

An Efficient Restoration Mechanism for Big Data Analysis on Cloud Using Compression Techniques

^[1] M. Bairavi, ^[2] M.Chandraleka, ^[3] A.Meenakshi, ^[4] S.RAMPRAKASH

^{[1][2][3]} Final Year Student, ^[4] Assistant Professor

^{[1][2][3][4]} Department of Computer Science and Engineering, University College of Engineering, Thirukkuvilai

Abstract - With the rapidly increasing amounts of data produced worldwide, networked and multi-user storage systems are becoming very popular. However, concerns over data security still prevent many users from migrating data to remote storage. The conventional solution is to encrypt the data before it leaves the owner's premises. While sound from a security perspective, this approach prevents the storage provider from effectively applying storage efficiency functions, such as compression and deduplication, which would allow optimal usage of the resources and consequently lower service cost. Client-side data deduplication in particular ensures that multiple uploads of the same content only consume network bandwidth and storage space of a single upload. Deduplication is actively used by a number of cloud backup providers as well as various cloud services. Unfortunately, encrypted data is pseudorandom and thus cannot be deduplicated: as a consequence, current schemes have to entirely sacrifice either security or storage efficiency. In this paper, we present schemes that permit a more fine-grained trade-off in data chunk similarity. The intuition is that outsourced data may require different levels of protection, depending on how popular it is: content shared by many users. Various deduplication schemes are analyzed and provide experimental results that shows proposed secure data chunk similarity provide improved results in real time cloud environments

Index Terms— Data chunks, Similarity matching, Parallel processing, Data security, Data compression

I. INTRODUCTION

Now a day there is growth in information. With infinite storage space provide by cloud service provider users tend to use as much space as they can and vendors constantly look for techniques aimed to minimize redundant data and maximize space savings. Users will access information according to their needs and most users access same information again and again, the cost of computation, application hosting, content storage and delivery is reduced significantly. The cloud makes it possible for you to access your information from anywhere at any time. Cloud provides benefits such as, flexibility, disaster recovery, software updates automatically, pay-per-use model and cost reduction.[3] While a traditional computer setup requires you to be in the same location as your data storage device, the cloud takes away that step. The cloud removes the need for you to be in the same physical location as the hardware that stores your data. Each provider serves a specific function, giving users more or less control over their cloud depending on the type. Your cloud needs will vary depending on how you intend to use the space and resources associated with the cloud. Cloud computing refers to the use of computers which access Internet locations for computing power, storage and

applications, with no need for the individual access points to maintain any of the infrastructure. Data deduplication is a technique for reducing the amount of storage space an organization needs to save its data. In most organizations, the storage systems contain duplicate copies of many pieces of data. For example, the same file may be saved in several different places by different users, or two or more files that aren't identical may still include much of the same data. Along with low ownership costs and flexibility, users require the protection of their data and confidentiality guarantees through encryption. To make data management scalable deduplication we are use Encryption for secure deduplication services.[1] Unfortunately, deduplication and encryption are two conflicting technologies. While the aim of deduplication is to detect identical data segments and store them only once, the result of encryption is to make two identical data segments indistinguishable after being encrypted. This means that if data are encrypted by users in a standard way as like shared authority, the cloud storage provider cannot apply deduplication since two identical data segments will be different after encryption. On the other hand, if data are not encrypted by users, confidentiality cannot be guaranteed and data are not protected against curious cloud storage providers. There are two types of deduplication in terms of the size: (i) file-

level deduplication, which discovers redundancies between different files and removes these redundancies to reduce capacity demands, and (ii) block level deduplication, which discovers and removes redundancies between data blocks. The file can be divided into smaller fixed-size or variable-size blocks. Using fixed-size blocks simplifies the computations of block boundaries, while using variable-size blocks. A technique which has been proposed to meet these two conflicting requirements is Tag generation and AES Scheme whereby the encryption key is usually the result of the hash of the data segment. Although encryption seems to be a good candidate to achieve confidentiality and deduplication at the same time, it unfortunately suffers from various well-known weaknesses. The confidentiality issue can be handled by encrypting sensitive data before outsourcing to remote servers. Along with low ownership costs and flexibility, users require the protection of their data and confidentiality guarantees through encryption. In this paper, we address the a for mentioned privacy issue to propose a shared authority to the files which Deduplicated based privacy preserving authentication for the cloud data storage, which realizes authentication and authorization without compromising a user's private information. The basic data chunk similarity is shown in fig 1.

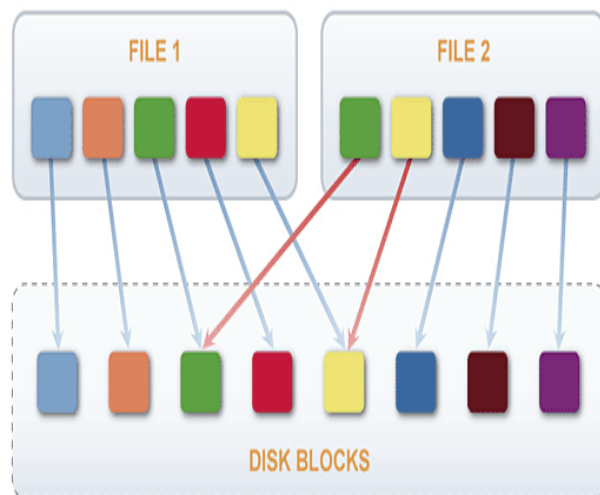


Fig 1: Data chunk similarity

II. RELATED WORK

L. Wang, et al., ... [1] proposed an innovative public cloud usage model for small-to medium scale scientific communities to utilize elastic resources on a public cloud site while maintaining their flexible system controls, i.e., create, activate, suspend, resume, deactivate, and destroy their high-level management entities—service management layers without knowing the details of management. Second, we design and implement an innovative system—DawningCloud, at the core of which are lightweight service management layers running on top of a common management service framework. The common management service framework of DawningCloud not only facilitates building lightweight service management layers for heterogeneous workloads, but also makes their management tasks simple. Third, we evaluate the systems comprehensively using both emulation and real experiments.

B. Li, et al., ... [2] took a step towards bringing the many benefits of the MapReduce model to incremental one-pass analytics. In the new model, the MapReduce system reads input data only once, performs incremental processing as more data is read, and utilizes system resources efficiently to achieve high performance and scalability. The goal is to design a platform to support such scalable, incremental one-pass analytics. This platform can be used to support interactive data analysis, which may involve online aggregation with early approximate answers, and, in the future, stream query processing, which provides near real-time insights as new data arrives. We argue that, in order to support incremental one-pass analytics, a MapReduce system should avoid any blocking operations and also computational and I/O bottlenecks that prevent data from “smoothly” flowing through map and reduce phases on the processing pipeline.

R. Kienzler, et al., ... [3] propose an incremental data access and processing approach for data-intensive cloud applications that can hide data transfer latencies while maintaining linear scalability. Similar in spirit to pipelined query evaluation in traditional database systems, data is accessed and processed in small increments, thereby propagating data chunks from one stage of the data analysis task to another as soon as they

are available instead of waiting until the whole dataset becomes available. This way we can process data mostly in memory (hence, reduce time-consuming I/O to local disk and cloud storage, and avoid storage costs) as well as achieving pipelined parallelism (in addition to the existing partitioned parallelism), leading to a reduction in overall task completion time.

C. Olston, et al., ... [4] proposed a system for Building and updating a search index from a stream of crawled web pages. Some of the numerous steps are deduplication, link analysis for spam and quality classification, joining with click-based popularity measurements, and document inversion. Processing semi-structured data feeds, e.g. news and (micro-)blogs. Steps include de-duplication, geographic location resolution, and named entity recognition. Processing along these lines is increasingly carried out on a new generation of flexible and scalable data management platforms, such as Pig/Hadoop. Hadoop is a scalable, fault-tolerant system for running individual map-reduce processing operations over unstructured data files. Pig adds higher-level, structured abstractions for data and processing. Despite the success of Pig/Hadoop, it is becoming apparent that a new, higher, layer is needed: a workflow manager that deals with a graph of interconnected Pig Latin programs, with data passed among them in a continuous fashion. Given that Pig itself deals with graphs of interconnected data processing steps, it is natural to ask why one would layer another graph abstraction on top of Pig.

K.H. Lee, et al., ... [5] implemented The programming model is inspired by the map and reduces primitives found in Lisp and other functional languages. Before developing the MapReduce framework, Google used hundreds of separate implementations to process and compute large datasets. Most of the computations were relatively simple, but the input data was often very large. Hence the computations needed to be distributed across hundreds of computers in order to finish calculations in a reasonable time. MapReduce is highly efficient and scalable, and thus can be used to process huge datasets. When the MapReduce framework was introduced, Google completely rewrote its web search indexing system to use the new programming model. The indexing system

produces the data structures used by Google web search. The parallelization doesn't necessarily have to be performed over many machines in a network. There are different implementations of MapReduce for parallelizing computing in different environments. Phoenix is an implementation of MapReduce, which is aimed at shared-memory, multi-core and multiprocessor systems, i.e. single computers with many processor cores.

III. DATA DEDUPLICATION TYPES

In this chapter, we can discuss the various data deduplication types as follows

File-level de-duplication

It is commonly known as single-instance storage, file-level data de-duplication compares a file that has to be archived or backup that has already been stored by checking all its attributes against the index. The index is updated and stored only if the file is unique, if not than only a pointer to the existing file that is stored references. Only the single instance of file is saved in the result and relevant copies are replaced by "stub" which points to the original file.

Block-level de-duplication

Block-level data deduplication operates on the basis of sub-file level. As the name implies, that the file is being broken into segments blocks or chunks that will be examined for previously stored information vs redundancy. The popular approach to determine redundant data is by assigning identifier to chunk of data, by using hash algorithm for example it generates a unique ID to that particular block. The particular unique ID will be compared with the central index. In case the ID is already present, then it represents that before only the data is processed and stored before. Due to this only a pointer reference is saved in the location of previously stored data. If the ID is new and does not exist, then that block is unique. After storing the unique chunk the unique ID is updated into the Index. There is change in size of chunk as per the vendor. Some will have fixed block sizes, while some others use variable block sizes

Variable block level de-duplication

It Compares varying sizes of data blocks that can reduce the chances of collision, stated Data links Orlandini. The difference between deduplication schemes are shown in fig 2.

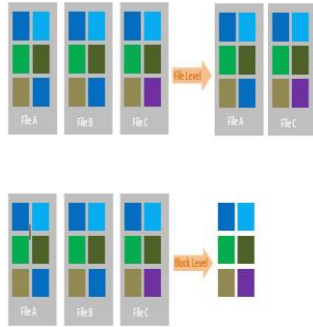


Fig 2: Deduplication schemes

IV. SECURE DEDUPLICATION ALGORITHMS

The main objective of this paper to analyze various encryption algorithms with deduplication schemes. The basic algorithms are shown as follows:

Traditional Encryption algorithm:

Although it is known that data deduplication gives more benefits, security and privacy concerns arise because the users sensitive data is susceptible to both the outsider as well as insider attacks. So, while considering the traditional encryption techniques to secure the users sensitive data there are many issues are associated. Traditional encryption provides data confidentiality but it is not compatible with deduplication. As in traditional encryption different users encrypt their data with their own keys. Thus, the identical data of the different users will lead to different ciphertext which is making the data deduplication almost impossible in this traditional approach. The basic step of the algorithm as shows:

KeyGenSE: k is the key generation algorithm that generates κ using security parameter I

EncSE (k, M): C is the symmetric encryption algorithm that takes the secret κ and message M and then outputs the ciphertext C ;

DecSE (k, C): M is the symmetric decryption algorithm that takes the secret κ and ciphertext C and then outputs the original message M .

Convergent Encryption algorithm: The convergent encryption techniques are those which provide the data confidentiality to the users outsourced data stored on the public clouds. These techniques while providing the confidentiality to the data are also compatible with the data deduplication process. In this algorithm the encryption key is itself derived from the message. So it supports data deduplication also, because the same file will give the same encryption key so it will generate the same ciphertext irrespective of users which makes data deduplication possible.

KeyGenCE(M) $\rightarrow K$ is the key generation algorithm that maps a data copy M to a convergent key K ;

EncCE(K, M) $\rightarrow C$ is the symmetric encryption algorithm that takes both the convergent key K and the data copy M as inputs and then outputs a ciphertext C ;

DecCE(K, C) $\rightarrow M$ is the decryption algorithm that takes both the ciphertext C and the convergent key K as inputs and then outputs the original data copy M ; and

TagGen (M) $\rightarrow T (M)$ is the tag generation algorithm that maps the original data copy M and outputs a tag $T (M)$.

Block cipher algorithm:

In cryptography, a block cipher is a deterministic algorithm operating on fixed-length groups of bits, called blocks, with an unvarying transformation that is specified by a symmetric key. Block ciphers operate as important elementary components in the design of many cryptographic protocols, and are widely used to implement encryption of bulk data. Iterated product ciphers carry out encryption in multiple rounds, each of which uses a different sub key derived from the original key. One widespread implementation of such ciphers is notably implemented in the DES cipher. Many other realizations of block ciphers, such as the AES, are classified as substitution-permutation networks. The publication of the DES cipher was fundamental in the public understanding of modern block cipher design. It also influenced the academic development of

cryptanalytic attacks. Both differential and linear cryptanalysis arose out of studies on the DES design. There is a palette of attack techniques against which a block cipher must be secure, in addition to being robust against brute force attacks. Even a secure block cipher is suitable only for the encryption of a single block under a fixed key. A multitude of modes of operation have been designed to allow their repeated use in a secure way, commonly to achieve the security goals of confidentiality and authenticity. However, block ciphers may also feature as building-blocks in other cryptographic protocols, such as universal hash functions and pseudo-random number generators.

One important type of iterated block cipher known as a substitution-permutation network (SPN) takes a block of the plaintext and the key as inputs, and applies several alternating rounds consisting of a substitution stage followed by a permutation stage—to produce each block of cipher text output. The non-linear substitution stage mixes the key bits with those of the plaintext, creating Shannon's confusion. The linear permutation stage then dissipates redundancies, creating diffusion. A substitution box (S-box) substitutes a small block of input bits with another block of output bits. This substitution must be one-to-one, to ensure invertibility (hence decryption). A secure S-box will have the property that changing one input bit will change about half of the output bits on average, exhibiting what is known as the avalanche effect—i.e. it has the property that each output bit will depend on every input bit.

Variable chunk similarity:

It requires more processing power than the file deduplication, since the number of identifiers that need to be processed increases greatly. Correspondingly, its index for tracking the individual iterations gets also much larger. Using of variable length blocks is even more source-intensive. Moreover, sometimes the same hash number may be generated for two different data fragments, which is called hash collisions. If that happens, the system will not save the new data as it sees that the hash number already exists in the index. The algorithm steps as follows

BlockTag(FileBlock) - It computes hash of the File block as file block Tag;

DupCheckReq(Token) - It requests the Storage Server for Duplicate Check of the file block.

FileUploadReq(FileBlockID, FileBlock, Token) – It uploads the File Data to the Storage Server if the file block is Unique and updates the file block Token stored.

FileBlockEncrypt(Fileblock) - It encrypts the file block with Convergent Encryption, where the convergent key is from SHA Hashing of the file block;

TokenGen(File Block, UserID) – the process loads the associated privilege keys of the user and generate token.

FileBlockStore(FileBlockID, FileBlock, Token) - It stores the FileBlock on Disk and updates the Mapping. The variable chunk similarity level deduplication is shown in fig 3.

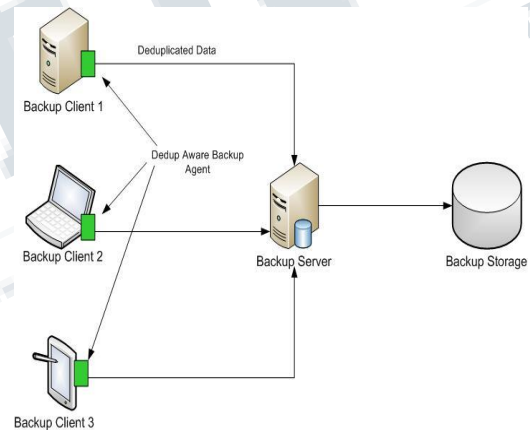


Fig 3: Variable chunk similarity backup server

The mathematical model as follows:

Let S be the system object.

It consist of following $S = \{U, F, CSP\}$

$U =$ no of users $U = \{u_1, u_2, u_3, \dots, u_n\}$

$F =$ no of files $F = \{f_1, f_2, f_3, \dots, f_n\}$

$B =$ no of blocks. $B = \{B_1, B_2, \dots, B_n\}$

$CSP = \{C, PF, V, POW\}$

$C =$ challenge

PF =proof by CSP
 V= verification by TPA
 POW= proof of ownership
 CSP= Cloud Service provider
 CSP={PF,F} PF=proof F=files

The proposed architecture is shown in fig 4.

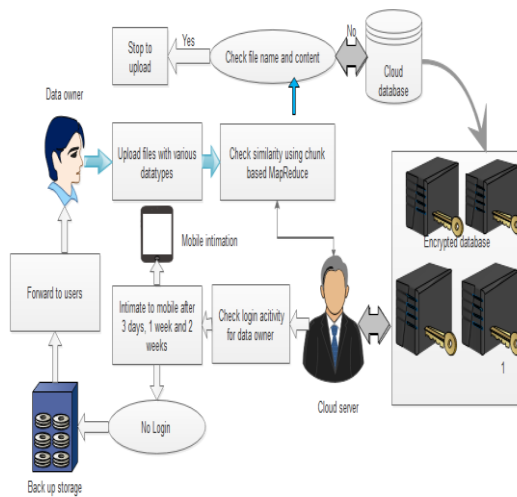


Fig 4: Proposed framework

V. EXPERIMENTAL RESULTS

The proposed algorithm is analyzed in terms of storage preserving and implemented in real time environments. The proposed result is shown in fig 4.

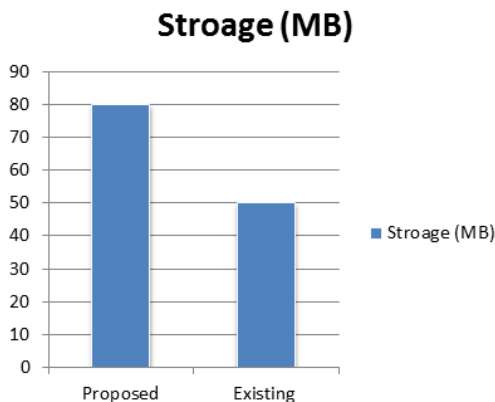


Fig 5: Experimental Results

The proposed system preserves the storage up to 80% storage with security.

VI. CONCLUSION

In this paper analyzed distributed deduplication systems to improve the reliability of data while achieving the confidentiality of the users and also shared authority outsourced data with an encryption mechanism. The constructions were proposed to support file-level and block-level data deduplication. The security of tag consistency and integrity were achieved. We implemented our deduplication systems using the block cipher scheme with variable chunk similarity and demonstrated that it incurs small encoding/decoding overhead compared to the network transmission overhead in regular upload/download operations. Data anonymity is achieved since the wrapped values are exchanged during transmission. User privacy is enhanced by access requests to privately inform the cloud server about the users access desires.

REFERENCES

- 1 L. Wang, J. Zhan, W. Shi and Y. Liang, "In cloud, can scientific communities benefit from the economies of scale?" IEEE Transactions on Parallel and Distributed Systems 23(2): 296-303, 2012.
- 2 B. Li, E. Mazur, Y. Diao, A. McGregor and P. Shenoy, "A platform for scalable one-pass analytics using mapreduce," in: Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD'11), 2011, pp. 985-996.
- 3 R. Kienzler, R. Bruggmann, A. Ranganathan and N. Tatbul, "Stream as you go: The case for incremental data access and processing in the cloud," IEEE ICDE International Workshop on Data Management in the Cloud (DMC'12), 2012
- 4 C. Olston, G. Chiou, L. Chitnis, F. Liu, Y. Han, M. Larsson, A. Neumann, V.B.N. Rao, V. Sankarasubramanian, S. Seth, C. Tian, T. ZiCornell and X. Wang, "Nova: Continuous pig/hadoop workflows," Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD'11), pp. 1081-1090, 2011.

- 5 K.H. Lee, Y.J. Lee, H. Choi, Y.D. Chung and B. Moon, "Parallel data processing with mapreduce: A survey," *ACM SIGMOD Record* 40(4): 11-20, 2012.
- 6 X. Zhang, C. Liu, S. Nepal and J. Chen, "An Efficient Quasiidentifier Index based Approach for Privacy Preservation over Incremental Data Sets on Cloud," *Journal of Computer and System Sciences (JCSS)*, 79(5): 542-555, 2013.
- 7 X. Zhang, T. Yang, C. Liu and J. Chen, "A Scalable Two-Phase Top-Down Specialization Approach for Data Anonymization using Systems, in MapReduce on Cloud," *IEEE Transactions on Parallel and Distributed*, 25(2): 363-373, 2014.
- 8 N. Laptev, K. Zeng and C. Zaniolo, "Very fast estimation for result and accuracy of big data analytics: The EARL system," *Proceedings of the 29th IEEE International Conference on Data Engineering (ICDE)*, pp. 1296-1299, 2013.
- 9 T. Condie, P. Mineiro, N. Polyzotis and M. Weimer, "Machine learning on Big Data," *Proceedings of the 29th IEEE International Conference on Data Engineering (ICDE)*, pp. 1242-1244, 2013.
- 10 Abounaga and S. Babu, "Workload management for Big Data analytics," *Proceedings of the 29th IEEE International Conference on Data Engineering (ICDE)*, pp. 1249, 2013
- 11 M. Abadi, D. Boneh, I. Mironov, A. Raghunathan, and G. Segev, "Message-locked encryption for lock-dependent messages," in *Advances in Cryptology - CRYPTO 2013 - 33rd Annual Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2013. Proceedings, Part I*, 2013, pp. 374-391.
- 12 T. Jiang, X. Chen, Q. Wu, J. Ma, W. Susilo, and W. Lou, "Secure and efficient cloud data Deduplication with randomized tag," *IEEE Trans. Information Forensics and Security*, vol. PP, no. 99,
- 13 Y. Zhou, D. Feng, W. Xia, M. Fu, F. Huang, Y. Zhang, and C. Li, "Secdep: A user-aware efficient fine-grained secure Deduplication scheme with multi-level key management," in *IEEE 31st Symposium on Mass Storage Systems and Technologies, MSST 2015, Santa Clara, CA, USA, May 30 - June 5, 2015*, 2015, pp. 1-14.
- 14 Z. Yan, W. Ding, X. Yu, H. Zhu, and R. H. Deng, "Deduplication on encrypted big data in cloud," *IEEE Trans. Big Data*, vol. 2, no. 2, pp. 138-150, 2016.
- 15 "Prism (surveillance program)," <https://www.theguardian.com/us-news/prism>.
- 16 R. Bhaskar, S. Guha, S. Laxman, and P. Naldurg, "Verito: A practical system for transparency and accountability in virtual economies," in *20th Annual Network and Distributed System Security Symposium, NDSS 2013, San Diego, California, USA, February 24-27, 2013*, 2013.
- 17 D. Boyd, K. Crawford, S. Shaikh, and V. Ravishanker, "Six provocations for big data," *Six provocations-for-Big-Data*.pdf.
- 18 X. Yang, R. Lu, H. Liang, and X. Tang, "SFPM: A secure and fine-grained privacy-preserving matching protocol for mobile social networking," *Big Data Research*, vol. 3, pp. 2-9, 2016.
- 19 A. K. Mohan, E. Cutrell, and B. Parthasarathy, "Instituting credibility, accountability and transparency in local service delivery?: helpline and aasthi in karnataka, india," in *International conference on information and communication technologies and development, ICTD 2013, Cape Town, South Africa, December 7-10, 2013, Volume 1: Papers*, 2013, pp. 238-247.
- 20 D. Boneh, E. Goh, and K. Nissim, "Evaluating 2-dnf formulas on ciphertexts," in *Theory of Cryptography, Second Theory of Cryptography Conference, TCC 2005, Cambridge, MA, USA, February 10-12, 2005, Proceedings*, 2005, pp. 325-341.
- 21 D. Boneh and B. Waters, "Conjunctive, subset, and range queries on encrypted data," in *Theory of Cryptography, 4th Theory of Cryptography Conference, TCC 2007, Amsterdam, The Netherlands, February 21-24, 2007, Proceedings*, 2007, pp. 535-554.
- 22 T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms* (3. ed.). MIT Press, 2009.
-

23 D. E. Knuth, The Art of Computer Programming, Volume III: Sorting and Searching. Addison-Wesley, 1973.

24 G. Ateniese, K. Fu, M. Green, and S. Hohenberger, "Improved proxyre-encryption schemes with applications to secure distributed storage," ACM Trans. Inf. Syst. Secur., vol. 9, no. 1, pp. 1–30, 2006.

25 B. Lynn, "Pbc library," <https://crypto.stanford.edu/pbc/>. "Openssl," <https://www.openssl.org/>. "Dropbox, a file-storage and sharing service," <http://www.dropbox.com>. "Google drive," <http://drive.google.com>.

26 "Mozy: A file-storage and sharing service." <http://mozy.com/>. "Spideroak," <https://www.spideroak.com/>.

27 D. T. Meyer and W. J. Bolosky, "A study of practical deduplication," TOS, vol. 7, no. 4, p. 14, 2012.

28 C. Policroniades and I. Pratt, "Alternatives for detecting redundancy in storage systems data," in Proceedings of the General Track: 2004 USENIX Annual Technical Conference, June 27 - July 2, 2004, Boston Marriott Copley Place, Boston, MA, USA, 2004, pp. 73–86.

29 M. Mulazzani, S. Schrittwieser, M. Leithner, M. Huber, and E. R. Weippl, "Dark clouds on the horizon: Using cloud storage as a attack vector and online slack space," in 20th USENIX Security Symposium, San Francisco, CA, USA, August 8-12, 2011, Proceedings, 2011.

30 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.