



ISSN(Online): 2320-9801  
ISSN (Print): 2320-9798

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 4, Issue 12, December 2016

## Survey on data Deduplication using Similarity and Locality based Approaches

Jayashri Patil, Sunita Barve, Mayura Kulkarni,

M. E Student, MIT Academy of Engineering, Alandi, Pune

Assistant Professor, Department of Computer Engineering, MIT Academy of Engineering, Alandi, Pune, India

Assistant Professor, Department of Computer Engineering, MIT Academy of Engineering, Alandi, Pune, India

**ABSTRACT:** In large scale storage system data deduplication has gained more popularity and attention. Deduplication is one such storage optimization technique that avoids storing duplicate copies of data and only one occurrence of the data is stored on storage media. It is essentially a compression method for removing redundant data. As a space efficient method, data deduplication is used in storage systems for data backup. Storage space is saved by removing redundant data and also in network storage systems the transmission of duplicate data is minimized. Scalability of fingerprint-index search for centralized data deduplication is a main challenge. For high throughput and performance, removing duplicate contents and balancing load by low RAM overhead SiLo scalable deduplication system is used. Similarity and locality exploit both the similarity and locality approaches which are complementary. In SiLo deduplication system, small files which are related are grouped into a segment and segmentation of large files is done. In SiLo RAM usage is reduced for index lookup.

**KEYWORDS:** Deduplication, Similarity, Locality, Storage system.

### I. INTRODUCTION

Only a single instance of data is stored in data deduplication instead of storing multiple copies of the same data. It is a method of removing duplicate copies of data and duplicate copies are replaced with pointers, which point to the identical copy which is stored in storage as a single instance of data. A data set or stream is examined at a sub-file level and only identical data is stored or saved. The workflow of data deduplication consists of input file, hash computation, computing hash with hash index table; whether a match is found or not, if yes set pointer to existing data location and if no save data to memory and its hash to hash index. Duplicate data segments in deduplication technology are detected with the help of fingerprints. Hash functions such as MD5, SHA-1 are used by fingerprints to identify identical segments. Based on granularity deduplication can be categorized as file level deduplication and block level deduplication. In case of file level deduplication the entire file is considered, thus even small append or update makes the file different from the previous version and thus reduces the deduplication ratio. No duplicate file exists at file level deduplication. Whereas in case of block level deduplication data blocks are considered for deduplication. In block level similar data segments of a file will be detected. In offline data deduplication technique, the deduplication process is done after storing the data on the storage disk. In online data deduplication, duplicate data is deleted before writing to the storage disk. Data deduplication can also be categorized as target based deduplication and source based deduplication. In target based deduplication the client does not modify and the client does not perform any deduplication which improves storage utilization and does not save bandwidth. In source based deduplication the client does the deduplication process only on identical data in backup, it saves bandwidth as well as storage space, but there is an extra computational load on the backup client.

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijirce.com](http://www.ijirce.com)

Vol. 4, Issue 12, December 2016

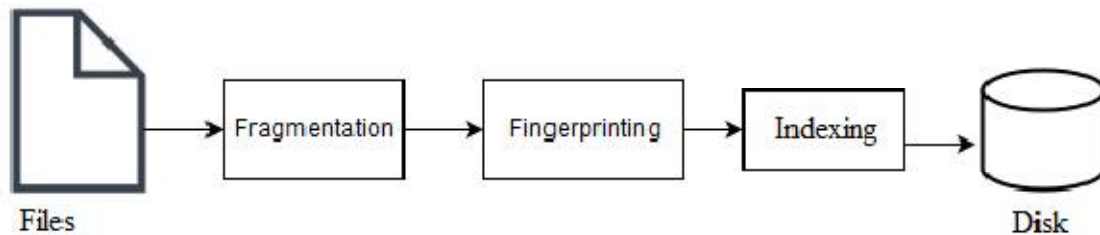


Fig. 1. Deduplication process

As the duplicate data increases the performance decrease. The volume of fingerprints grows with increase in backup data, on disk drive more amount of space is acquired by fingerprint. Thus fingerprint indexing leads to bottleneck performance. Due to this there is a frequent disk access to locate fingerprints which blocks the process of data deduplication. The fingerprints of the same file are stored separately on disk drivers. Whenever the fingerprints are referred there is significant performance degradation. Many approaches have been proposed to address the performance bottleneck. Two primary approaches for data deduplication are similarity-based deduplication and locality-based deduplication.

## A. LOCALITY APPROACH

In locality based approach the order of the backup stream is same for each backup with high probability. Due to which there is increase in the RAM utilization and accesses to on-disk index is reduced, which alleviate disk bottleneck problem. Normally chunk lookups are one by one but some backup streams have high locality between the first, second, and next backups have a very high probability that chunks are in the same order. However this approach shows low speed on backup stream with weak locality.

## B. SIMILARITY APPROACH

Similarity based approaches are designed to overcome the problem encountered by In locality-based approaches Backup streams that either have weak locality or backup stream that lack locality in that case locality approach have problem. Similarity based approach overcome this problem of locality. In backup stream, instead of locality they exploit data similarity from the backup extract similar characteristics and reduce the RAM usage. Instead of lookups per chunks or per local chunks (locality) the lookups are per files. Although is much faster than locality approach it can sacrifice the duplication accuracy. The main idea behind SiLo is that for small files combine into segments to reduce number of fingerprints. For large files divide into segments to increase the similarity detection. Group similar segments order into blocks (preserve locality).

## II. LITERATURE SURVEY

Bo Mao, Hong Jiang Suzhen Wu, Lei Tian (2014) have proposed Performance Oriented I/O Deduplication approach [1]. If data deduplication is directly applied on primary storage then it will cause two problems, fragmentation of data on disks and space contention in memory. Due to this they proposed Performance-Oriented I/O Deduplication. Two approach namely have been considered in POD namely selective dedupe and iCache. Selective dedupe is consider to remove data fragmentation problem and iCache is consider for memory management. POD support features like capacity saving, performance enhancement, small writes elimination, large writes elimination and cache partitioning strategy. POD achieves comparable or better capacity saving than idedupe. I/O performance of primary storage is improved by POD.

Mazhar Ali, Kashif Bilal, Samee U. Khan, Bharadwaj Veeravalli, Keqin Li, Albert Y. Zomaya, (2015) have proposed T-coloring [2]. They have consider Security and performance. In this methodology file is divided into fragment and each



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 4, Issue 12, December 2016

fragment is replicated to different storage node. Only single fragment of particular file is stored on each node. Fragment of particular file is stored on different storage node to increase the surface area for attacker, if in case the attacker accesses the one fragment of file he/she is unable to access the another fragment. Because fragments are stored at a centrality distance from each other which is difficult to guess. In DROP methodology the fragment is also replicated to provide data availability, reliability and improve the data retrieval time. DROP also performs a controlled replication to increase the data availability, reliability, and improve data retrieval time. For reconstruction of file, it provides improved retrieval time for accessing particular file fragment, fragments are placed on the node in such a way that they provide the decreased access cost.

Wen Xia, Hong Jiang, Dan Feng, and Lei Tian, (2015) proposed two data reduction approaches duplicate detection and resemblance detection [3]. In resemblance detection detects similar data object. Granularity at byte-level. Scalability is weak. It is Delta compression method based on super-feature. Duplicate detection is deduplication method based on Secure-Fingerprint. It detects duplicate data object. Granularity is at chunk level. Scalability is strong. Approaches are twofold

: Memory overhead and computation overhead. When this segment is loaded into the locality cache, the two pointers are associated to doubly linked list, the doubly linked list is freed when the segment is removed from cache. Computation overhead is removed by confirming the similarity degree of the DupAdj-detected chunks. To detect similar data chunks DARE efficiently exploits existing duplicate-adjacency information, this achieves highest throughput, data reduction approach.

T. Yang, H. Jiang, D. Feng, Z. Niu, K. Zhou and Y. Wan, (2010) proposed DEBAR, a scalable and high performance deduplication storage architecture for backup and archiving [4]. DEBAR improved capacity, throughput and scalability for deduplication. DEBAR is compared with DDFS in this paper. More backup clients are supported by DEBAR as compared to DDFS. Various applications are supported by DEBAR such as geographic information system grid, WAN, data sharing platform for scientific and engineering applications. In DDFS bloom filter is used to reduce disk index access, it improved deduplication but there is poor scalability. For avoiding fingerprint lookup disk bottleneck in data deduplication DEBAR uses sparse index which exploits inherent locality in backup stream. The main advantage of using DEBAR is that it required half memory space for deduplication compared to DDFS. For high throughput DEBAR can simultaneously run multiple backup servers. DDFS performs to data deduplication scheme two dedupe, in dedupe-1 data chunks are collected and dedupe-2 new data chunk is identified.

M. Fu et al (2016), proposed reducing fragmentation for inline deduplication backup storage [5]. Two drawbacks of fragmentation, first is restore performance is decreased and second it results in invalid chunk. For reducing fragmentation problem two algorithms are proposed History-Aware Rewriting algorithm and Cache-Aware Filter. Two container sparse container and out of order container decreases the restore performance. Fragmentation is in sparse and out of order container. To identify and reduce sparse containers HAR is used which exploits historical information and CAF exploits cache knowledge to identify and reduce out of order container.

C. Li, S. Wang, Xiaochunyun, X. Zhou and G. Wu (2014), have proposed MMD [6], multiple disks are used to boost the reading performance, each disk is used independently as logical device. Due to fragmentation in data deduplication system, reading performance is decreased. For this reason MMD storage approach is used which increases read performance and it is different from RAID. Two algorithms are used, algorithms are used to assign the container to disk. MMD performance is higher compared to RAID0.

J. Liu, Y. Chai, C. Yan and X. Wang (2016), propose a new Delayed Container Organization [7], to increase the restore performance in data deduplication system. The construction of containers is delayed after assigning data chunk in nonvolatile memory. DCO has higher restore speed, better optimization based on a large amount of information, space saving is medium. DCO has three advantages Higher UDRs Containers are produced, More data is duplicated, Restore is speed up.



ISSN(Online): 2320-9801  
ISSN (Print): 2320-9798

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 4, Issue 12, December 2016

## III. CONCLUSION

In this paper we studied different deduplication approaches. Similarity and locality approaches overcome the shortcomings of various approaches in this paper. This combined approach reduces RAM usage, keeps duplication accuracy, and it also increases throughput. Silo approach can effectively improve the disk bottleneck with adequate overhead of CPU, memory, and storage when performing fingerprint lookup, thus improving the throughput of data deduplication. There are multiple approaches and methods suggested by different authors for data deduplication in large storage systems. Various methods for data reduction, data compression, data encoding, data deduplication have been examined to improve performance. Restore performance is also increased.

## ACKNOWLEDGMENT

I take immense pleasure in expressing my humble note of gratitude to my project guide Mrs. Mayura Kulkarni and Dr. Sunita Barve, Assistant Professor, Department of Computer Engineering, MIT Academy of Engineering, Alandi, Pune, for their remarkable guidance and useful suggestions, which helped me in completing the paper before the deadline.

## REFERENCES

- [1] B. Mao, H. Jiang, S. Wu and L. Tian, POD: Performance Oriented I/O Deduplication for Primary Storage Systems in the Cloud, (2014) pp. 767-776.
- [2] M. Ali, K. Bilal, S. Khan, B. Veeravalli, K. Li, A. Zomaya, DROPS: Division and Replication of Data in the Cloud for Optimal Performance and Security, IEEE Transactions on Cloud Computing, vol. PP, no. 99, pp. 1-1.
- [3] W. Xia, H. Jiang, D. Feng and L. Tian, DARE: A Deduplication-Aware Resemblance Detection and Elimination Scheme for Data Reduction with Low Overheads, IEEE Transactions on Computers, vol. 65, no. 6, pp. 1692-1705, June 1 2016.
- [4] T. Yang, H. Jiang, D. Feng, Z. Niu, K. Zhou and Y. Wan, DEBAR: A Scalable high-performance de-duplication storage system for backup and archiving, (2010) pp. 1-12.
- [5] M. Fu et al., Reducing Fragmentation for In-line Deduplication Backup Storage via Exploiting Backup History and Cache Knowledge, IEEE Transactions on Parallel and Distributed Systems, vol. 27, no. 3, pp. 855-868, March 1 2016.
- [6] C. Li, S. Wang, Xiaochunyun, X. Zhou and G. Wu, MMD: An Approach to Improve Reading Performance in Deduplication Systems, (2014) pp. 93-97.
- [7] J. Liu, Y. Chai, C. Yan and X. Wang, A Delayed Container Organization Approach to Improve Restore Speed for Deduplication Systems, IEEE Transactions on Parallel and Distributed Systems, vol. 27, no. 9, pp. 2477-2491, Sept. 1 2016.
- [8] Jingwei Li, Jin Li, Dongqing Xie and Zhang Cai, Secure Auditing and Deduplicating Data in Cloud. IEEE Transactions on Computers (2015).
- [9] D. Meyer and W. Bolosky, A study of practical deduplication, Feb. 2011, pp. 229241.
- [10] F. Guo and P. Efstathopoulos, Building a High-Performance Deduplication System (2011).
- [11] K. Srinivasan, T. Bisson, G. Goodson, and K. Voruganti, iDedup: Latency-Aware, Inline Data Deduplication for Primary Storage, (2012).
- [12] W. Xia, H. Jiang, D. Feng, L. Tian, M. Fu, and Z. Wang, PDedupe: Exploiting Parallelism in Data Deduplication System, pp. 338-347 (2012).
- [13] M. Lillibridge, K. Eshghi, and D. Bhagwat, Improving restore speed for backup systems that use inline chunk-based deduplication, Feb. 2013, pp. 183197.
- [14] A. F. Banu and C. Chandrasekar, A survey on deduplication methods, vol. 3, no. 3, pp. 364368, 2012.
- [15] G. Lu, Y. J. Nam, and D. H. Du, BloomStore: Bloom-filter based memory-efficient key-value store for indexing of data deduplication on flash, Apr. 2012, pp. 111.
- [16] J. Li et al, Secure deduplication with efficient and reliable convergent key management, vol. 25, no. 6, pp. 111, 2013.