International Journal of Computer Science Engineering & Technology

APPROVED BY
NATIONAL SCIENCE LIBRARY (NSL)
NATIONAL INSTITUTE OF SCIENCE-COMMUNICATION AND INFORMATION RESOURCES (NISCAIR)
COUNCIL OF SCIENTIFIC AND INDUSTRIAL RESEARCH (CSIR)– NEW DELHI INDIA .

ISSN :2455-9091

# PREPROCESSING EVALUATION FOR DEDUPLICATION FOR BIG DATA IN CLOUD SERVERS

[1] S. Usharani, [2] K. Kungumaraj
[1] Research Scholar, [2] Assistant Professor,
[1] Department of Computer Science, [2] PG Department of Computer Science,
[1] Mother Teresa Women's University, [2] Arulmigu Palaniandavar Arts College for Women,
[1] Kodaikanal, [2] Palani.

_____

**ABSTRACT:** Data de-duplication is normally received in cloud stockpiling administrations to improve capacity usage and decrease transmission transfer speed. It, in any case, clashes with the prerequisite for data confidentiality offered by data encryption. Various leveled approved de-duplication reduces the pressure between data de-duplication and confidentiality and enables a cloud client to perform benefit based copy checks before transferring the data. Existing various leveled approved de-duplication frameworks grant the cloud server to profile cloud clients as per their benefits. In this paper, we propose a protected preprocessing de-duplication framework to help benefit based copy checks and furthermore avoid benefit based client profiling by the cloud server.

**Keywords:** [Preprocessing, De-Duplication, Cloud Servers, Big Data.]
_____

## 1. INTRODUCTION

The missing office that limits a bigger gathering of clouds for legitimate preparing is data organization, as a result of the nonattendance of specific help for data-heightened sensible work forms. At this moment, work process data dealing with in the clouds is practiced using either some application specific overlays that guide the yield of one undertaking to the commitment of another in a pipeline shape, or, even more starting late, using the MapReduce programming model. Such applications require better stockpiling systems that engage VMs than get to shared data all the while. In any case, the present reference business clouds simply give dissent stores, for example, S3 or Azure Blobs got to through high-torpidity

REST (HTTP) interfaces. Besides, conditions may rise where applications may need to change the manner in which data is managed with a particular ultimate objective to acclimate to the certified access system. The necessity for gainful capacity for data-genuine outstanding tasks at hand. A first methodology for regulating data would contain in depending on such open cloud question stores in the manner the application would use a progressively standard parallel record system. Regardless, in the present cloud structures, computational center points are discrete from the capacity centers and correspondence between the two demonstrates a high latency because of the recently referenced data get to traditions. Besides, as these organizations mainly target stockpiling, they simply help

data exchange as a side effect, which infers that they don't engage exchanges between optional VMs without go-between securing the data. Likewise, customers need to pay for securing and moving data in/out of these chronicles despite the expense of leasing the VMs. Cloud providers starting late exhibited the decision of interfacing the cloud stockpiling as virtual volumes to the register center points: Amazon EBS or Azure Drives. Other than being obligated to unclear high latencies from the default stockpiling get to, this decision moreover displays flexibility and sharing containments as only a solitary VM can mount without a moment's delay such a volume. Another alternative to the cloud stockpiling is pass on a parallel record structure on the procedure center points, remembering the true objective to abuse data territory while securing and exchanging work process data. Dispersed capacity game plans, for example, Gfarm were passed on in a register cloud — Eucalyptus, anyway work in the host OS of the physical center point remembering the ultimate objective to store the data in the close-by capacity circles of the machine.

## 2. LITERATURE SURVEY

[1] **Sang-Hyun Lee, Kyung-Wook Shin** describes a design of SHA processor executing three hash algorithms of SHA-512, SHA-512/224 and SHA-512/256. The SHA processor creates summaries of three unique lengths with 512, 224, and 256 bits as indicated by hash algorithms. It was designed that the underlying hash estimations of SHA512/224 and SHA-512/256 were produced utilizing SHA-512 and it depended on 32-bit data path, bringing about a zone effective execution. The SHA processor designed with HDL was checked by FPGA execution. The SHA processor incorporated with 0.18μm CMOS cell library involves 27,368 gate equivalents (GEs) and can work up to 185 MHz clock recurrence. The SHA processor can be utilized for IoT security applications. Catchphrases: Hash, SHA, respectability, IoT security.

[2] **IMTIAZ AHMAD AND A. SHOBA DAS** described a detailed investigation of the engendering of blunders in an equipment usage of SHA-512 is contemplated. This examination included single, transient just as lasting shortcomings infused at all phases of hash esteem calculation. It is discovered that even a solitary blunder infused brought about a large portion of the bits of hash esteem being in mistake and the blunders are spread over the processed hash esteem. In-depth investigation of the individual activities of SHA-512 guided us in proposing appropriate equality forecast plans. They proposed a mistake detection conspire dependent on equality codes and equipment excess and our experiments led on an expansive number of experiments demonstrate that our plan has 100% blame inclusion on account of single blunders. The calculation of execution measurements, for example, region and throughput demonstrate that our plan has just constrained equipment overhead and short delay overhead.

[3] **Alavi Kunhu, Hussain Al-Ahmad and Fatma Taher** proposed algorithm, the vigorous proprietorship watermark data of a patient's mobile number is embedded in hybrid discrete wavelet change and discrete cosine change domain. They introduced a visually impaired computerized watermarking method for the possession assurance and substance confirmation of medical pictures. The delicate SHA256 hash-key data is embedded in hybrid domain is utilized for the substance confirmation of ROI area; in the interim the spatial domain embedded MD5 hash-key watermark is utilized for the verification of RONI district of watermarked medical pictures. It deals with a visually impaired advanced watermarking method for the possession insurance and substance validation of X-ray and MRI medical pictures. Extreme consideration is required before installing watermarking data in medical pictures, to secure the picture quality to

maintain a strategic distance from the wrong finding. The proposed watermarking strategy contains a vigorous watermark for the possession insurance and delicate watermarks for the substance verification. In the watermarking system, the medical picture is divided into areas and the watermark data is embedded in both the change domain and the spatial domain. Advanced watermarking offers new and motivating highlights to the administration frameworks of medical pictures. These medical pictures should be verified and shielded from any controls. Notwithstanding that, the medical picture ought to be related with one of a kind data in regards to the patient, for example, a patient's mobile telephone number. Medical picture watermarking requires extreme consideration while implanting watermarking data inside the medical pictures on the grounds that the extra data must not influence the picture quality as this may cause misdiagnosis. The proposed algorithm has been effectively tried on a number of X-ray and MRI medical pictures. The delicate watermarking is touchy to any slight alteration to the medical pictures, while the proprietorship watermark data is hearty against JPEG, channel and pivot assaults.

[4] **Mochamad Vicky Ghani Aziz, Rifki Wijaya, Ary Setijadi Prihatmanto, Diotra Henriyan** directed on the expansion of the capacity signature, with the reason for data uprightness. Algorithms MD5 execute on 8-bit microcontroller-board, in light of Arduino Uno pack. At that point the next development phase of this examination is, the data sent from a remote terminal units can be scrambled by the microcontroller, utilizing either the algorithms AES, DES, 3DES or TwoFish encryption algorithm, among others. The aftereffect of this underlying examination, MD5 hash algorithm can be actualized in a 8-bit microcontroller with 100% exactness. Be that as it may, it has a few impediments on the issue among them, the data can be prepared to a maximum of 15 (fifteen) characters, data input utilizing keypad matrix 4x3, MD5 hash yield is shown on the LCD illustrations 128x64 and can just enter data input capital letters as it were. The utilization of a microcontroller is regularly executed in world assembling industry or structure some portion of an electronic item, for example, a remote terminal unit on the framework. Remote terminal unit here is an electronic device that is in charge of recovering data utilizing sensors are then sent to a server through link organize or a remote system. The Sensor can be utilized to peruse the data of temperature, altitude, wind speed, etc as required. With respect to the system that is utilized as a mechanism for data delivery might be through link fiber-optic, modem GPRS, VSAT satellite RF stations, etc. The issue emerges when data that was sent was a plaintext that has not been verified (scrambled). Such Data can be perused by unapproved parties in different ways.

[5] **Eko Sediyono, Kartika Imam Santoso and Suhartono** developed the login security framework utilizing OTP that is encoded with MD5 Hash, and the OTP is sent consequently to the enrolled client phone cell number. The mix of One Time Password (OTP), SMS Gateway, and MD5 Hash encryption algorithm are utilized to develop a more verified login method to get to the electronic Academic Information System. The code to be encoded comprises of Student ID, phone number, and access time. The System needs three minutes for security login with SMS-based OTP. The limitation is narrowing the time for programmers to tap and invade. This delay time is a normal gotten from the study among a few specialist organizations in Indonesia. The code produced from the framework is superior to Pseudo Random Number Generator (PRNG) in that the subsequent code is never the equivalent. The upside of this framework is the utilization of MD5 Hash to encode a lot of Student ID, Phone Number, and Time stamp (date and hour of access). MD5 Hash makes results that never been the equivalent with the recently produced OTP. Contrasted with the OTP produced with Pseudo Random Number

Generator (PNRG) may make similar codes. With this condition, it is unthinkable for programmer to break the code and penetrate to the framework. The time delay for dynamic OTP is set as long as 3 minutes. It is unreasonably short for programmers to perhaps break the code. This setup time is likewise short enough contrasted with different applications, for example Facebook and Google that utilization 20 minutes to hold up the client key in the OTP.

## 3. PROPOSED WORK

To address these issues for overseeing data in the clouds, we propose a design for simultaneousness advanced, PaaS-level cloud stockpiling utilizing virtual plates, called Novel preprocessing framework for deduplication(NPFD). For an application comprising of an expansive arrangement of VMs, it federates the neighborhood circles of those VMs into an all inclusive shared data store. Henceforth, applications straightforwardly utilize the nearby circle of the VM occurrence to share input records and spare the yield documents or middle of the road data. As demonstrated by the outcomes introduced in this part, this methodology builds the throughput multiple times over remote cloud stockpiling. Additionally, the advantages of the Novel preprocessing framework for deduplication(NPFD) approach were approved in the context of MapReduce, by structure an Azure model which executes this calculation worldview and utilizes the proposed stockpiling approach as data the board back-end.
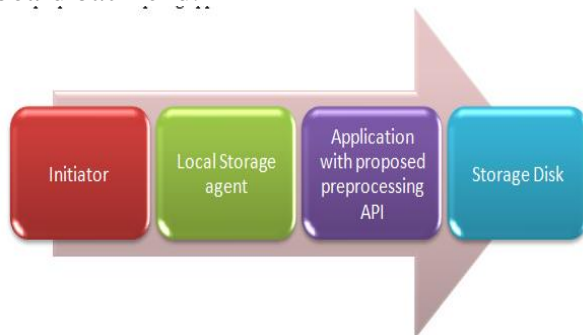


**Figure 1- Novel preprocessing framework for deduplication**

Novel preprocessing framework for deduplication(NPFD): Associating Virtual Disks for a Communication Efficient Storage This fragment introduces the Novel preprocessing framework for deduplication(NPFD) approach for Asscociating the virtual plates of VMs. The framework watches out for the guideline necessities of data genuine applications, point by point in Section 4.2, by giving low-latency data stockpiling upgraded for synchronization. We start from the observation that the plates secretly associated with the VMs, with capacity points of confinement of a few GBs available at no extra expense, are not mishandled to their maximum limit in various cloud associations. Therefore, we propose to add up to pieces of the extra room from the virtual plates in a shared essential pool that is managed in a spread shape. This pool is used to store application-level data. With a particular ultimate objective to change the stack and as such to engage versatility, data is secured in a striped way, for example split into little lumps that are similarly scattered among the area plates of the capacity. Each knot is recreated on different neighborhood plates remembering the true objective to endure dissatisfactions. With this methodology, read and make get to performance under concurrence is colossally upgraded, as the worldwide I/O remaining task at hand is similarly scattered among the close-by circles. Also, this arrangement diminishes latencies by engaging data zone and has a potential for high flexibility, as a developing number of VMs thus prompts a bigger stockpiling framework.

Novel preprocessing framework for deduplication(NPFD) is arranged as requirements be to the going with game plan of blueprint guidelines. These models were picked to such an extent, that they agree to the for the most part normally nullifying necessities of cloud providers and intelligent prevalent enrolling. Data district. Getting to data from remote zones expands the expense of taking care of data monetarily and process

time astute. In any case, in the present cloud appear, the figuring routinely uses the cloud stockpiling for I/O while the locally and uninhibitedly available virtual circles from each VM stay, as it were, unused. This applies despite for moderate delayed consequences of gigantic scale intelligent dealing with, as MapReduce or general work forms. This lessens the general execution performance. Thus, we will probably utilize this free neighborhood VM extra room by aggregating and supervising it in a passed on form and making it available to the applications. Additionally, the general expense is lessen as we decrease the use of the for the most part payable cloud stockpiling. No change of the cloud middleware. Our methodology centers around the business open clouds. It is therefore mandatory that its structure squares don't require any remarkable or raised advantages. As our data organization approach is passed on inside the VMs, the cloud middleware isn't changed in at any rate. This is a key difference from the past endeavors to gather the capacity physical circles of the register center points. These endeavors forced changes to the cloud establishment, so they just worked with open source cloud units. Therefore, our answer is proper for both open and private clouds. It keeps an eye on standard cloud customers, for example, scientists which don't have the stuff or approval to organize and deal with the cloud middleware toolbox. Free coupling among capacity and applications. The Novel preprocessing framework for deduplication(NPFD) cloud data organization is for the most part centered at (anyway not limited to) extensive scale sensible applications executed like the MapReduce computations. Therefore, we propose a deliberate, stub-based building, which can without a doubt be acclimated to other getting ready ideal models.

# 4. EXPERIMENTAL RESULTS

## I/O Time Per Job

| P1 | P2 | P3 | Proposed |
|---|---|---|---|
| **0.581** | 0.6 | 0.475 | 0.299 |
| **0.699** | 0.71 | 0.5 | 0.35 |
| **0.726** | 0.789 | 0.654 | 0.391 |
| **0.78** | 0.812 | 0.699 | 0.423 |
| **0.81** | 0.856 | 0.765 | 0.501 |

**Table 1: Comparison Table of I/O Time Per Job**

Correlation table of I/O time per work describes three existing strategies (P1,P2,P3) and one proposed technique. Contrasted with existing techniques the proposed strategies esteems are low. Proposed technique esteems begins from 0.299 to 0.501.
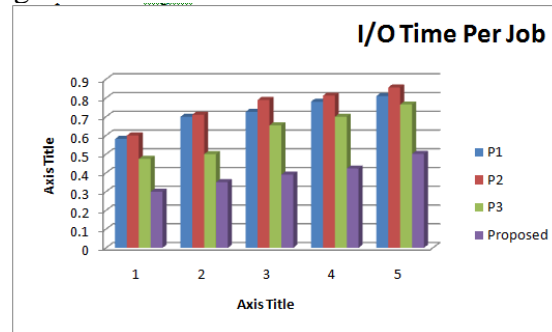


**Figure 2: Comparison Chart of I/O Time Per Job**

The Comparison graph of I/O time per work demonstrates the distinctive benefits of existing techniques and proposed strategy. No of records in x axis and arrangement level in Y axis. At the point when looked at existing strategy and proposed technique the proposed technique esteems are lower than other. Proposed strategy esteems begins from 0.299 to 0.501.

## Computing Time

| P1 | P2 | P3 | Proposed |
|---|---|---|---|
| **0.388** | 0.499 | 0.555 | 0.167 |
| **0.459** | 0.49 | 0.541 | 0.201 |
| **0.501** | 0.569 | 0.615 | 0.28 |
| **0.575** | 0.599 | 0.629 | 0.31 |
| **0.59** | 0.62 | 0.65 | 0.388 |

**Table 2: Comparison Table of Computing Time**

Examination table of processing time describes three existing techniques (P1,P2,P3) and one proposed strategy. Contrasted with existing techniques the proposed strategies esteems are low. Proposed strategy esteems begins from 0.167 to 0.388.
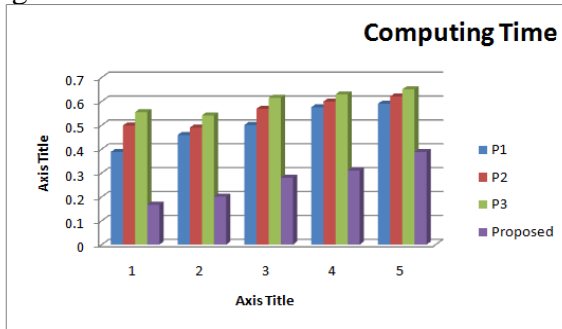


**Figure 3: Comparison Chart of Computing Time**

The Comparison graph of figuring time demonstrates the diverse benefits of existing techniques and proposed strategy. No of records in x axis and grouping level in Y axis. At the point when looked at existing strategy and proposed technique the proposed strategy esteems are lower than other. Proposed technique esteems begins from 0.167 to 0.388.

**Read/Write Processing Throughput**

| P1 | P2 | P3 | Proposed |
|---|---|---|---|
| **0.145** | 0.423 | 0.286 | 0.5 |
| **0.193** | 0.457 | 0.345 | 0.543 |
| **0.222** | 0.494 | 0.377 | 0.588 |
| **0.276** | 0.52 | 0.434 | 0.621 |
| **0.301** | 0.552 | 0.471 | 0.666 |

**Table 3: Comparison Table of Read/Write Processing Throughput**

Correlation table of Read/Write Processing Throughput describes three existing techniques (P1,P2,P3) and one proposed strategy. Contrasted with existing strategies the proposed techniques esteems are high. Proposed technique esteems begins from 0.5 to 0.666.
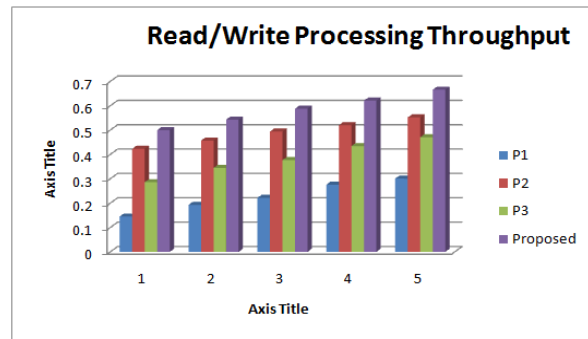


**Figure 4: Comparison Chart of Read/Write Processing Throughput**

The Comparison outline of Read/Write Processing Throughput demonstrates the diverse benefits of existing techniques and proposed strategy. No of records in x axis and grouping level in Y axis. At the point when thought about existing technique and proposed strategy the proposed technique esteems are higher than other. Proposed technique esteems begins from 0.5 to 0.666.

**Transfer Time**

| P1 | P2 | P3 | Proposed |
|---|---|---|---|
| **0.489** | 0.23 | 0.301 | 0.123 |
| **0.522** | 0.259 | 0.388 | 0.188 |
| **0.555** | 0.29 | 0.439 | 0.225 |
| **0.594** | 0.333 | 0.456 | 0.258 |
| **0.635** | 0.358 | 0.492 | 0.282 |

**Table 4: Comparison Table of Transfer Time**

Comparison table of Transfer Time describes three existing methods (P1,P2,P3) and one proposed method. Compared to existing methods the proposed methods values are low. Proposed method values starts from 0.123 to 0.282.
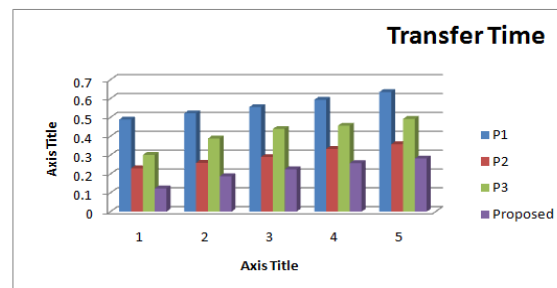


**Figure 5: Comparison Chart of Transfer Time**

The Comparison chart of Transfer Time shows the different values of existing methods and proposed method. No of records in x axis and sequence level in Y axis. When compared existing method and proposed method the proposed method values are lower than other. Proposed method values starts from 0.123 to 0.282.

**Stream Processing**

| P1 | P2 | P3 | Proposed |
|---|---|---|---|
| **0.322** | 0.243 | 0.099 | 0.391 |
| **0.388** | 0.276 | 0.145 | 0.45 |
| **0.432** | 0.289 | 0.172 | 0.489 |
| **0.499** | 0.323 | 0.198 | 0.524 |
| **0.521** | 0.35 | 0.222 | 0.556 |

**Table 5: Comparison Table of Stream Processing**

Comparison table of Stream processing describes three existing methods (P1,P2,P3) and one proposed method. Compared to existing methods the proposed methods values are high. Proposed method values starts from 0.381 to 0.556.
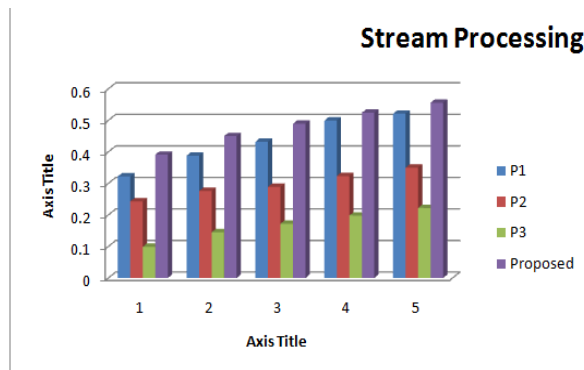


**Figure 6: Comparison Chart of Stream Processing**

The Comparison chart of Transfer Time shows the different values of existing methods and proposed method. No of records in x axis and sequence level in Y axis. When compared existing method and proposed method the proposed method values are higher than other. Proposed method values starts from 0.391 to 0.556.

# CONCLUSION

The thought of approved data de-duplication system is specific data pressure procedure which disposes of excess data just as improves stockpiling and transfer speed use. Merged encryption procedure is proposed to enforce confidentiality amid de-duplication, which scramble data before re-appropriating. Security examination demonstrates that the plans are secure as far as insider and outsider assaults. To all the more likely ensure data security, we present Two Factor Authentication plan of client alongside PoW of documents, to address issue of approved data de-duplication, in which the copy check tokens of records are created by the private cloud server with private keys.

# REFERENCES

[1]. **Sang-Hyun Lee, Kyung-Wook Shin**, "An Efficient Implementation of SHA processor Including Three Hash Algorithms (SHA-512, SHA-512/224, SHA-512/256)", **Published in:** 2018 International Conference on Electronics, Information, and Communication (ICEIC); **Publisher:** IEEE; **Print on Demand(PoD) ISBN:** 978-1-5386-4754-7

[2]. **IMTIAZ AHMAD AND A. SHOBA DAS**, "Analysis and Detection Of Errors In Implementation Of SHA-512 Algorithms On FPGAs", **Published in:** The Computer Journal ( Volume: 50, Issue: 6, Nov. 2007 ); **Publisher:** IEEE; **Date of Publication:** Nov. 2007; Print ISSN: 0010-4620; Electronic ISSN: 1460-2067

[3]. **Alavi Kunhu, Hussain Al-Ahmad and Fatma Taher**, "Medical Images Protection and Authentication using hybrid DWT-DCT and SHA256-MD5 Hash Functions", **Published in:** 2017 24th IEEE International Conference on Electronics, Circuits and Systems (ICECS); **Publisher:** IEEE; Electronic ISBN: 978-1-5386-1911-7; Print on Demand(PoD) ISBN: 978-1-5386-1912-4

[4]. **Mochamad Vicky Ghani Aziz, Rifki Wijaya, Ary Setijadi Prihatmanto, Diotra Henriyan**, "HASH MD5 Function

Implementation at 8-bit Microcontroller", **Published in:** 2013 Joint International Conference on Rural Information & Communication Technology and Electric-Vehicle Technology (rICT & ICeV-T); **Publisher:** IEEE; Electronic ISBN: 978-1-4799-3365-5; Print ISBN: 978-1-4799-3363-1;CD-ROM ISBN: 978-1-4799-336

[5]. **Eko Sediyono, Kartika Imam Santoso and Suhartono**, "Secure Login by Using One-time Password Authentication Based on MD5 Hash Encrypted SMS", **Published in:** 2013 International Conference on Advances in Computing, Communications and Informatics (ICACCI), **Publisher:** IEEE; Electronic ISBN: 978-1-4673-6217-7; Print ISBN: 978-1-4799-2432-5; CD-ROM ISBN: 978-1-4799-2659-6

[6]. **Hongwei Wua, Xiangnan Liua, Weibin Tang**, "A Fast GPU-based Implementation for MD5 Hash Reverse", **Published in:** 2011 IEEE International Conference on Anti-Counterfeiting, Security and Identification, **Publisher:** IEEE; Electronic ISSN: 2163-5056; Print ISSN: 2163-5048

[7]. **Won-Bin Kim, Im-Yeong Lee and Jae-Cheol Ryou**, "Improving dynamic ownership scheme for data Deduplication", **Published in:** 2017 4th International Conference on Computer Applications and Information Processing Technology (CAIPT); **Publisher:** IEEE; Electronic ISBN: 978-1-5386-0600-1; Print on Demand(PoD) ISBN: 978-1-5386-0601-8

[8]. **Anand Bhalerao, Ambika Pawar**, "A Survey: On Data Deduplication for Efficiently Utilizing Cloud Storage for Big Data Backups", **Published in:** 2017 International Conference on Trends in Electronics and Informatics (ICEI); **Publisher:** IEEE; Electronic ISBN: 978-1-5090-4257-9; Print on Demand(PoD) ISBN: 978-1-5090-4258-6

[9]. **Yongtao Zhou , Yuhui Deng, Laurence T. Yang , Ru Yang , and Lei Si**, "LDFS: A Low Latency In-line Data Deduplication File System', **Published in:** IEEE Access ( Volume: 6 ); **Date of Publication:** 01 February 2018 ;

**Publisher:** IEEE; **Electronic ISSN:** 2169-3536

[10]. **Naresh Kumar, Shobha Antwal, Ganesh Samarthyam and S.C Jain**, "Genetic Optimized Data Deduplication for Distributed Big Data Storage Systems", **Published in:** 2017 4th International Conference on Signal Processing, Computing and Control (ISPCC); **Publisher:** IEEE; Electronic ISBN: 978-1-5090-5838-9; Print on Demand(PoD) ISBN: 978-1-5090-5839-6

[11]. **Chia-Mu Yu, Sarada Prasad Gochhayat, Mauro Conti, and Chun-Shien Lu** , "Privacy Aware Data Deduplication for Side Channel in Cloud Storage", **Published in:** IEEE Transactions on Cloud Computing ( Early Access ); **Publisher:** IEEE; **Date of Publication:** 17 January 2018; Electronic ISSN: 2168-7161; CD-ROM ISSN: 2372-0018

[12]. **Zheng Yan, Wenxiu Ding, Xixun Yu, Haiqi Zhu, and Robert H. Deng**, "Deduplication on Encrypted Big Data in Cloud", **Published in:** IEEE Transactions on Big Data; **Date of Publication:** 13 July 2016; **Publisher:** IEEE; Electronic ISSN: 2332-7790; CD-ROM ISSN: 2372-2096

[13]. **Akash Dave , Prof Bhargesh Patel , Prof. Gopi Bhatt , Yash Vora**, " Load balancing in cloud Computing Using Particle Swarm Optimization on Xen Server", **Published in:** 2017 Nirma University International Conference on Engineering (NUiCONE); **Publisher:** IEEE; Electronic ISBN: 978-1-5386-1747-2; Print on Demand(PoD) ISBN: 978-1-5386-1748-9.