

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/358872324>

# Deduplication of IoT Data in Cloud Storage

Chapter · February 2022

DOI: 10.1007/978-981-16-5090-1\_13

CITATIONS

7

READS

256

3 authors:



**Chilukuri Prathima**

Mohan Babu University

8 PUBLICATIONS 22 CITATIONS

[SEE PROFILE](#)



**Naresh Babu Muppalaneni**

National Institute of Technology, Silchar

96 PUBLICATIONS 211 CITATIONS

[SEE PROFILE](#)



**Kabir Kharade**

Shivaji University, Kolhapur

72 PUBLICATIONS 261 CITATIONS

[SEE PROFILE](#)

# Chapter 13

## Deduplication of IoT Data in Cloud Storage



Ch. Prathima, Naresh Babu Muppalaneni, and K. G. Kharade

### Introduction

Distributed computing has lately risen as a very much preferred plan for utility of IoT data. The possibility of cloud is to create processing assets as an utility or an administration on interest to clients over the sensor data. The prospect of distributed computing is kind of the equivalent as matrix registering, which plans to accomplish virtualization of IoT data [1]. In frame computing, the associations sharing their computing assets, similar to processors, in order to accomplish the most processing ability, while distributed computing intends to deliver processing assets as an utility on interest, which may extent to different clients. This makes distributed computing assume a genuine job inside the business space, though grid framework is very much loved in education, science, and engineering [2]. Many meanings of distributed computing are outlined, relied on the individual reason for read or innovation utilized for framework improvement. We trend to plot distributed computing as a plan of action that give processing assets as an administration on interest to clients over the sensor data [3].

Cloud providers pool figuring assets along serve clients by means of a multi-occupant sensor data Fig. 13.1. Registering assets are conveyed over the IoT wherever clients will get to them through various customer stages. Clients will get to the assets on request whenever while no human cooperation with the cloud provider. From a

---

Ch. Prathima (✉)

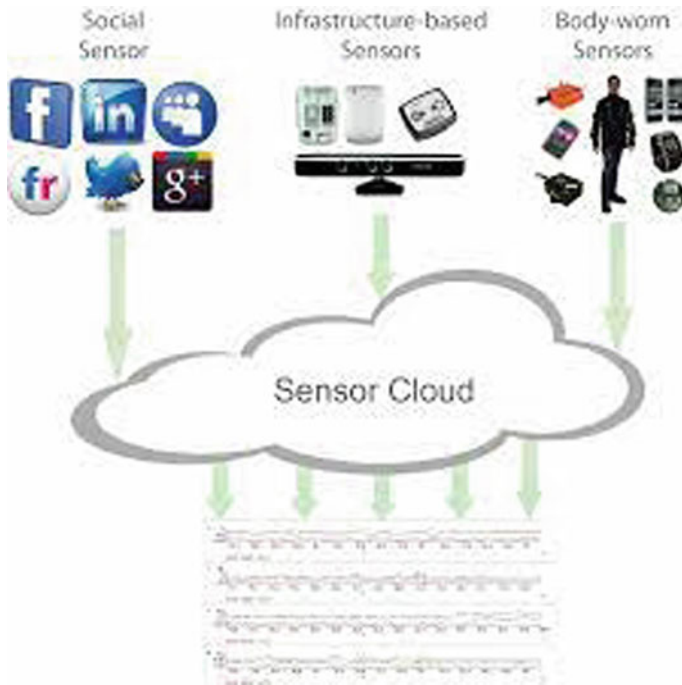
Department of IT, Sree Vidyanikethan Engineering College (Autonomous), Tirupathi, India  
e-mail: [prathima.ch@vidyanikethan.edu](mailto:prathima.ch@vidyanikethan.edu)

N. B. Muppalaneni

Department of CSE, National Institute of Technology Silchar, Silchar, India

K. G. Kharade

Department of Computer Science, Shivaji University, Kolhapur, Maharashtra, India  
e-mail: [kgk\\_csd@unishivaji.ac.in](mailto:kgk_csd@unishivaji.ac.in)



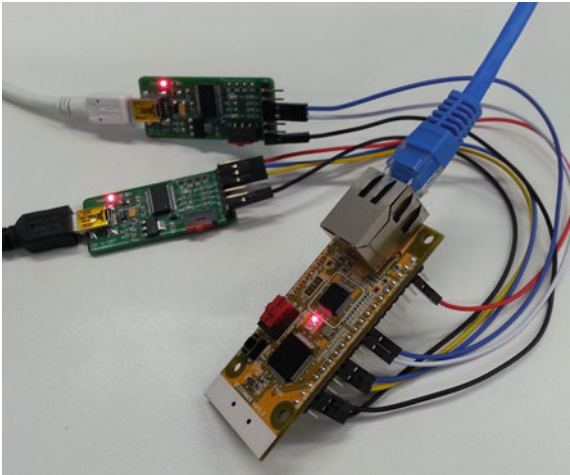
**Fig. 13.1** Enormous data collected from sensors and stored in cloud

client's motivation, registering assets are huge, and customer requests on sensor data are adjusted to fulfill business goals. This is frequently assisted by the adaptability by cloud administrations to scale assets all over on interest contributing the office of virtualization. Also, cloud providers can screen and administrate the use of sensor data for each customer for charge capacities, enhancement assets, and ability to plan and diverse undertakings.

Cloud storage is one among the administrations in distributed computing that gives virtualized capacity on request to clients. Cloud storage might be used in numerous different ways that [4]. For instance, clients will utilize IoT storage as a reinforcement benefit, as against keeping up their very own cache. Associations will move their sensor data to the cloud Fig. 13.2 that they will achieve extra ability at the moderate cost, rather than looking for further physical capacity. Applications running within the cloud conjointly require change or changeless sensor data store in order to help the client applications.

As the quantity of IoT data inside the cloud is progressively expanding, clients hope to succeed in requesting cloud usage whenever, though providers whereas suppliers are needed to take care of system handiness and method an outsized quantity of sensor data. Suppliers want some way to dramatically cut back knowledge volumes, in order that they will cut back prices whereas to spare vitality utilization

**Fig. 13.2** Sensor data transferred to cloud

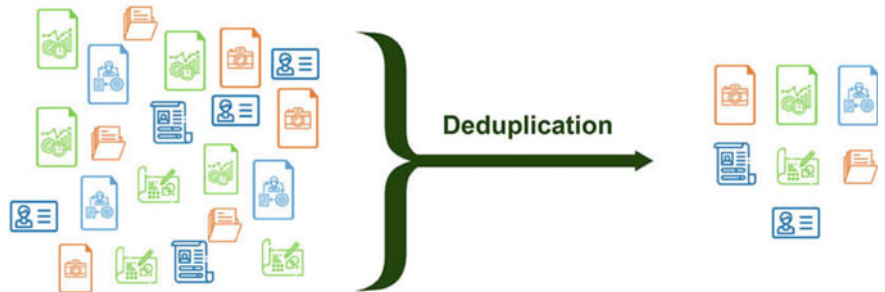


for running huge capacity frameworks. The equivalent as various caches, cache in cloud situations may likewise utilize learning deduplication strategy.

Deduplication might be a procedure whose goal is to help cache power. With the intend to increase sensor data, in ancient deduplication frameworks, copied learning pieces establish and store only one proliferation of the data. Deduplication back each space for storing and network information measure [5]. However, such methods may end up with a negative effect on framework adaptation to internal failure. Thanks to this drawback, several approaches and techniques are planned that not solely give solutions to attain storage potency, however conjointly to boost and help its adaptation to non-critical failure. These procedures give excess of sensor data block after deduplication process.

However, current knowledge deduplications in Fig. 13.3 instruments in distributed storage are static plans connected are applied to all or any knowledge eventualities. For instance, IoT data use in cloud changes in the course of due, some data blocks

Deduplication reduces the amout of stored data



**Fig. 13.3** Deduplication reducing the storage

could likewise be used at a time, anyway probably would not be used in another period of time. On account of setbacks in static plans, that cannot address consistently changing client conduct, deduplication in cloud cache needs a dynamic idea that has the adaptability to adjust to designs and regularly changing client conduct on sensor data in cloud cache.

The contributions of this chapter might be a dynamic learning deduplication subject for distributed storage for IoT data, in order to satisfy a harmony between capacity intensity and adaptation to internal failure necessities, and conjointly to support execution in cloud storage frameworks.

## Background and Related Work

### A. *Deduplication in Cloud storage:*

Deduplication might be a system to downsize space for IoT data. By characteristic excess knowledge, victimization hash esteems to coordinate sensor blocks only one, and making legitimate tips to various duplicates instead of putting away extraordinary genuine duplicates of the sensor information [6, 7]. Deduplication cut backs sensor data volume in this way circle space and system data volume might be decreased that lessen costs. Deduplication might be connected at almost each purpose that sensor data are hold on or sent in cloud storage. Few cloud providers supply no recovery in IoT data [5] and deduplication will not recovery of disaster more functional by duplicating data block once deduplication for increasing mirroring time and data measure savings. Reinforcement and safe caching in clouds may apply facts deduplication in order to reduce physical ability and system movement [8, 9]. Additionally, in movement strategy, we get to send an outsized IoT data of copied picture learning [10]. There are three significant measurements of relocation to consider: knowledge completely transferred, movement time and fix period. Longer movement time and period would be because resources fail. Hence, deduplication will aid movement [11]. To boot, Mandagere et al. express that deduplication calculations replicate the execution of deduplicated sensor data is an issue.

### B. *Problems of Reliability in IoT data*

Acting as deduplication, fewer bits of IoT data are much more important than others. Earlier deduplication approaches do not actualize repetition of sensor data blocks. Subsequently, the sensor blocks should be recreated over the less vital blocks in order to support the architecture. The creators in [12] consider the outcomes of deduplication on the obligation of the cloud-based IoT. They arranged partner degree to deal with boost responsibility by building up a method to mass and experience the significance of every block by inspecting to measure the sensor data blocks that share the same block and utilize this mass to recognize the degree of repetition required for guarantee quality of service.

### C. *Work Done*

Application aware source dedupe [13], causality-based dedupe, and scalable hybrid hash cluster [14]. The majority of existing outcomes that utilize deduplication innovation specialize in spending significant time in the decrease of reinforcement time while overlooking the reclamation time in IoT data. The creators arranged causality-based dedupe, an execution supporter for each cloud reinforcement and cloud reestablish activities that might be a middleware that is symmetrical and might be coordinated into any current reinforcement framework in sensor data. The fundamental point of those associated works is accompanying: surface-to-air missile aims to attain associate degree best exchange off between deduplication strength and deduplication overhead, causality-based dedupe diminishes every reinforcement time and reclamation time. Application aware source dedupe [13] expects to scale back the procedure, increment outturn.

Droplet [15], a shared deduplication cache framework, intended for good outcomes and measurability. It comprises of three segments: one meta server that screens the total sensor data, numerous process servers that run deduplication on computer file stream and different cache junction that store sensor data and deduplicated sensor data chunks.

## Proposed System Model

### Overall Systematic design

Our framework is presently supported customer side deduplication victimization entire record hashing on IoT data. Hashing method is executed at the customer and interfaces with anyone of deduplicators with regards to their time and IoT data. The deduplicator at that point recognizes the duplication by examination with the prevailing hash index in referential server. In early deduplication frameworks, if it is a substitution of hash index, it will be recorded in referential server, and furthermore, the IoT data will be transferred to file servers, and its intelligent way will record in referential server. If it is present, the measure of index for the record will be exaggerated.

Few frameworks might be a scope of duplicates of each record of sensor data with a static number. If the sensor data with an outsized scope of indexing might need addition framework in order to boost handiness. To overcome this issue, some current works brought dimension of repetition into deduplication frameworks. In any case, characteristic dimension of repetition by range of indexing might be a poor measure as a result of files with less indexing could be essential IoT data.

To boost handiness whereas maintaining IoT data, we have a tendency to implement a deduplication framework that thinks about each the changing and using quality of service of the cloud storage. In our framework, when duplication id identified, the redundancy manager at that point figures associate degree best range of duplicates for the record IoT data supported range of references and dimension of quality of

service is important. The ranges of duplicates are modified on the regularly changing sensor data, dimension of quality of service are provided for the sensor data records. The progressions are observed, for instance, when a document is erased by a client, or the degree of quality of service of the record has been refreshed, this may trigger the repetition administrator to re-compute associate the best range of duplicates on the IoT data.

Our arranged frameworks demonstrate is appeared in Fig. 13.4. The frameworks consist of the subsequent parts:

**Balancing of Load:** once hashing strategy with Secure Hash Algorithm 1, buyers send a unique Thumb impression (hash esteem) to a deduplicator by means of the balancing of sensor data. The balancing of load takes demands from purchases causing to anybody of deduplicators in keeping their loads at that instant.

**Referential Deduplicators:** An element intended for characteristic the referral duplication by examination with the prevailing hash esteem hold on in referential server. Distributed storage: An information server to store sensor data and assortment of file servers to store genuine IoT data records and their duplicates.

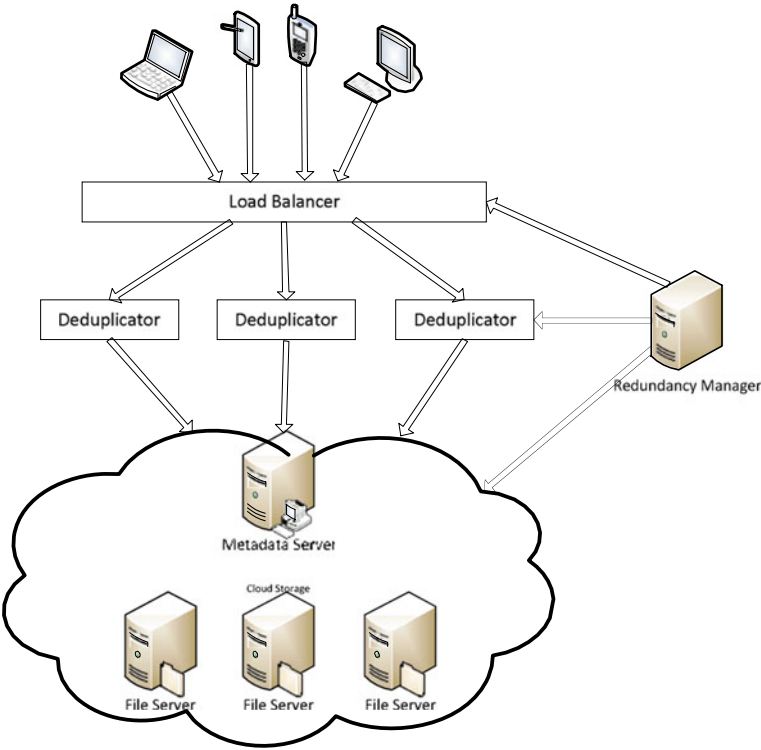


Fig. 13.4 Planned system model

**Replication Manager:** An element to recognize the underlying scope of duplicates of the IoT data and screen the consistently changing dimension of quality of service.

## Results After Experiments Done

Performed tests on the imitations of our planned model. The tests are conducted for 01, 05, and 10 deduplicators.

Every IoT data used in the tests are made with random substance and properties. There are various sizes of sensor records used in this test: one hundred PC memory unit, 150, 200, 250, 300, 500, 800 Kilobytes, 1,2 Megabytes. Uploading, updating and deleting occasions on 10 documents, 100 records, cardinal records, and 10,000 documents of IoT data. For testing the regularly changing dimension of quality of service, each document has been discretionarily distributed its dimension of quality of services (one-five). One quality of service worth of one-five demonstrates the degree of repetition of each sensor document. Records with top-notch of quality of service will be repeated over the lower sensor data.

When isolated deduplicator is utilized, the framework expandability measurability issues take an all-inclusive time once the measure of records exceeds as appeared in Fig. 13.5. This is frequently because of underneath the critical load with extra demands and extra clients, isolated deduplicator cannot keep up the execution of the IoT data framework. When the measure of deduplicators is exceeds to 5 and 10, the outcomes demonstrate that it scale back the range.

For uploading, every sensor data are transferred to the framework initially, and the amount of time taken is observed on 5 folds deduplicator and 10 folds deduplicators. Including extra deduplicators once the uploading IoT documents increment may facilitate to reduce execution time. The test results in Table 13.1. When the quantities of exchange sensor records are 10 and 100 IoT data documents, victimization 5 deduplicators will scale up to 85.76% and 94.25% of the executed taken by isolated deduplicator, though 10 deduplicators will spare longer at 90.86% and 97.57%. When the measure of exchange IoT data documents has been extended to cardinal records, 5 and 10 deduplicators will facilitate to scale back the interval; however, they are small to 91.45% and 95.60% severally. Time efficient impressively lowers once the uploading IoT data records are higher to 10,000 thousands as 5 and 10 deduplicators will reduce 60.15% and 79.76% of interval.

When IoT data records have been transferred to the framework, we tend to perform tests for the case once there's a consistently changing dimension of quality of service, which implies amount of duplicates of sensor documents within the framework likely could be adjusted with regards to range of quality of service. The consequences of updating IoT data records demonstrate that once the measure of sensor data documents increment, adding extra deduplicators will facilitate to scale the interval. When the quantities of records are 10, 100 documents, 1000 and 10,000 documents, victimization 5 deduplicators will reduce 41.77% and 61.78%, 63.79%, and 75.27% of the



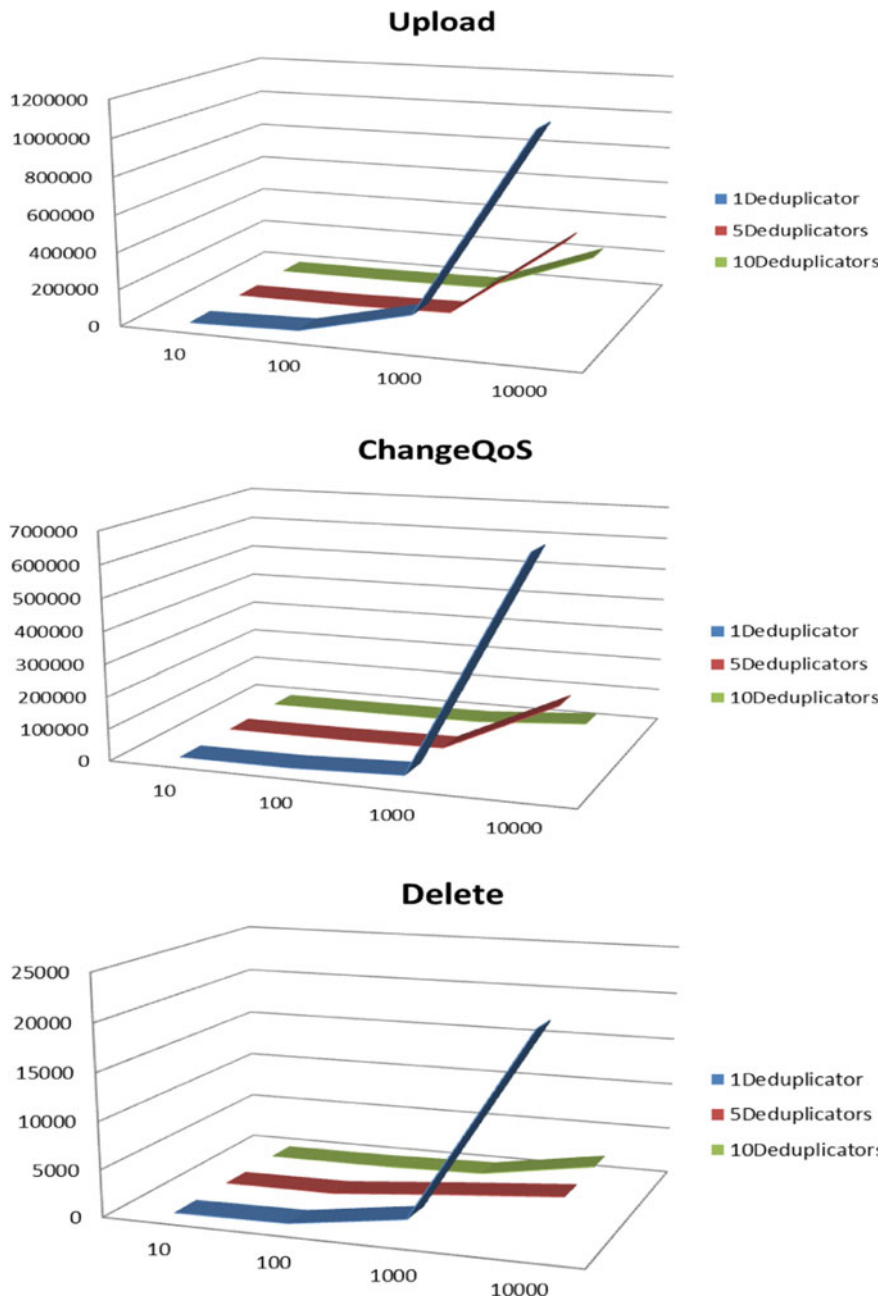


Fig. 13.5 Experimental results

**Table 13.1** Five and ten duplications percentage of time saving

Quantity of files	Uploading of data		Updating of data		Deletion of data	
	5	10	5	10	5	10
10	85.76	90.86	41.77	75.02	93.44	98.68
100	94.26	97.57	61.78	75.34	69.33	90.58
1000	91.45	95.60	63.79	82.09	40.77	85.88
10,000	60.15	79.76	75.27	96.17	90.29	90.05

execution process by Isolated deduplicator, though 10 deduplicators will last longer at 75.02%, 75.34%, 82.09%, and 96.17%. We tend to establish that, when the quantities of IoT data records are 10, 100 and 1000, efficient by including extra deduplicators is nevertheless efficient for the uploading cases. When the quantities of IoT data documents are exaggerated to something like 1000 and 10,000 sensor data records, the efficiency of 5 and 10 deduplicators still increment, in distinction to the uploading cases on the sensor data.

We perform analyses for deleting of IoT data records. Adding extra deduplicators may likewise reduce execution interval; anyway, the consequences of deletion sensor data documents are marginally totally unique in relation to the updating and uploading cases. The outcomes demonstrate that once the quantities of records are 10, 100 documents, 1000 and 10,000 records, exploitation 5 deduplicators will decrease 93.44% and 69.33%, 40.77%, and 90.29% of the intervals taken by isolate deduplicator, though 10 deduplicators will spare longer at 98.68%, 90.58%, 85.88%, and 90.05%. We can see that, for the deleting cases, efficient by including extra deduplicators is small once the quantities of IoT data documents are exaggerated from 10 to 100 and 1000 sensor data records. Nonetheless, when the quantities of IoT data records are exaggerated to 10 thousands, extra deduplicators facilitate to expand efficient.

The test outcomes do not seem to be essentially shocking. Adding additional deduplicators will facilitate to scale back the interval. However, we have a tendency to still establish what are the best range deduplicators to be valuable into the framework with regards to the situations and furthermore the range of IoT data records by then. Also, the outcomes to be assessed against constant data collected through sensors in IoT.

Conclusion

Distributed storage administrations gave distributed computing has been expanding in quality. It offers on interest virtualized capacity assets and clients exclusively get the volume of IoT data that they really wanted. Since the expanding request and IoT data store within the cloud, facts deduplication is one among the procedures to enhance caching potency. In any case, current facts deduplication procedures in

distributed caching are constant topic that restrains their full persistence in effective sensor data in distributed caches.

In this chapter, we tend to propose an effective IoT data deduplication subject for distributed caching, in order to meet a harmony between regularly changing capacity intensity and adaptation to non-critical failure necessities, and conjointly to boost execution in distributed caching frameworks. We tend to dynamically modify the amount of duplicates of sensor data documents with regards to the regularly changing dimension of quality of service. The test outcomes demonstrate that our arranged framework is acting admirably and may deal with measurability drawbacks. We tend to conjointly decide to screen the consistently changing of clients' interest of sensor data documents. Additionally, we tend to choose to assess availability and execution of the framework on the IoT data.

## References

1. Qabil, S., Waheed, U., Awan, S.M., Mansoor, Y., Khan, M.A.: A survey on emerging integration of cloud computing and internet of things. *Int. Conf. Inf. Sci. Commun. Technol. (ICISCT)* **2019**, 1–7 (2019). <https://doi.org/10.1109/CISCT.2019.8777438>
2. Dillon, T., Chen, W., Chang, E.: Cloud computing: issues and challenges. In: 2010 24th IEEE International Conference on Advanced Information Networking and Applications (AINA), pp. 27–33 (2010)
3. Abdelwahab, S., Hamdaoui, B., Guizani, M., Rayes, A.: Enabling smart cloud services through remote sensing: an internet of everything enabler. *Internet of Things J. IEEE* **1**(3), 276–288 (2014)
4. SNIA Cloud Storage Initiative. In: *Implementing, Serving, and Using Cloud Storage. Whitepaper* (2010)
5. Aazam, M., Khan, I., Alsaffar, A.A., Huh, E.-N.: Cloud of things: integrating internet of things and cloud computing and the issues involved. In: *Proceedings of 2014 11th International Bhurban Conference on Applied Sciences and Technology (IBCAST) Islamabad, 14th-18th January, 2014* (2014)
6. Harnik, D., Pinkas, B., Shulman-Peleg, A.: Side channels in cloud services: deduplication in cloud storage. *Secur. Privacy IEEE* **8**, 40–47 (2010)
7. Guo-Zi, S., Yu, D., Dan-Wei, C., Jie, W.: Data backup and recovery based on data deduplication. In: 2010 International Conference on Artificial Intelligence and Computational Intelligence (AICI), pp. 379–382 (2010)
8. Kumar Bose, S., Brock, S., Skeoch, R., Shaikh, N., Rao, S.: Optimizing live migration of virtual machines across wide area networks using integrated replication and scheduling. In: 2011 IEEE International Systems Conference (SysCon), pp. 97–102 (2011)
9. Bose, S.K., Brock, S., Skeoch, R., Rao, S.: CloudSpider: combining replication with scheduling for optimizing live migration of virtual machines across wide area networks. In: 2011 11th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid), pp. 13–22 (2011)
10. Mandagere, N., Zhou, P., Smith, M.A., Uttamchandani, S.: Demystifying data deduplication. In: *Proceedings of the ACM/IFIP/USENIX Middleware '08 Conference Companion, Leuven, Belgium* (2008)
11. Bhagwat, D., Pollack, K., Long, D.D.E., Schwarz, T., Miller, E.L., Paris, J.F.: Providing high reliability in a minimum redundancy archival storage system. In: 14th IEEE International Symposium on Modeling, Analysis, and Simulation of Computer and Tele Communication Systems, 2006. MASCOTS 2006, pp. 413–421 (2006)

12. Hartman, R.D.: Architecture and measured characteristics of a cloud based internet of things. In: 2012 International Conference on Collaboration Technologies and Systems (CTS), IEEE, pp. 6–12 (2012)
13. Yinjin, F., Hong, J., Nong, X., Lei, T., Fang, L.: AA-Dedupe: an application-aware source deduplication approach for cloud backup services in the personal computing environment. In: 2011 IEEE International Conference on Cluster Computing (CLUSTER), pp. 112–120 (2011)
14. Lei, X., Jian, H., Mkandawire, S., Hong, J.: SHHC: a scalable hybrid hash cluster for cloud backup services in data centers. In: 2011 31st International Conference on Distributed Computing Systems Workshops (ICDCSW), pp. 61–65 (2011)
15. Yang, Z., Yongwei, W., Guangwen, Y.: Droplet: a distributed solution of data deduplication. In: 2012 ACM/IEEE 13th International Conference on Grid Computing (GRID), pp. 114–121 (2012)