client-side CoRE solution assumes a single TCP flow from a server to each client. Because a server typically serves numerous clients continuously, we're curious how the server cost varies with the number of clients. We can see from the design that the server runs an instance of CoRE for each client; the cost of CoRE at the server should be additive and will increase linearly with the number of clients served by the server.

4.4. Architecture Specifications



CoRE contains two TRE modules, one for short-term redundancy and one for long-term redundancy. A two-layer redundancy detection system incorporates two TRE modules. The first-layer TRE module detects long-term redundancy in any outgoing traffic from the server. If no short-term redundancy is identified, it switches to the second-layer TRE module to look for it at a finer resolution.

4.5. Hardware Requirements

- **System** - Intel Core i3 or higher
- **Hard Disk** - 120 GB.
- **Monitor** - 15'' LED
- **Input Devices** - Keyboard, Mouse
- **RAM** - 4 GB

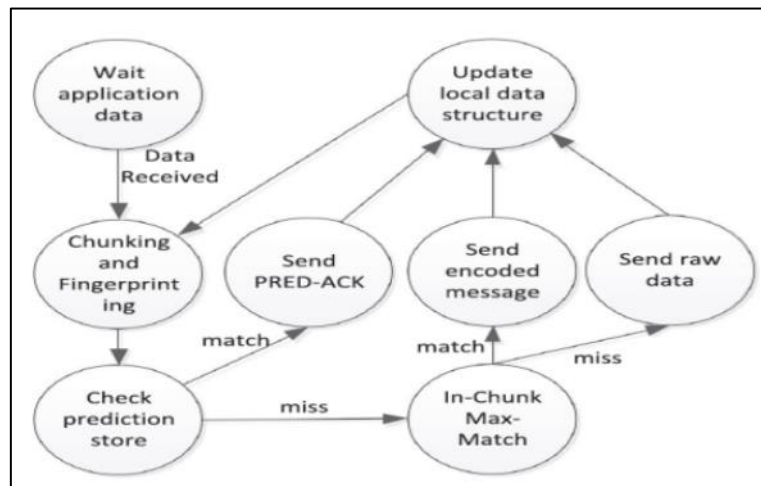
4.6. Software Requirements

- **Operating System** - Windows 7 or higher
- **Coding Language** - Java
- **Tools** - JAVA JDK1.8, NetBeans IDE 8.2

5. System Design

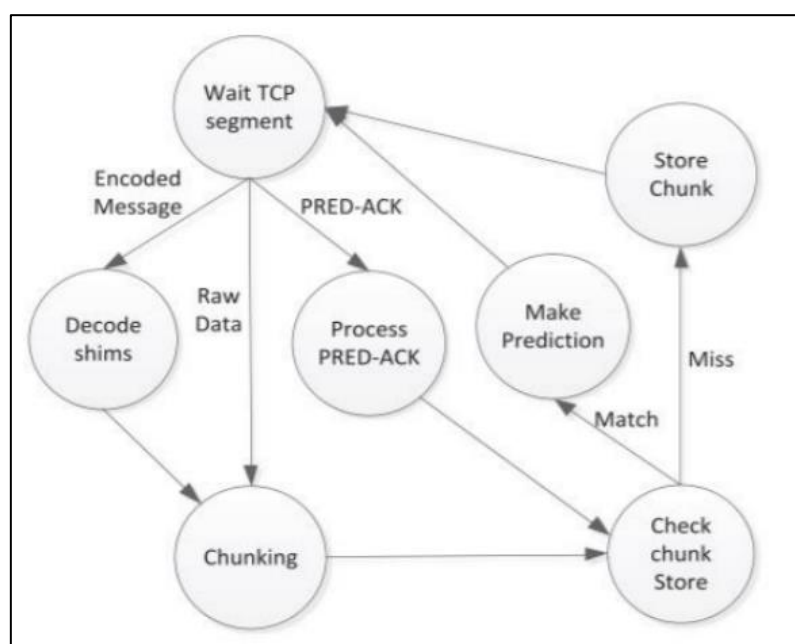
5.1. High Level Design

5.1.1. Sender Side



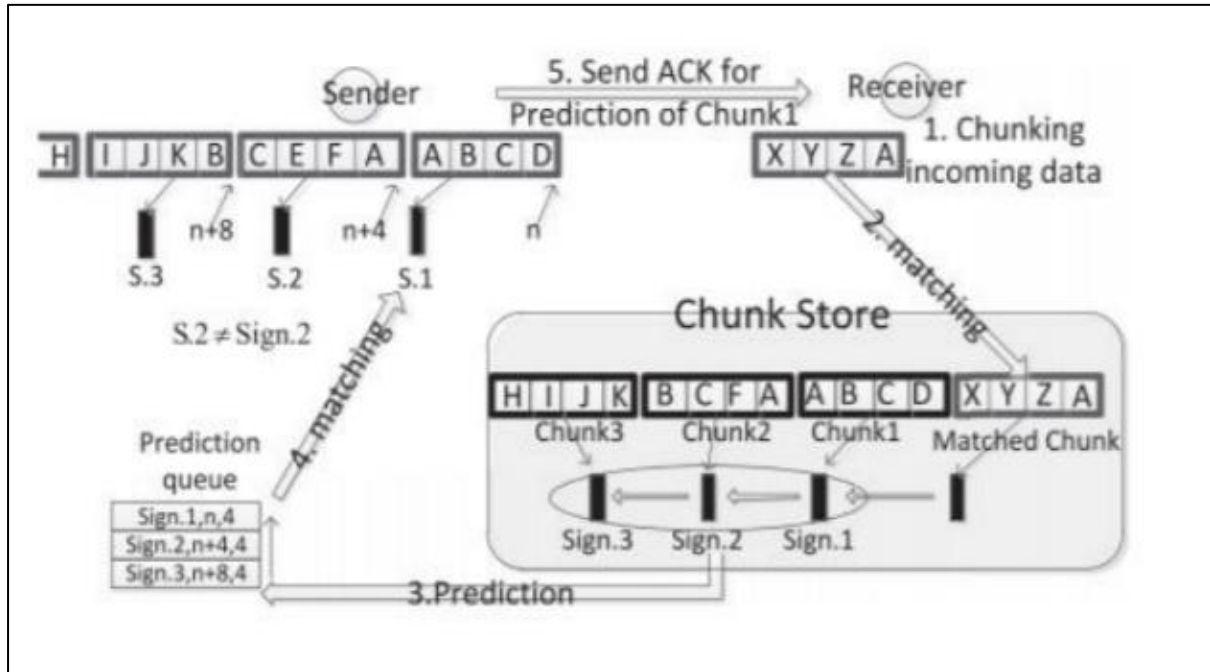
The sender employs a prediction store to cache the most recent predictions received from the receiver. Each prediction comprises the SHA-1 signature of a projected chunk as well as its expected offset, i.e., TCP sequence number in the TCP byte stream. The sender also has a chunk cache, which stores recently delivered chunks. A fingerprint store stores the meta-data of each representative fingerprint for each cached chunk, which contains the fingerprint value, the address of the chunk in the cache addressed by the fingerprint, and the byte offset of the window in the chunk across which the fingerprint is generated.

5.1.2. Receiver Side



The receiver separates the incoming data stream into pieces for each TCP connection and maintains a local prediction store that stores recent predictions for the TCP connection. To reconstitute the original data stream, the receiver sorts incoming TCP segments based on their kind. TCP segments from the sender to the receiver are classified into three types: PRED-ACK messages, shim-encoded messages, and bare data. The chunk predictions are sent in the order of their projected offsets within a PRED message.

5.2. Low Level Design

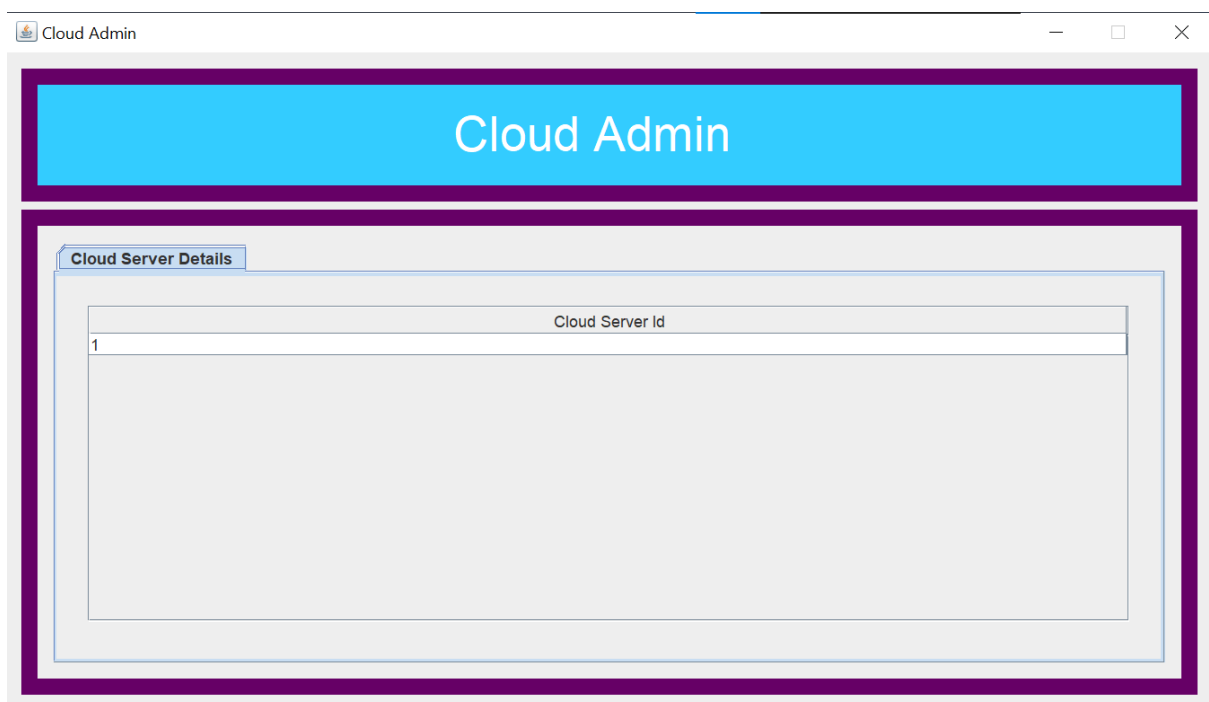
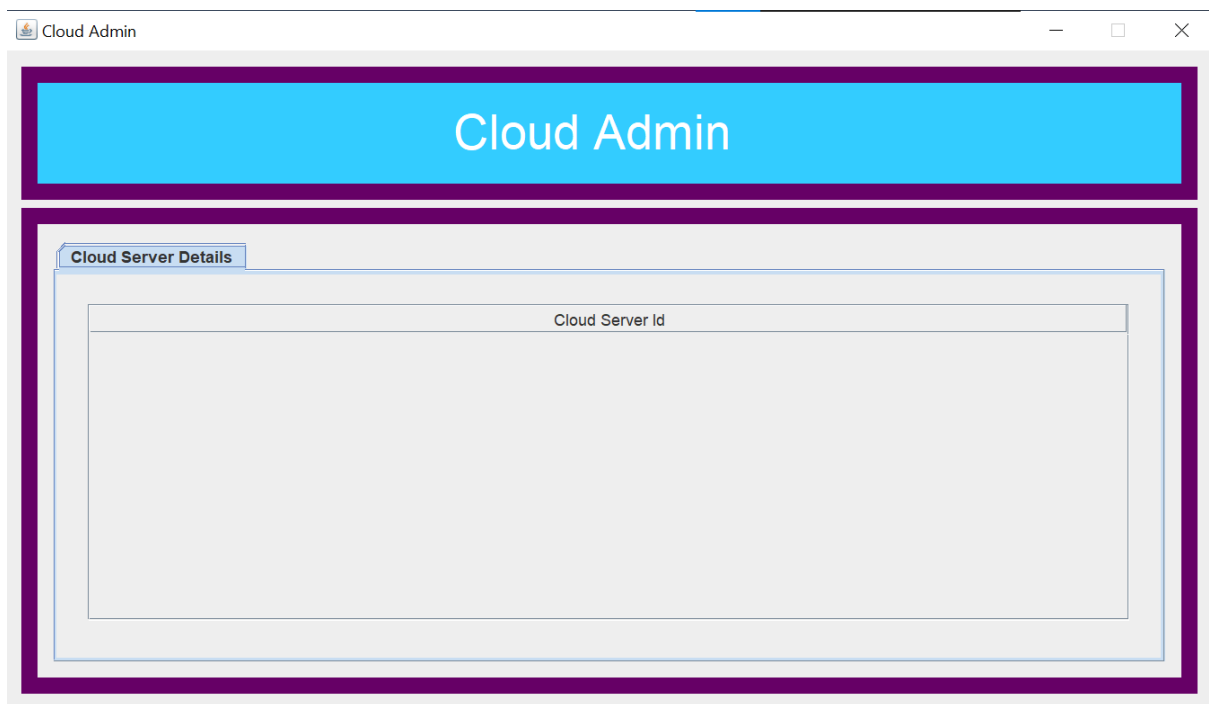


A system has been created that takes data from the sender and forwards it to the recipient. The entering TCP data stream is broken into pieces at the receiver. The chunks are linked in sequence to form a chain and are then stored in a local chunk storage. Each arriving chunk is compared to the chunk storage by the receiver. Once a matched chunk is found on a chain, it retrieves a number of consecutive chunks down the chain as projected chunks in future incoming data.

The recovered chunks' signatures and predicted offsets in the incoming data stream are provided to the sender in a PRED message as a forecast for the sender's upcoming outgoing data. To match data with a prediction, the sender computes SHA-1 across the outgoing data at the predicted offset with the length specified by the prediction, then compares the result to the signature in the prediction. When a signature match is detected, the sender sends a PRED-ACK message to the receiver, instructing it to transfer the matched data from its local storage. This approach eliminates the extra computational and storage expenses paid by TRE in the cloud, which would otherwise negate the bandwidth savings advantages.


6. Implementation: 70%

6.1. Admin Side



6.2. Server Side

Input ×

 Enter the Cloud Server Id:


1Cloud Server - — □ ×

Cloud Server - 1

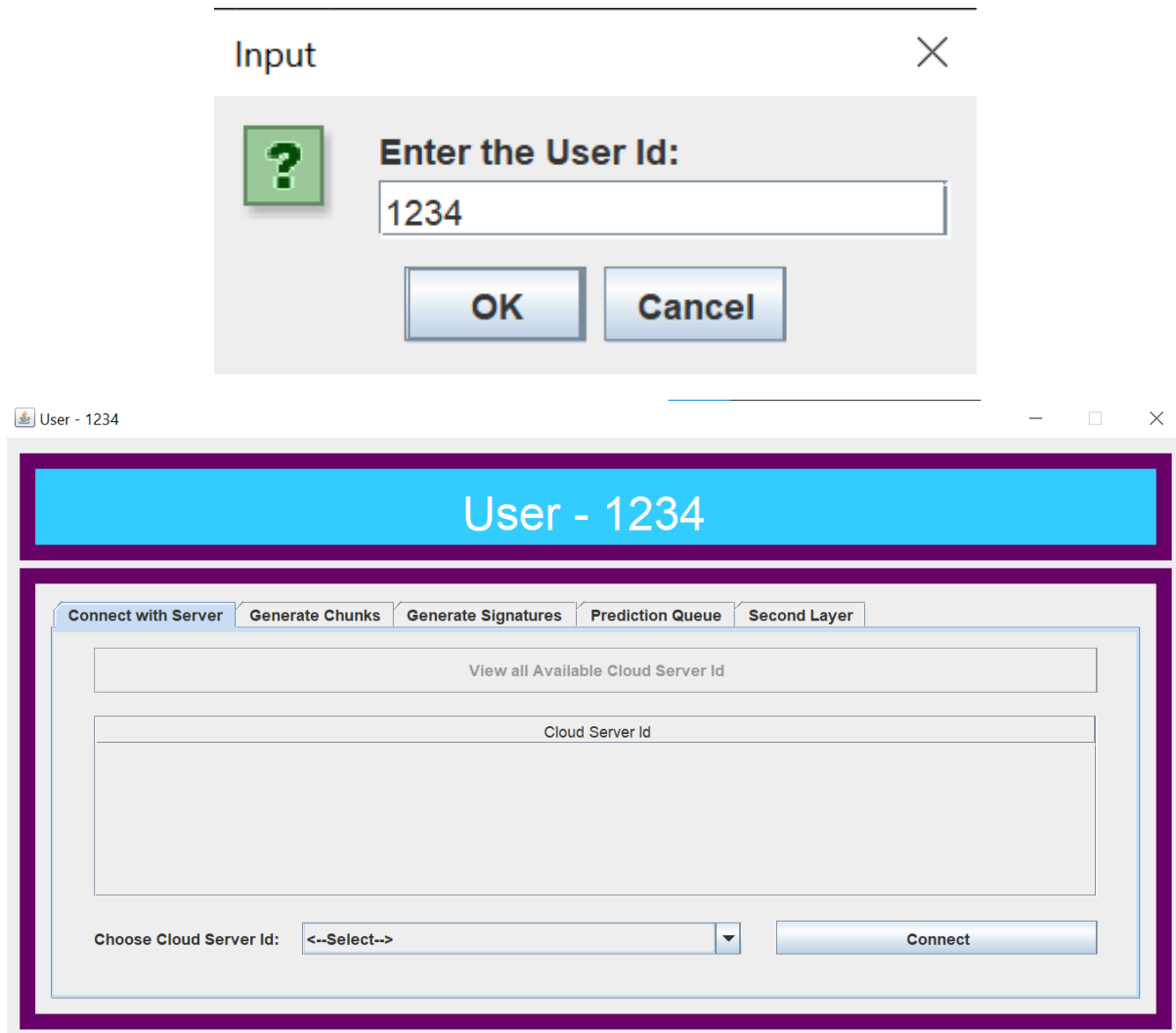
1Cloud Server - — □ ×

Cloud Server - 1

Message ×

 Connected Successfully!

6.3. User Side



```
run:
Port is 5000

Exception in thread "AWT-EventQueue-0" java.lang.NumberFormatException: For input string: "<--Select-->"
    at java.base/java.lang.NumberFormatException.forInputString(NumberFormatException.java:67)
    at java.base/java.lang.Integer.parseInt(Integer.java:668)
    at java.base/java.lang.Integer.parseInt(Integer.java:786)
    at coreuser.UserFrame jButton2ActionPerformed(UserFrame.java:582)
    at coreuser.UserFrame.access$100(UserFrame.java:22)
    at coreuser.UserFrame$2.actionPerformed(UserFrame.java:171)
    at java.desktop/javax.swing.AbstractButton.fireActionPerformed(AbstractButton.java:1972)
    at java.desktop/javax.swing.AbstractButton$Handler.actionPerformed(AbstractButton.java:2313)
    at java.desktop/javax.swing.DefaultButtonModel.fireActionPerformed(DefaultButtonModel.java:405)
    at java.desktop/javax.swing.DefaultButtonModel.setPressed(DefaultButtonModel.java:262)
    at java.desktop/javax.swing.plaf.basic.BasicButtonListener.mouseReleased(BasicButtonListener.java:279)
    at java.desktop/java.awt.Component.processMouseEvent(Component.java:6626)
    at java.desktop/javax.swing.JComponent.processMouseEvent(JComponent.java:3389)
    at java.desktop/java.awt.Component.processEvent(Component.java:6391)
    at java.desktop/java.awt.Container.processEvent(Container.java:2266)
    at java.desktop/java.awt.Component.dispatchEventImpl(Component.java:5001)
    at java.desktop/java.awt.Container.dispatchEventImpl(Container.java:2324)
    at java.desktop/java.awt.Component.dispatchEvent(Component.java:4833)
    at java.desktop/java.awt.LightweightDispatcher.retargetMouseEvent(Container.java:4948)
    at java.desktop/java.awt.LightweightDispatcher.processMouseEvent(Container.java:4575)
    at java.desktop/java.awt.LightweightDispatcher.dispatchEvent(Container.java:4516)
    at java.desktop/java.awt.Container.dispatchEventImpl(Container.java:2310)
    at java.desktop/java.awt.Window.dispatchEventImpl(Window.java:2780)
    at java.desktop/java.awt.Component.dispatchEvent(Component.java:4833)
    at java.desktop/java.awt.EventQueue.dispatchEventImpl(EventQueue.java:773)
    at java.desktop/java.awt.EventQueue$4.run(EventQueue.java:722)
    at java.desktop/java.awt.EventQueue$4.run(EventQueue.java:716)
    at java.base/java.security.AccessController.doPrivileged(AccessController.java:399)
    at java.base/java.security.ProtectionDomain$JavaSecurityAccessImpl.doIntersectionPrivilege(ProtectionDomain.java:86)
```

User - 1234

User - 1234

Connect with Server

Generate Chunks

Generate Signatures

Prediction Queue

Second Layer

Data:

A

A

A

A

A

A

A

A

A

A

A

A

View Selected Data:

Arihant is a good boy

Split Data into Chunks

Arihant
is a g
ood boy

Clear

User - 1234

User - 1234

Connect with Server

Generate Chunks

Generate Signatures

Prediction Queue

Second Layer

Generate Signatures of each Chunks

Chunk - 1: Arihant
Signature: c03c1efe2696aaaf798b86cfc2b2d7d830709fb9

Chunk - 2: is a g
Signature: 3e5d7d1ccadd0a224f8ae5025a092169cc1c0f53

Chunk - 3: ood boy
Signature: c29c731a1968d407872082d8c5f10a1708c993e1

Send Prediction Message with One Chunks Signature to Server

7. Summary

The idea of this project is to develop a system that transmits data from the transmitter to the receiver. At the receiver, the TCP data stream that is entering is divided into pieces. The chunks are placed in a local chunk storage after being joined together in a chain. The receiver compares each arriving piece to the chunk storage. When a matching chunk is located on a chain, the following chunks are retrieved as predicted chunks for incoming data in the future. In a PRED message, the sender is informed of the signatures and anticipated offsets in the incoming data stream for the retrieved chunks as a forecast for the sender's impending outgoing data.

For the same idea we have developed different modules –

1. Admin

The admin will show all the active servers that have been deployed with the id of the servers.

2. Server

Server is basically forming the cloud which will store all the chunks of the data that has been sent by the user.

3. User

User is basically any pay as you go model or the service provider which wants to reduce its cloud bandwidth by reducing the similar chunks of the data by matching the signatures of the chunks created.

Till now we have developed the UIs for all of these modules with the backend created for Admin-Server Connection. We are still working on the backend connection for Server and User.

By the next review, we will be able to send the chunks of message or the data(the data will be in textual form and no other formats like images or videos because that will contain an overhead of extra storage) from the user to the server which will be matched in a 2 layered system.

8. References

- [1] Chhabra N, Bala M. A Comparative study of data deduplication strategies. In 2018 First International Conference on Secure Cyber Computing and Communication (ICSCCC) 2018 Dec 15 (pp. 68-72). IEEE.
- [2] Balasundaram A, Kothandaraman D, Kumar PS, Ashokkumar S. An Approach to Secure Capacity Optimization in Cloud Computing using Cryptographic Hash Function and Data De-duplication. In 2020 3rd International Conference on Intelligent Sustainable Systems (ICISS) 2020 Dec 3 (pp. 1256-1262). IEEE.
- [3] Shin H, Koo D, Shin Y, Hur J. Privacy-preserving and updatable block-level data deduplication in cloud storage services. In 2018 IEEE 11th International Conference on Cloud Computing (CLOUD) 2018 Jul 2 (pp. 392-400). IEEE.
- [4] Pooranian Z, Chen KC, Yu CM, Conti M. RARE: Defeating side channels based on data-deduplication in cloud storage. In IEEE INFOCOM 2018-IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS) 2018 Apr 15 (pp. 444-449). IEEE.
- [5] Leesakul W, Townend P, Xu J. Dynamic data deduplication in cloud storage. In 2014 IEEE 8th International Symposium on Service Oriented System Engineering 2014 Apr 7 (pp. 320-325). IEEE.
- [6] Tian G, Ma H, Xie Y, Liu Z. Randomized deduplication with ownership management and data sharing in cloud storage. *Journal of Information Security and Applications*. 2020 Apr 1;51:102432.
- [7] Kim WB, Lee IY. Survey on data deduplication in cloud storage environments. *Journal of Information Processing Systems*. 2021;17(3):658-73.
- [8] Harnik D, Pinkas B, Shulman-Peleg A. Side channels in cloud services: Deduplication in cloud storage. *IEEE Security & Privacy*. 2010 Dec 3;8(6):40-7
- [9] Fu Y, Xiao N, Jiang H, Hu G, Chen W. Application-aware big data deduplication in cloud environment. *IEEE transactions on cloud computing*. 2017 May 31;7(4):921-34.
- [10] Yuan H, Chen X, Li J, Jiang T, Wang J, Deng RH. Secure cloud data deduplication with efficient re-encryption. *IEEE Transactions on Services Computing*. 2019 Oct 17;15(1):442-56.
- [11] Yan Z, Zhang L, Wenxiu DI, Zheng Q. Heterogeneous data storage management with deduplication in cloud computing. *IEEE Transactions on Big Data*. 2017 May 4;5(3):393-407.

- [12] Ng WK, Wen Y, Zhu H. Private data deduplication protocols in cloud storage. In Proceedings of the 27th Annual ACM Symposium on Applied Computing 2012 Mar 26 (pp. 441-446).
- [13] Reddy BT, Vaishnavi M, Lalitha M, Poojitha P, Kanthi VB. Privacy Preserving Data Deduplication in cloud using Advanced Encryption Standard. In 2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS) 2021 Mar 25 (pp. 1205-1210). IEEE.
- [14] Wen Z, Luo J, Chen H, Meng J, Li X, Li J. A verifiable data deduplication scheme in cloud computing. In 2014 International Conference on Intelligent Networking and Collaborative Systems 2014 Sep 10 (pp. 85-90). IEEE.
- [15] Yingdan S, Huiba L. Data deduplication in cloud computing systems. In 1st International Workshop on Cloud Computing and Information Security 2013 Nov (pp. 483-486). Atlantis Press.
- [16] H. Lou and F. Farnoud, "Data Deduplication With Random Substitutions," in IEEE Transactions on Information Theory, vol. 68, no. 10, pp. 6941-6963, Oct. 2022, doi: 10.1109/TIT.2022.3176778.
- [17] Mr. M. A. R KUMAR, & Mrs. SRILATHA PULI. (2022). FINDING DATA DEDUPLICATION USING CLOUD. ymer, 21(5), 136–142. YMER || ISSN : 0044-0477.
- [18] Lior Aronovich, Ron Asher, Danny Harnik, Michael Hirsch, Shmuel T. Klein, & Yair Toaff. (2016). Similarity based deduplication with small data chunks. Discrete Applied Mathematics, 212, 10–22, ISSN 0166-218X.
- [19] Gagandeep Kaur, Mandeep Singh Devgan, "Data Deduplication Methods: A Review", International Journal of Information Technology and Computer Science(IJITCS), Vol.9, No.10, pp.29-36, 2017. DOI: 10.5815/ijitcs.2017.10.03
- [20] J. Xiong, Y. Zhang, S. Tang, X. Liu and Z. Yao, "Secure Encrypted Data With Authorized Deduplication in Cloud," in IEEE Access, vol. 7, pp. 75090-75104, 2019, doi: 10.1109/ACCESS.2019.2920998.
- [21] Md. Jareena Begum, & B. Haritha. (2020). Data Deduplication Strategies in Cloud Computing. International Journal of Innovative Science and Research Technology, 5, 734–738.
- [22] M. Adithya, Dr. B. Shanthini, 'Security Analysis and Preserving Block-Level Data DEDuplication in Cloud Storage Services by ', Journal of trends in Computer Science and Smart technology (TCSST) (2020) Vol.02/ No. 02 Pages: 120-126

- [23] X. Liu, T. Lu, X. He, X. Yang and S. Niu, "Verifiable Attribute-Based Keyword Search Over Encrypted Cloud Data Supporting Data Deduplication," in *IEEE Access*, vol. 8, pp. 52062-52074, 2020, doi: 10.1109/ACCESS.2020.2980627.
- [24] Q. He et al., "Research on Data Routing Strategy of Deduplication in Cloud Environment," in *IEEE Access*, vol. 10, pp. 9529-9542, 2022, doi: 10.1109/ACCESS.2021.3139757.
- [25] W. Shen, Y. Su and R. Hao, "Lightweight Cloud Storage Auditing With Deduplication Supporting Strong Privacy Protection," in *IEEE Access*, vol. 8, pp.44359-44372, 2020, doi: 10.1109/ACCESS.2020.2977721.
- [26] Stanek, J., Sorniotti, A., Androulaki, E., Kencl, L. (2014). A Secure Data Deduplication Scheme for Cloud Storage. In: Christin, N., Safavi-Naini, R. (eds) *Financial Cryptography and Data Security. FC 2014. Lecture Notes in Computer Science()*, vol 8437. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-662-45472-5_8
- [27] Reddy, B. T., & Rao, M. C. S. (2018). Filter based data deduplication in cloud storage using dynamic perfect hash functions. *International Journal of Simulation Systems, Science & Technology*.
- [28] Priteshkumar Prajapati, Parth Shah, A Review on Secure Data Deduplication: Cloud Storage Security Issue, *Journal of King Saud University - Computer and Information Sciences*, Volume 34, Issue 7, 2022, Pages 3996-4007, ISSN 1319-1578.
- [29] Xuewei Ma, Wenyuan Yang, Yuesheng Zhu, Zhiqiang Bai, 'A Secure and Efficient Data Deduplication Scheme with Dynamic Ownership Management in Cloud Computing'. 2208.09030v3 Wed, 31 Aug 2022 15:47:52 UTC
- [30] P. G., S., R. K., N., Menon, V.G. et al. A secure data deduplication system for integrated cloud-edge networks. *J Cloud Comp* 9, 61 (2020). <https://doi.org/10.1186/s13677-020-00214-6>
- [31] Balu, U., & Resmi, S. (n.d.). Traffic Redundancy and Elimination approach to Reduce cloud Bandwidth and Costs. [Unpublished manuscript].
- [32] Yu, L., Shen, H., Sapra, K., Ye, L., & Cai, Z. (2016). CoRE: Cooperative end-to-end traffic redundancy elimination for reducing cloud bandwidth cost. *IEEE Transactions on Parallel and Distributed Systems*, 28(2), 446-461.