# Optimized Pricing Scheme in Cloud Environment Using Dedupication

R. Menaga, R. Sangoli, M. Priyadharshini, S. Jayaraj
Computer science and engineering,
UCE-Thirukkuvalai,
Nagapattinam,India,
*Menagadoss96@gmail.com, Sangoli96@gmail.com, priyakannan1396@gmail.com, Jayaraj3486@gmail.com*

**Abstract**: IAAS environment is referred as resources with VM instanSces. Customers can't utilize all resource, but provide full charge for allocated storage.And in server side, storage are not utilized, so scalability become degraded. Implement best billing cycle for access and utilize the resources. Data Deduplication is becoming increasingly popular in storage systems as a space-efficient approach to data backup. Present SiLo, a near-exact deduplication system.That effectively and complementarily exploits similarity and locality to achieve high duplicate elimination. The data secure storing and sharing of the files.

*Keywords: pricing schme, sharing, deduplication, access control.*

_____*****_____

## I. INTRODUCTION

Cloud computing is basically a service oriented architecture rendering easy access to all who make use of it. The need of computation power rendered by the machines is on a continuous hike nowadays. The CPU computation owner is boosting twice for every 3 years. However size of the files keeps increasing also in an amazing rate. 20 years ago, the common format is only text file. In recent year, human is not satisfied in the quality of DVD, and introduce the Blue-ray disk. The file is changed form a few KBs to now nearly 1 TB. The varying characteristics of cloud making it different from other computing technologies are on-demand self-service,agility,autonomic computing, virtualization, parallel and distributed architecture and pay for use model. A Cloud is said to be parallel and distributed architecture as it is consisting of a family of interconnected and virtualized computers.

An organization moving from traditional model to cloud model while been    switched from dedicated hardware and depreciate the service for certain period of time to use shared resources in cloud infrastructure and pay based on the usage. Adjust resources to meet inconsistent and capricious business demand and allow the user to differentiate their business instead of on infrastructure. The cloud providers use pay-as-you-go model. If the cloud model is not used which unexpectedly leads to high charges. First the system is defined as physical machine PM as cloud The cloud computing reduces the ground work of infrastructure costs, run the application faster with upgraded manageability and less maintenance, to rapid servers, and instance or virtual machine VM as the virtual server provided to the users. The resource allocation is cloud computing is much additionalcomplex than in other spread systems like Grid computing environment. In a Grid system

it is improper to share the computer possessions among the multiple applications concurrently running atop it due to the expected mutual performance interference among itself. Whereas, cloud systems usually do not provision corporeal hosts directly to users, but influence virtual resources isolated byVM technology not only can such an elastic storage usage way adapt to user's specific demand, but it can also maximize resource utilization in fine granularity and isolate the anomalous environments for safety purpose. Such a deadline-guaranteed resource distribution with minimized payment is rarely studied in literatures. Moreover, foreseeable errors in predict task workloads will certainly make the problem harder. Based on the elastic storage usage model, aim to design a resource allocation algorithm with high guess error tolerance ability, also minimizing users payments subject to their expected deadlines. Since the idle physical resources can be illogicallydivider and allocated to new tasks, the VM-based divisible resource allocation could be very flexible.

This implies the probability of finding the optimal solution through convex optimization strategy unlike the established Grid model that relies on the indivisible.Resources like the number of physical cores. However, found that it is enviable to directly solve the necessary and sufficient condition to find the optimal solution Karis-Kuhn-Tucker KKT conditions.

In the paper, first develop an optimal resource rental planning model for elastic applications in a cloud environment. In, given known demand patterns over a specific time period, the ASP needs to at regular intervals review the application progress so that no cost is wasted on excessive computation, data transfer and storage. show that proposed planning model is above all useful for high-cost Virtual Machine VM class. This is because cost saving from

the proposed planning model primarily comes from eliminategratuitous data generation cost by decreasing VM rental frequency over the planning horizon. From this perspective, formulation of the deterministic planning problem is consistent with the dynamic lot-sizing problem commonly met in the field of production planning. The second major contribution of this paper is that further address the brave of price uncertainty in optimization model. In fastidious, two possible solution are jointly explored in this paper. first systematically analyse the predictability of spot price in Amazon EC2 and show that the prediction results cannot at hand adequate approximation to be used in the deterministic optimization model.

The solution to this multistage reformulation takes full advantage of the specific structure of the stochastic planning problem.
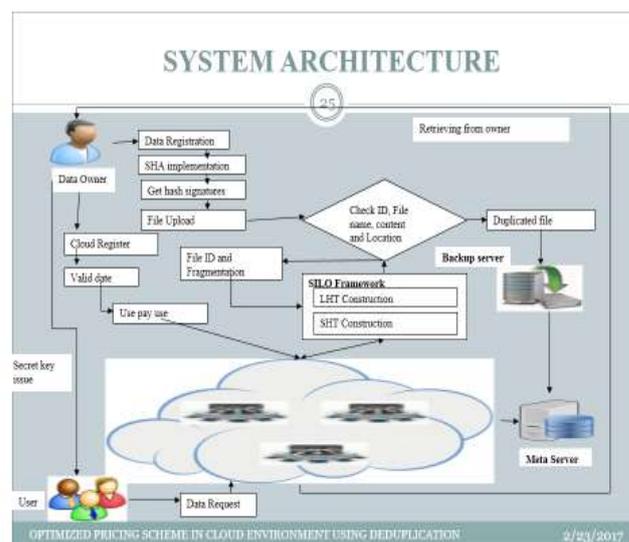
In order to evaluate the tradeoff in space savings connecting whole-file and block-based deduplication, conducted a large-scale study of file system filling on desktop Windows machines at Microsoft. It includes results from a broad cross-section of employees, including software developers, testers, management, sales & marketing, tech-naval support, documentation writers and legal staff. The find that while block-based deduplication of dataset can lower storage blazing up to as little as 32% of its original requirements, nearly three quarters of the development observed could be captured through whole-file deduplication and sparseness. also explore the restriction space for deduplication systems, and quantify the relative been-fits of sparse file support.study of file content is larger and more detailed than any previously available effort, which promises to inform the design of space-efficient storage systems.The customers may want their data encrypted, for reasons ranging from personal privacy to communal policy to legal system. A client could encrypt its file, under a user's key, before storing it. But common encryption modes are randomized, assemblydeduplication impossible since the SS Storage Service effectively always sees different cipher textsanyway of the data. The encryption is deterministic so that the same file will always map to the same cipher text deduplication is possible, but only for that user. Cross-user deduplication, which allows more resources savings, is not possible because encryptions of changed clients, being under different keys, are usually different. It provides high availability and reliability while executing on a substrate of inherently unreliable machinesprincipally through a high degree of duplication of both file content and directory communications.

Since the disk space of desktop computers is mostly unused and apposite less used over time reclaim disk space might not seem to be an important issue.

However, Far sitelike most peer-to-peer systems relies on voluntary participation of the client machines owners.The stage this reclamation in Far site requires solutions to four problems: Enabling the classification and coalescing of identical files when these files are for security reasons encrypted with the keys of different users. Identifying, in a decentralized, scalable, fault understanding manner, files that have identical content. Relocating the replicas of files with identical satisfied to a common set of storage machines. Coalescing the identical files to reclaim storagespace, while maintain the semantics of separate files.

## RELATED WORK:

In the paper, first develop an optimal resource rental planning model for elastic applications in a cloud environment. In, given known demand patterns over a specific time period, the ASP needs to at regular intervals review the application progress so that no cost is wasted on excessive computation, data transfer and storage. first major part is the formulation of a deterministic optimization model that successfully solves the resource rental planning problem. show that proposed planning model is above all useful for high-cost Virtual Machine VM class. This is because cost saving from the proposed planning model primarily comes from eliminategratuitous data generation cost by decreasing VM rental frequency over the planning horizon. From this perspective, formulation of the deterministic planning problem is consistent with the dynamic lot-sizing problem commonly met in the field of production planning. The second major contribution of this paper is that further address the brave of price uncertainty in optimization model. first systematically analyse the predictability of spot price in Amazon EC2 and show that the prediction results cannot at hand adequate approximation to be used in the deterministic optimization model.

Having shown the limitation of prediction, propose a multistage recourse model for stochastic optimization. The model decomposes the stochastic process into sequential pronouncement processes with learning of the random data at various stages. As such, the stochastic optimization problem is transformed into a large-scale linear programming problem solvable by industrial optimization packages. The solution to this multistage reformulation takes full advantage of the specific structure of the stochastic planning problem, and thus provides a near-optimal cost for resource rental under stochastic cost parameters.

In order to evaluate the tradeoff in space savings connecting whole-file and block-based deduplication, conducted a large-scale study of file system filling on desktop Windows machines at Microsoft. It includes results from a broad cross-section of employees, including software developers, testers, management, sales & marketing, tech-naval support, documentation writers and legal staff. The find that while block-based deduplication of dataset can lower storage blazing up to as little as 32% of its original requirements, nearly three quarters of the development observed could be captured through whole-file deduplication and sparseness. For four weeks of full backups, whole file deduplicationwhere a new backup image contains a reference to a duplicate file in an old backup achieves 87% of the savings of block-based. also explore the restriction space for deduplication systems, and quantify the relative been-fits of sparse file support.study of file content is larger and more detailed than any previously available effort, which promises to inform the design of space-efficient storage systems.The customers may want their data encrypted, for reasons ranging from personal privacy to communal policy to legal system. A client could encrypt its file, under a user's key, before storing it. But common encryption modes are randomized, assemblydeduplication impossible since the SS Storage Service effectively always sees different cipher textsanyway of the data. The encryption is deterministic so that the same file will always map to the same ciphertext deduplication is possible, but only for that user. Cross-user deduplication, which allows more resources savings, is not possible because encryptions of changed clients, being under different keys, are usually different. Sharing a single key across a group of users makes the system brittle in the face of client compromise.Far site's intended purpose is to provide the compensation of a central file server a global name space, location transparency, reliability, availability, and security without the attendant disadvantage additional expense, physical plant, administration, and vulnerability to geographically localized faults. It provides high availability and reliability while executing on a substrate of inherently unreliable

machinesprincipally through a high degree of duplication of both file content and directory communications.

Since this intentional and controlled duplication causes a dramatic increase in the space consumed by the file system, it is precious to reclaim storage space due toincidental and erratic duplication of file content. Since the disk space of desktop computers is mostly unused and apposite less used over time reclaim disk space might not seem to be an important issue.

However, Far sitelike most peer-to-peer systems relies on voluntary participation of the client machines owners, who may be reluctant to let their machines participate in a distributed file system that substantially reduces the disk space available for their use. the stage this reclamation in Far site requires solutions to four problems: Enabling the classification and coalescing of identical files when these files are for security reasons encrypted with the keys of different users. Identifying, in a decentralized, scalable, fault understanding manner, files that have identical content. Relocating the replicas of files with identical satisfied to a common set of storage machines. Coalescing the identical files to reclaim storagespace, while maintain the semantics of separate files.

## II.    PROBLRM STATEMENTS:

**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 Sinews in the Settable. If it hits in SHTable, SiLo checks if the block Bbkcontaining Snew's similar segment is in the cache. If it is not in the cache, SiLo will load Bbkfrom 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 Sneware detected and eliminated by checking the fingerprint sets of Snewwith LHTable fingerprints index of Bbkin the cache.

Step 3: If Snewmisses in SHTable, it is then checked against recently accessed blocks in the read cache for potentially similar segment i.e., locality-enhanced similarity detection.
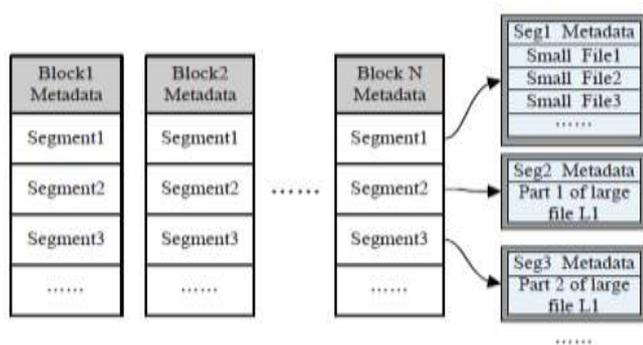
Step 4: Then SiLo will construct input segments into blocks to retain access locality of the input backup stream. For an input block B new, SiLo does following

Step 5: The representative fingerprint of Bnewwill be examined to determine the stored backup nodes of data block Bnew.

123

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 Bnewand then written to the disk.

After the process of deduplication indexing, SiLo will record the chunk-to-file mapping information as the reference for each file, which is managed by the Job Metadata of the Backup Server. For the read operation, SiLo will read the referenced metadata of each target file in the Job Metadata that allows the corresponding data chunks to be read from the data blocks in the Storage Server. These data chunks will then be used to reconstruct the target files in the File Daemon according to the index mapping relationship between files and deduplicated data chunks.

Data structure is shown in below:



### III. SYSTEM ANALYSIS:

In optimal pricing general, there are two serious issues in deploying and provisioning virtual machine (VM) instance over IaaS environment, refined storage allocation and specific pricing for resource renting. Refined resource allocation is usefully implemented VM instance and customizes their storage on demand, which impacts the performance of VMs to complete customers' workload. Precise is also identified as Pay-as-you-go, which involves many types of capital like CPU, memory, and I/O devices. Pricing is a significant component of the cloud computing than it directly affects providers' revenue and customers' budget.

For VM picture backup, file level semantics are normally not provided. Snapshot operation take place at the virtual device level, which means no fine-grained file system metadata can be used to finish the changed data. Backup systems have been residential to use content signature to identify duplicate content. Offline deduplication is used toward remove formerly written duplicate blocks during idle time. Several techniques have been proposed to speedup searching of duplicate signature. Existing approaches have focused on such inline duplicate revealing in which deduplication of anpersonality block is on the critical write path. In existing work, this constraint is complex and there is

no waiting time for many duplicate detection requests. This reduction is improper because in context, difficult to last the backup of required VM images within a logical time window.

In this paper proposed, An appropriate optimal pricing scheme which can make both providers and customer's satisfied services is becoming a major concept in IaaS environment. In Amazon EC2, for example, the minimum pricing time unit of an on-demand instance is one hour. Such a coarse-grained hourly pricing is expected to be reasonably inefficient for short-job users. For instance, users have to pay a for full hour using

cost even their jobs just consumed the storage with a small portion of the one-hour period. Such an incident is called partial usage waste, which appears very often as cloud jobs are quite short in common. Based on the recent characterization of Cloud environment against Grid systems, cloud computing jobs are usually much shorter than Grid jobs. This will induce serious partial usage waste issue. This hourly pricing scheme probably induce idle charge property especially for short jobs, while the fine-grained optimal pricing method not only makes users pay less but also makes provider gain more due to the optimization of unit calculate price time and more users served. Execute a novel optimized fine-grained fair pricing scheme method used by taking into relation the VM maintenance overhead. And find a best-fit bill cycle to reach the maximize social benefit. The designed optimized fine-grained pricing scheme should also satisfy provider, yet providers may suffer higher in the clouds due to finer price rates. In reduplication framework, propose system put into action block level deduplication system and named as similarity and locality based deduplication. This framework that is a scalable and short overhead near-exact reduplication system, to defeat the aforementioned shortcoming of existing schemes. The main implementation of SiLo is to consider both similarity and locality in the backing stream concurrently. Specifically, expose and use more similarity through grouping strongly connected small files into a distribution and segmenting large files, and power region in the backup stream by grouping closest segments into blocks to confine similar and duplicate data miss by the probabilistic comparison detection. By keeping the parallel index and preserve spatial locality of help streams in RAM, SiLo is able to exterminate titanic amounts of disused data, severely reduce the numbers of access to on-disk index, and substantially increase the RAM utilization. This approach divide a large file into many small segment to better representation comparison along with large files while growing the efficiency of the reduplication pipeline.

## IV.    PERFORMANCE EVALUATION:

System testing is a critical aspect of Software Quality Assurance and represents the ultimate review of specification, design and coding. Testing is a process of executing a program with the intent of finding an error. A good test is one that has a probability of finding an as yet undiscovered error. The purpose of testing is to identify and correct bugs in the developed system. Nothing is complete without testing. Testing is the vital to the success of the system. In the code testing the logic of the developed system is tested. For this every module of the program is executed to find an error. To perform specification test, the examination of the specifications stating what the program should do and how it should perform under various conditions.

Unit Testing

Unit testing is a software development process in which the smallest testable parts of an application, called units, are individually and independently scrutinized for proper operation. Unit testing is often automated but it can also be done manually. This testing mode is a component of Extreme Programming XP, a pragmatic method of software development that takes a meticulous approach to building a product by means of continual testing and revision.In this project each and every modules tested separately, whether the user receive the transaction, withdraw and deposit results correctly or not. If the admin update the userdetails correctly modify or not. Each and Every module is check by developer.

Integration Testing:Integration testing is used to verify the combining of the software modules. Integration testing addresses the issues associated with the dual problems of verification and program construction. System testing is used to verify, whether the developed system meets the requirements.The purpose of integration testing is to verify functional, performance, and reliability requirements placed on major design items. These "design items", i.e. assemblages or groups of units, are exercised through their interfaces using black box testing, success and error cases being simulated via appropriate parameter and data inputs. Simulated usage of shared dataareas and inter-process communication is tested and individual subsystems are exercised through their input interface. Tests are constructed to test whether all the components within assemblages interact correctly, for example across procedure calls or process activations, and this is done after testing individual modules, i.e. unit testing. The overall idea is a building block approach, in which verified assemblages are added to a verified base which is then used to support the integration testing of further assemblages. Some different types of integration testing are big bang, top-down, and bottom-up. Other Integration Patterns are: Collaboration Integration, Backbone Integration, Layer Integration, Client/Server Integration, Distributed Services Integration and High-frequency Integration

System testing:

It is a critical aspect of Software Quality Assurance and represents the ultimate review of specification, design and coding. Testing is a process of executing a program with the intent of finding an error. A good test is one that has a probability of finding an as yet undiscovered error. The purpose of testing is to identify and correct bugs in the developed system.The developer check the software whether the program was successfully run in all the operating systems.

Validation Testing:

It's the process of using the new software for the developed system in a live environment i.e., new software inside the organization, in order to find out the errors. The validation phase reveals the failures and the bugs in the developed system. It will be come to know about the practical difficulties the system faces when operated in the true environment.The validation testing was mainly tested in each and every project. For example, in login form, the valid user only allowed to view the website.

## V.    CONCLUSION:

In cloud many data are stored in cloud computing once more and again by user. So the user needs more spaces resourceone more data. That will remove the memory space of the cloud for the users. To defeat this problem uses the deduplication concept. Data deduplication is a technique for sinking the amount of storage space an association wants to save its data. In many associations, the resources systems surround duplicate copies of many sections of data. For instance, the similar file might be keep in several divergent places by dissimilar users, two or extra files that aren't the same may still contain much of the similar data. Deduplicationtake away these extra copies by saving just one copy of the data and return the other copies with pointer that lead repeal to the unique copy. So proposed Block-level deduplication frees up more spaces and demanding category documented as variable block or variable length deduplication has become very accepted. In cloud using the SHT and LHT tables the user easily searches the data and retrieves the searched data from the cloud. And implemented heart beat protocol to recover the data from corrupted cloud server. New metrics are proved that our proposed draw near provide improved results in deduplication process.

_____

In future can extend our work to handle multimedia data for deduplication storage. The multimedia data includes audio, image and videos. And also implement heart beat protocol recover each data server and increase scalability process of system.

### REFERENCES:

[1] Y. Kouki and T. Ledoux, "Rightcapacity: Sla-driven cross-layer cloud elasticity management," IJNGC, vol. 4, no. 3, 2013.

[2] .R. Stefan, K. Holger, M. Pascal, B. Andrew, and L. Miroslaw, "Sizing the cloud," Forrester Research, 2011.

[3] .M. Bellare, S. Keelveedhi, and T. Ristenpart, "Dupless: Server aided encryption for deduplicated storage," in USENIX Security Symposium, 2013.

[4] J. Stanek, A. Sorniotti, E. Androulaki, and L. Kencl, "A secure data deduplication scheme for cloud storage," in Technical Report, 2013.

_____