# SIMILARITY AND LOCATION AWARE SCALABLE DATA CLEANING AND BACKUP SYSTEM IN CLOUD COMPUTING USING ATTRIBUTE BASED ENCRYPTION

## G.Jayapriya[1], K.Ragini[2], J.Vinothini[3]

Student, B.E. Computer Science and Engineering , Anand Institute of Higher Technology, Chennai, India[1,2]

Assistant Professor, CSE, Anand Institute of Higher Technology, Chennai, India[3]

**Abstract:** massive records is widely considered as potentially the following dominant technology in IT enterprise. It gives simplified machine renovation and scalable aid control with garage structures. As a essential era of cloud computing, storage has been a warm studies subject matter in current years. The high overhead of virtualization has been properly addressed through hardware development in CPU enterprise, and by means of software program implementation improvement in hypervisors. however, the high demand on garage picture garage stays a challenging hassle. existing systems have made efforts to lessen garage picture storage consumption via deduplication inside a storage location network system. nonetheless, storage place community can't satisfy the increasing demand of big-scale garage web hosting for cloud computing due to its value quandary. on this challenge, we suggest SILO, a scalable deduplication report device that has been in particular designed for huge- scale garage deployment. Its design presents rapid storage deployment with similarity and locality based totally fingerprint index for records transfer and low garage consumption by way of deduplication on storage photos. It additionally affords a complete set of storage capabilities including on the spot cloning for storage pics, on-call for fetching through a network, and caching with local disks by way of copy-on-examine strategies. Experiments display that SILO functions perform properly and introduce minor overall performance overhead.

**Keywords:** Deduplication ,Storage ,Security ,SILO function.

## I. INTRODUCTION

Storing large amounts of statistics efficient, in phrases of each time and area, is of paramount difficulty within the design of backup and restore structures. users would possibly wish to periodically (e.g., hourly, every day or weekly) backup records that is stored on their computer systems as a precaution in opposition to viable crashes, corruption or unintentional deletion of vital information. It usually occurs that most of the records has no longer modified since the ultimate backup has been performed, and therefore a whole lot of the current records can already be located inside the backup repository, with best minor adjustment. If the information, in the repository, that is much like the cutting-edge backup statistics, may be positioned efficient, then there's no want to keep the records again, as a substitute, simplest the modifications need be recorded.. This technique of storing commonplace statistics once best is called information deduplication. statistics deduplication is much less difficult to attain with disk based storage than with tape backup. The era bridges the fee hole among disk based totally backup and tape primarily based backup, making disk based totally backup low-cost. Disk based backup has several one of a kind blessings over tape backup in phrases of decreasing backup home windows and improving repair reliability and speed. In a backup and repair system with deduplication it is very probable that a brand new input data flow is similar to statistics already inside the repository, but many exclusive styles of modifications are viable.. Given the ability size of the repository which may additionally have hundreds of terabytes of data, figuring out the regions of similarity to the new incoming statistics is a main venture. in addition, the similarity matching need to be completed quickly in an effort to hold excessive backup bandwidth necessities. This challenge concept aims to alleviate the disk bottleneck of fingerprint lookup, reduce the time of fingerprint research, and improve the throughput of data deduplication. which means that deduplication is simplest executed within individual servers due to overhead concerns, which leaves move-node redundancy untouched. hence, facts routing, a technique to concentrate records redundancy within individual nodes, lessen flow-node redundancy and stability load, turns into a key problem inside the cluster deduplication design.. second, for the intra-node state of affairs, it suffers from the disk bite index research bottleneck. instead, entire-record deduplication is easier and gets rid of record-fragmentation concerns, though on the fee of a few in

any other case reclaimable garage. because the disk era trend is closer to stepped forward sequential bandwidth and decreased in step with-byte price with little or no development in random get admission to velocity, it's now not clear that trading away sequentiality for area financial savings makes sense, as a minimum in number one garage. Complicating topics, these files are in opaque unstructured formats with complicated get right of entry to styles. on the same time there are increasingly many small files in an an increasing number of complicated report device tree.

## 1.1 OBJECTIVE

- In computing data deduplication is a technique for eliminating duplicate copies of repeating data.
- Successfull Implementation of the Technique can improve Storage Utilization ,which may in turn lower capital Expenditure by reducing the overall amount of Storage media required to meet Storage capacity needs.

## 1.2 SCOPE

- Data deduplication- often called intelligent compression or single-instances storage is a process that eliminates redundant copies of data and reduces storage overhead.

- Data duplication techniques ensures that only unique instances of data is retained on storage media ,such as disk ,flash or tape

## II. ANALYSIS

## 2.1 SYSTEM ANALYSIS

System Analysis is a combined process dissection the system responsibilities that are based on problem domain characteristics and user requirement.

### 2.1.1 PROBLEM DEFINITION:

In computing, information deduplication is a specialized statistics compression method for doing away with reproduction copies of repeating information. associated and really synonymous phrases are sensible (records) compression and unmarried-instance (records) storage. This approach is used to enhance garage usage and can also be carried out to community statistics transfers to reduce the number of bytes that should be dispatched. in the deduplication process, specific chunks of facts, or byte patterns, are recognized and stored in the course of a process of analysis. as the evaluation keeps, other chunks are compared to the saved copy and every time a healthy takes place, the redundant chew is changed with a small reference that factors to the saved chunk.

### 2.1.2 EXISTING SYSTEM:

For STORAGE snapshot backup, file level semantics are normally not provided. Snapshot operations take place at the virtual device driver level, which means no fine-grained file system metadata can be used to determine the changed data. Backup systems have been developed to use content fingerprints to identify duplicate content. Offline deduplication is used to remove previously written duplicate blocks during idle time. Several techniques have been proposed to speedup searching of duplicate fingerprints. Existing approaches have focused on such inline duplicate detection in which deduplication of an individual block is on the critical write path.

### ALGORITHM USED:

**Whole File Hashing:** In a whole file hashing (WFH) technique, the whole file is directed to a hashing function. The hashing function is always cryptographic hash like MD5 or SHA-1. The cryptographic hash is used to find entire replicate files. This approach is speedy with low computation and low additional metadata overhead. It works very well for complete system backups when total duplicate files are more common. However, the larger granularity of replicate matching stops it from matching two files that only differ by one single byte or bit of data.

### DISADVANTAGES:

- Traditional data backup approaches often result in a large volume of redundant data.
- File name-based deduplication has relatively higher throughput, but it can not find the redundant data inside files.

### 2.1.3 PROPOSED SYSTEM:

deduplication framework, propose machine implement block degree deduplication system and named as similarity and locality primarily based deduplication framework that could be a scalable and quick overhead close to-precise deduplication device, to defeat In the aforementioned shortcomings of existing schemes. the primary concept of T3S is to recall both similarity and locality inside the backup move concurrently. specifically, disclose and make use of extra similarity via grouping strongly correlated small files into a division and segmenting huge files, and leverage locality within the backup circulation with the aid of grouping closest segments into blocks to confine similar and replica information overlooked by the probabilistic similarity detection. via retaining the parallel index and maintaining spatial locality of assist streams in RAM (i.e., hash desk and locality cache), T3S is capable of do away with massive quantities of redundant statistics, dramatically reduce the numbers of accesses to on-disk index, and drastically increase the RAM utilization.

## ALGORITHM USED:

**SILO similarity algorithm**: Files in the backup stream are first chunked, fingerprinted, and packed into segments by grouping strongly correlated small files and segmenting large files in the File Agent. For an input segment Snew,

Psuedocode for SiLo:

Step 1: Check to see if Snew is in the SHTable. If it hits in SHTable, SiLo checks if the block Bbk containing Snew's similar segment is in the cache. If it is not in the cache, SiLo will load Bbk from the disk to the Read Cache according to the referenced block ID of Snew's similar segment, where a block is replaced in the FIFO order if the cache is full.

Step 2: The duplicate chunks in Snew are detected and eliminated by checking the fingerprint sets of Snew with LHTable (fingerprints index) of Bbk in the cache.

Step 3: If Snew misses in SHTable, it is then checked against recently accessed blocks in the read cache for potentially similar segment (i.e., locality-enhanced similarity detection).

Step 4: Then SiLo will construct input segments into blocks to retain access locality of the input backup stream. For an input block Bnew, SiLo does following:
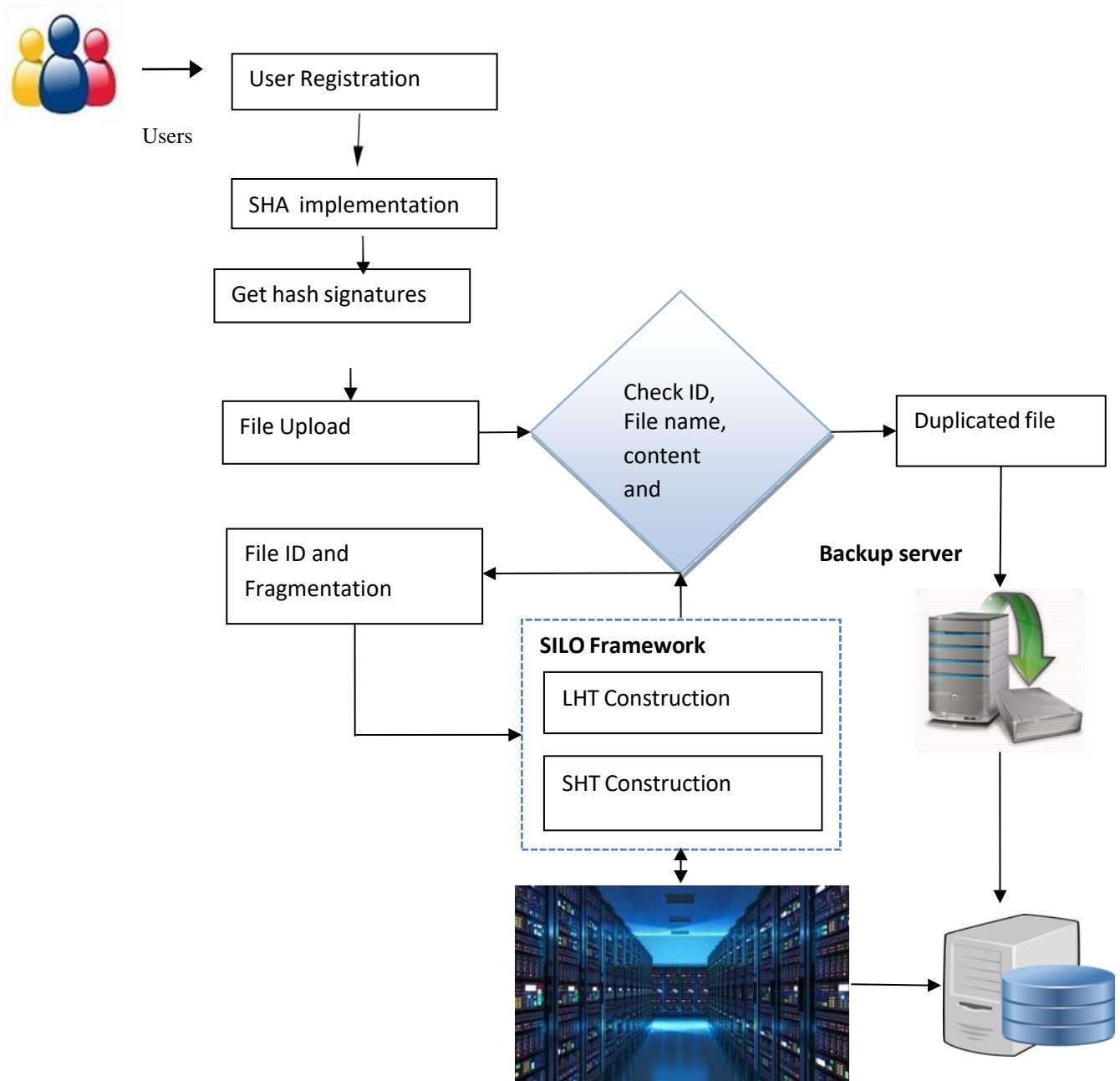
Step 5: The representative fingerprint of Bnew will be examined to determine the stored backup nodes of data block Bnew.

Step 6: SiLo checks if the Write Buffer is full. If the Write Buffer is full, a block there is replaced in the FIFO order by Bnew and then written to the disk.

### ADVANTAGES:

* Able to remove large amounts of redundant data, dramatically reduce the numbers of accesses to on-disk index.
* Maintain a very high deduplication throughput.
* Reduce the time based on index search.

## III.SYSTEM DESIGN

## IV.MODULES

4.1 Cloud resource allocation

4.2 Deduplication scheme

4.3     File system analysis

4.4     Data sharing components

4.5     Evaluation criteria

**4.1 Cloud resource allocation:**

The virtualization is being used to provide ever-increasing number of servers on virtual machines (STORAGEs), reducing the number of physical machines required while preserving isolation between machine instances. This approach better utilizes server resources, allowing many different operating system instances to run on a small number of servers, saving both hardware acquisition costs and operational costs such as energy, management, and cooling. Individual STORAGE instances can be separately managed, allowing them to serve a wide variety of purposes and preserving the level of control that many users want. In this module, clients store data into data servers for future usages. Then data servers stored data in Meta servers.

### 4.2 Deduplication scheme:

Deduplication is a technology that can be used to reduce the amount of storage required for a set of files by identifying duplicate "chunks" of data in a set of files and storing only one copy of each chunk. Subsequent requests to store a chunk that already exists in the chunk store aredone by simply recording the identity of the chunk in the file's block list; by not storing the chunka second time, the system stores less data, thus reducing cost. In this module, we implement fingerprint scheme to identifying chunks differ, both fixed-size and variable-size chunking use cryptographically secure content hashes such as MD5 or SHA1 to identify chunks, thus allowing the system to quickly discover that newly-generated chunks already have stored instances

### 4.3    File system analysis:

In this module, we first broke STORAGE disk images into chunks, and then analyzed different sets of chunks to determine both the amount of deduplication possible and the source of chunk similarity. We use the term disk image to denote the logical abstraction containing all of thedata in a STORAGE, while image files refers to the actual files that make up a disk image. A diskimage is always associated with a single STORAGE; a monolithic disk image consists of a singleimage file, and a spanning disk image has one or more image files, each limited to a particular size. Files are stored in data server with block id and this can be monitored by Data servers. Data servers are mapped by using Meta servers.

### 4.4    Data sharing components:

In this module, we can analyze data sharing components and Meta server in SILO responsible for managing all data servers. It contains SHT and LHT table for indexing each files details for improving search mechanisms. A dedicated background daemon thread will immediately send a heartbeat message to the problematic data server and determines if it is alive.This mechanism ensures that failures are detected and handled at an early stage. The stateless routing algorithm can be implemented since it could detect duplicate data servers even if no one is communicating with them.

### 4.5    Evaluation criteria:

Deduplication is an efficient approach to reduce storage demands in environments with large numbers of STORAGE disk images. As we have shown, deduplication of STORAGE disk images can save 80% or more of the space required to store the operating system and application environment; we explored the impact of many factors on the effectiveness of deduplication. We showed that data localization have little impact on deduplication ratio. However, factors such as the base operating system or even the Linux distribution can have a major impact on deduplication effectiveness. Thus, we recommend that hosting centers suggest "preferred" operating system distributions for their users to ensure maximal space savings. If this preference is followed subsequent user activity will have little impact on deduplication effectiveness.

## III.    RESULTS AND DISCUSSION

Deduplication as an additional post-processing step reduces associated overhead and minimizes impact on write performance. However, in the future even inline deduplication might become applicable to reduce the IO requirements. This might be the case if the computing capacity continuesgrowing while the performance penalty of storing data increases. large chunks: Using large chunks can reduce the overhead of deduplication, but there is a trade-off as larger chunk sizes detect less redundancy. The optimal setting for a deduplicating HPC system is probably a chunk size much larger than the 8 KB, which has been used in this study. small files: Files smaller than 1 MB account for an insignificant amount of the used disk space (though their number is very large). It might be worthwhile to avoid deduplicating these files using such resource-intensive approaches as content- defined chunking. For example, small files could be excluded from the redundancy detection or they could only be deduplicated at a full file level.

## IV.    CONCLUSION

In cloud many data are stored again and again by user. So the user need more spaces storeanother data. That will reduce the memory space of the cloud for the users. To overcome this problem uses the deduplication concept. Data deduplication is a method for sinking the amount ofstorage space an organization wants to save its data. In many associations, the storage systems surround duplicatecopies of many sections of data. For instance, the similar file might be keep in several dissimilar places by dissimilar users, two or extra files that aren't the same may still includemuch of the similar data.

## REFERENCES

[1] D. Meyer and W. Bolosky, "A study of realistic deduplication," in court cases of the ninth USENIX convention on report and storage technologies, 2011.

[2] B. Debnath, S. Sengupta, and J. Li, "Chunkstash: speeding up inline storage deduplication using flash reminiscence," in complaints of the 2010 USENIX convention on USENIX annual technical convention. USENIX association, 2010.

[3] W. Dong, F. Douglis, k. Li, H. Patterson, S. Reddy, and P. Shilane, "Tradeoffs in scalable statistics routing for deduplication clusters," in proceedings of the ninth USENIX convention on file and garage technologies. USENIX affiliation, 2011.

[4] E. Kruus, C. Ungureanu, and C. Dubnicki, "Bimodal content described chunking for backup streams," in lawsuits of the eighth USENIX convention on report and storage technology. USENIX affiliation, 2010.

[5] G.Wallace, F. Douglis, H. Qian, P. Shilane, S. Smaldone, M. Chamness, and W. Hsu, "traits of backup workloads in manufacturing structures," in proceedings of the tenth USENIX conference on file and storage technologies, 2012.

[6] A. Broder, "at the resemblance and containment of files," in Compression and Complexity of Sequences 1997.

[7] D. Bhagwat, k. Eshghi, and P. Mehra, "content-based totally report routing and index partitioning for scalable similarity-based totally searches in a huge corpus," in complaints of the thirteenth ACM SIGKDD global convention on know-how discovery and information mining. ACM, 2007, pp. one zero five–112.

[8] Y. Tan, H. Jiang, D. Feng, L. Tian, Z. Yan, and G. Zhou, "SAM: A Semantic-aware Multi- Tiered supply De-duplication Framework for Cloud Backup," in IEEE thirty ninth global convention on Parallel Processing. IEEE, 2010, pp. 614–623.

[9] M. Lillibridge, okay. Eshghi, D. Bhagwat, V. Deolalikar, G. Trezise, and P. Camble, "Sparse indexing: huge scale, inline deduplication the use of sampling and locality," in Proccedings of the seventh convention on record and garage technologies, 2009, pp. 111–123.

[10] D. Bhagwat, okay. Eshghi, D. long, and M. Lillibridge, "excessive binning: Scalable, parallel deduplication for chew-primarily based document backup," in IEEE international Symposium on Modeling, analysis & Simulation of computer and Telecommunication structures. IEEE, 2009