



ISSN: 2454-9940



**INTERNATIONAL JOURNAL OF APPLIED
SCIENCE ENGINEERING AND MANAGEMENT**

E-Mail :
editor.ijasem@gmail.com
editor@ijasem.org

www.ijasem.org

SECURE STORAGE AND DATA SHARING SCHEME USING PRIVATE BLOCKCHAIN-BASED HDFS DATA STORAGE FOR CLOUD COMPUTING

**RAVI KIRAN KUMAR TERA 1, K GANGOTHRI 2, DHANAVATH ANIL KUMAR 3,
DODDAMANI SHYAMKIRAN 4, B KALPANA 5,**

ABSTRACT – The storage of a vast quantity of data in the cloud, which is then delivered via the internet, enables Cloud Computing to make doing business easier by providing smooth access to the data and eliminating device compatibility limits. Data that is in transit, on the other hand, may be intercepted by a man-in-the-middle attack, a known plain text assault, a selected cypher text attack, a related key attack, or a pollution attack. Uploading data to a single cloud might, as a result, increase the likelihood that the secret data would be damaged. A distributed file system extensively used in huge data analysis for frameworks such as Hadoop is known as the Hadoop Distributed File System, more commonly referred to as HDFS. Because with HDFS, it is possible to manage enormous volumes of data while using standard hardware that is not very costly. On the other hand, HDFS has several security flaws that might be used for malicious purposes. This highlights how critical it is to implement stringent security measures to make it easier for users to share files inside Hadoop and to have a reliable system in place to validate the shared files' validity claims. The major focus of this article is to discuss our efforts to improve the security of HDFS by using an approach made possible by blockchain technology (hereafter referred to as BlockHDFS). To be more precise, the proposed BlockHDFS uses the Hyperledger Fabric platform, which was developed for business applications, to extract the most value possible from the data inside files to provide reliable data protection and traceability in HDFS. In the results section, the performance of AES is superior to that of other encryption algorithms because it ranges from 1.2 milliseconds to 1.9 milliseconds. In contrast, DES ranges from 1.3 milliseconds to 3.1 milliseconds, three milliseconds to 3.6 millimetres, RC2 milliseconds to 3.9 milliseconds, and RSA milliseconds to 1.4 milliseconds, with data sizes ranging from 910 kilos.

Index Terms – Cloud Computing, Hadoop Distributed File System, Blockchain, Authenticity, Data Security, DES, AES.

1. INTRODUCTION Over the last decade, research consortiums have achieved significant technical developments to adopt Data-sharing methodologies. Collaboration with others and making informed choices are two ways research-based activities might become more effective. The exchange of data is the first necessary step toward deriving the possible advantage from advances in

research [1]. However, it is also essential to be aware of the "three W's" for sharing, which are "what," "when," and "where." Before beginning the data-sharing process, these questions must be resolved to everyone's satisfaction.

ASSISTANT PROFESSOR 1, UG SCHOLAR 2,3,4&5

DEPARTMENT OF CSE, MNR COLLEGE OF ENGINEERING AND TECHNOLOGY, MOHD.SHAPUR, TELANGANA 502285

Cloud servers, a kind of centralised storage, are responsible for storing excessive data. One of the many hazards connected with centralised authorities is the possibility of a single point of failure. Third parties are brought in to offer data backups so that incidents like this may be avoided [2]. A blockchain gives trust and transparency, which helps remove the need for a third party to build a trust-based paradigm. Decentralised storage is a system that allows data to be stored independently on separate network nodes in the form of a distributed ledger. This may be accomplished via the use of a distributed ledger. Decentralised storage is a solution that was developed by blockchain technology. The issue is that network nodes have a limited capacity for both storage and computation. Interplanetary File System (IPFS), an architecture based on peer-to-peer communication, has been adopted for this purpose [3]. There is no possibility of failure at a single site [4]. It is very much like web3, but with some significant differences. It operates like a bit torrent and can perform content addressing [5]. Storing data on IPFS, a decentralised network, and guarantees that it will always be available when required. Hadoop is an open-source distributed file system, often known as HDFS [6,7], and is one of the business world's most well-liked choices for the batch processing of enormous volumes of

data. [6,7] It has a stellar reputation thanks to its low latency and fast throughput for data access. There are various native file systems besides HDFS, such as Ceph [8], GPFS [9], and Hydra [10]. HDFS comes preconfigured with the MapReduce programming framework. Users may use one or more files while working with MapReduce to map and reduce (sort) the data that is included inside those files to produce the required output. However, to finish MapReduce [11] tasks, tiny amounts of Java code need to be written. After MapReduce has finished reading a file from an input directory, it may be told to generate the required output in a separate output directory. HDFS may have been successful because of the way it was designed. While HDFS takes its cues from the UNIX file system, it differs because the information is spread over several drives rather than a single one. HDFS also offers the extra advantage of operating on affordable, commonly accessible hardware. Hadoop Distributed File System is the foundation upon which the powerful free and open-source framework known as Apache Hadoop is constructed (HDFS). Hadoop was developed to create a system with high throughput and resistance to failure. Data may be loaded into HDFS in either the command line or application programming interfaces (APIs). Both the command line and the API contain a variety of commands

that may be used to carry out a range of file activities, and if a remote connection is required, SSH may be used to establish it. To mention just a few of HDFS's features, users may create, read, delete, rename, list, and check the status of files. After HDFS is used to store data, several ecosystem frameworks, like Spark [12], Hive (<https://hive.apache.org/>) [13], and HBase [14], may be used to analyse the data. Remember that ecosystem apps often use the Hadoop Main API to access and manipulate data stored in HDFS. Integrating an encryption technique into IPFS's hashes of uploaded data is how the distributed file system (IPFS) ensures that data is secure. In a later step, the owner encrypts these hashes by using the Shamir secret sharing (SSS) [10] technique. This scheme breaks the hash up into an n-number of encrypted shares. The smart contract is where the encrypted shares are kept safe. The consumer initiates the process by submitting a request to view the data, which is subsequently validated using digital signatures. After the consumer has successfully downloaded the data using the decrypted shares, they must register reviews regarding the data using the review system. An incentive, which is the reimbursement of 10% of the real money placed by the client, has been proclaimed to encourage consumers to leave reviews and register for the contest

2. LITERATURE WORK Data security issues from storing information over several clouds are the topic of this article, which proposes a method called Proficient Security over Distributed Storage to address these issues (PSDS). PSDS classifies the information as either normal or sensitive. The private information is also partitioned into two parts. Data is typically encrypted and uploaded to a single cloud; however, with this setup, individual components and the cloud are encrypted and distributed among several clouds. Sensitive information from many clouds is pooled when encryption is complete. Evidence has been shown that the PSDS is secure against a wide variety of attacks, such as related key attacks, pollution attacks, chosen ciphertext attacks, and attacks based on previously known plain texts. These results came from testing the PSDS against a variety of threats. Compared to STTN and RFD, PSDS encrypts data significantly faster due to its shorter computation time. [15] A secure slicing-based resource orchestration (SS-RO) method has been developed to reduce the impact of sliceinitiated attacks by optimising for low latency and low resource consumption. In this procedure, the interslice orchestration result is obtained by the Benders decomposition, while the interslice orchestration solution is obtained via the quadratic transformation technique.

Both of these approaches will be discussed in detail below. The findings of the experiments show that the suggested SS-RO algorithm beats the baseline approaches in terms of the percentage of attacking jobs it accepts, the amount of energy it consumes, and the quantity of work it gets via the system. [16] Electronic medical records, often known as EMRs, are kept at several different institutions and managed by a cloud service provider. Patients, who are the genuine proprietors of their private and sensitive EMRs, might lose track of them. To guarantee that all components of smart health-care systems can share electronic medical records (EMRs), this article aims to construct an access control framework that uses blockchain technology and smart contracts based on distributed ledger technology. For the objectives of user authentication, access authorisation, the detection of improper conduct, and the termination of access, we recommend using four smart contracts. After being encrypted using ECC and EdDSA, electronic medical records (EMRs) are kept in the cloud, and their hashes are published to a blockchain. A private Ethereum network is used to analyse real-time smart health care access control architecture. [17] The Internet of Things affects everything. IoT enables remote health monitoring. Medical data is plentiful. IoT has boosted volume. To have useful health data, we must preserve them

effectively. We propose a cloudlet-enabled IoT e-health framework. This e-Health platform makes real-time cloudlet data retrieval easy. We design a healthcare data management system to store and handle end-user requests. NoSQL stores patient data. The proposed model analyses data transmission time, energy utilisation, query response time, and packet loss. We proved our model's superiority by comparing its findings to those competing for cloud-based e-Health systems. [18] IoT technologies like smart meters, appliances, and grids may improve energy efficiency and customer service. Under standard AIoT paradigms, users' IoT energy data must be transported to a central repository (such as the cloud or an edge device) for knowledge extraction. Data exploitation and privacy issues may result. They provide an AIoT solution that works with edge clouds to share energy usage data in smart grids safely. Simulations show that the suggested strategy may increase communication and motivate EDOs to upgrade their local models. [19] Significant changes have been made to load balancing outsourced in the cloud. The deploying blocks use a design for a distribution record that is shared across all servers to send data in a time-consent fashion without interfering with the usual functioning background modules of the record base a state machine. External data may be integrated and analysed. With

cloud load balancing, several servers, networks, or individual PCs may share the processing of a single task. DopCloud optimises the packet flow routing in the cloud by analysing it and traversing several stack layers at a microsecond scale. The server and client codes must be updated with each deployment to maintain reliable load balancing and cloud backups. Due to an abrupt in data storage, electronic health record (EHR) systems in the medical profession were forced to cease operations, resulting in data loss. Adjustments are piled to control the cloud network, making those controls available for the diagnostic data load. Here, RIFT and VNF cloud provide a more trustworthy protocol that promises to protect massive data loads. It is the packet-flowing route optimisation (DOPER) system that does the bulk of the planning for the containerisation of the state machine, and no communication is allowed between DopCloud peer nodes until all contracts are made accessible to all nodes. [20] Cloud networks safeguard the processing of data and its transmission to users with appropriate permissions. "cloud computing" refers to a data storage and management model that uses several computers in different parts of the world. There is a need for load balancing. Within the scope of this research project, a hybrid heuristic mathematic algorithm (HHMA) is proposed for IaaS computing networks to

resolve issues with resource allocation. To speed up the design process and download of files, improved K-means clustering was used to divide the cluster into many little parts (heuristic). MCSO uses real-time constraints to determine a determined node's load ratio. Following the execution of the MCSO algorithm, the optimal value for the complete evaluation is determined and used to allocate controller nodes to storage nodes. Two strategies distribute request data to relevant nodes for efficient processing. In comparison to earlier approaches, this one uses far fewer resources. The simulation findings imply that the suggested technique would lessen the burdens on memory, reaction time, and network overhead. Investors in cloud computing may benefit from the use of the recommended heuristic method. [21] The dawn of the information age in China resulted in a rise in the amount of documentation produced by end users. This came with it the difficulty of finding out how to store enormous volumes of data in a secure and straightforward manner. The technology of the cloud could solve this issue. People have been able to reduce their demand for hardware by adopting cloud storage, yet there is still a worry over hardware security. This work discusses dividing and merging source documents, encrypting the division table, recovery, and backup. [22] WukaStore can provide a

huge data storage solution that is scalable, adaptable, and reliable by combining non-volatile storage, such as that provided by the cloud, with volatile storage, such as that which is harvested from unused storage space on desktop computers via the internet. By customising a number of different storage methods, WukaStore offers storage that is efficient and affordable. In order to investigate how to guarantee the high availability and durability of WukaStore, we use trace-driven simulations. The open-source Big Data middleware BitDew serves as the foundation for WukaStore. On the experimental platform that France has developed, Grid'5000, we evaluate how well WukaStore performs. [23] Customers preserve several copies of cloud-hosted data to increase its availability and durability. Multi-copy data's integrity is ensured via PDP methods. All PDP copies are stored on a single cloud server. This shouldn't be duplicated. Many PDP protocols need expensive and insecure public key infrastructure (PKI). We present an identity-based PDP system that stores data on various cloud servers to increase safety and efficiency. [24] The cloud is being used by many corporations, academic institutions, government agencies, and other organisations because of its cheap initial cost, scalability, and other advantages. The cloud offers a great many

advantages, but it also has a great many disadvantages. Data protection is given top priority by both information security and cloud computing. This problem can be solved in a few different ways. Because the already available solutions have not been subjected to a comprehensive analysis, it is required to conduct research, classify, and evaluate the work that has already been done to determine whether or not it meets the requirements. This article compares and analyses the primary cloud data sharing and protection methods. The discussion of each specific method contains the following topics: data security functions, prospective and novel solutions in the field, workflow, accomplishments, scope, gaps, future directions, etc. In addition, a comparative analysis of the approaches has been presented. After that, a discussion of the applicability of the techniques follows, after which research gaps and potential future projects are outlined. The writers of this paper believe that it will inspire the next generations of researchers to investigate the subject. [25] Users may easily back data to the cloud and save archives with MCS. Here, we focus on how frequently data is accessed to create a private, secure, and effective cloud storage solution for mobile devices. As the central mechanism of the mobile cloud storage system, we propose an OSU protocol. The client's computation and communication overheads are reduced

since they only have to construct a tiny encrypted vector to get encrypted data from the cloud and update it. In contrast to prior research, our approach provides a finegrained data structure with a tiny item size, requiring just a few additively homomorphic operations on the client side and a constant communication cost. Because of the usage of verification chunks, our technique is immune to attacks from malicious clouds. Our method is more efficient when evaluating client and cloud workloads than the current storage solutions. [26] Access control methods in cloud storage are becoming important as cloud computing evolves. In business, the Chinese Wall is a tried-and-true solution for dealing with the CoI issue. The capacity to store conflict-proof data in the cloud might be useful. Users' interests, investing choices, and the access control mechanism may reveal other personal information since it does not provide anonymity. Problems arise while building the Chinese Wall without compromising users' privacy. Cloud storage is implemented utilising the Chinese Wall technique to protect user access patterns. The C-Wall protocol will be discussed next once the treebased Chinese Wall access control has been introduced. Protecting the user's anonymity, our C-Wall is impenetrable to hostile actors and may be built anywhere in the world. We further extend C-Wall to privacy-protecting

cloud storage with C2-Wall, which keeps all the advantages of C-Wall and stops "honest but inquisitive" cloud servers from accessing private data. These two additions represent significant advances in the discipline. When it comes to our C2-Wall, we test and check it thoroughly. According to experiments, it has real-world applications. [27] Computing on the cloud presents an opportunity for storage solutions. Concerns about the safety of cloud storage might limit its expansion. Cloud-stored data is vulnerable to malicious changes and data loss. Recent research outlines a three-layer fog server architecture for cloud storage. Hash algorithm and Hash-Solomon code are modified. It didn't improve the ability to identify changes or recover lost data, but it lost less to cloud servers. This article describes fogcentric secure cloud storage to prevent unauthorised data access, modification, and deletion. The recommended system hides data using XorCombination and XorCombination. Block Management outsources XorCombination results to prevent malicious retrieval and aid data recovery. A hashbased technique improves change detection. The system's security is analysed. Experiments show the proposed approach is faster than existing options. [28] There is a growing need for object storage in the cloud to manage massive amounts of large binary

objects (BLOBs), which include movies, images, and documents. The vast majority of the data encoding strategies being used are not ideal for the architecture of cloud object storage, despite the way that several companies and organisations want to utilise public cloud object storage services like Amazon Simple Storage Service (S3). In this investigation, we provide a technique called dynamic extreme erasure encoding, often called DexEncoding. Its purpose is to enhance client utility by dynamically optimising encoding sites between gateway and storage servers in a cloud environment that varies over time. The utility gauges the degree to which customers are happy with the speed and fairness of data storage. To efficiently relieve resource limits, DexEncoding can adapt to the availability of the network, processing, and storage resources, as well as storage requests. Simulations driven by real-world measurements demonstrate that DexEncoding provides higher levels of customer satisfaction than other cutting-edge object storage systems. [29] IoT and smart diagnostic implants improve medical systems by accessing remote patient data and screening for health risks. IoT medical gadgets monitor patients' health and transfer data to the cloud for doctors. When it comes to the exchange of data and outsourcing, the privacy and secrecy of patients' electronic medical records are very

important considerations that must be considered. In this study, a solution for securing data storage and regulating access is proposed for use in intelligent healthcare systems. This article includes a password, a psychometric slot enabled by deep learning, and an evaluation based on deep learning. The proposed solution has undergone testing and demonstrated that it is feasible in comparison to the existing access control systems. [30] Together with the surge in the popularity of cloud computing in general, the demand for cloud-based block storage services has increased. Performance optimisation becomes more difficult when resource demand is not constant throughout a cloud block storage system. [31] Understanding how to verify the integrity of data stored in the cloud is becoming an increasingly vital skill. ID-based proved data possession (PDP) is an auditing method for cloud storage that does not need a certificate and does not require user data. Cloud users must offshore data blocks, authenticators, and a small file tag to use an existing ID-PDP system. Additionally, bilinear pairing and elliptic curve cryptography are required. These drawbacks would increase storage, connectivity, and computing expenses, which was not anticipated for users of cloud services with limited resources. This is made possible by the absence of substantial cryptographic processes, contributing to the

protocol's overall effectiveness. The recommended protocol may give additional helpful functions thanks to its usage of primitive replacement. The suggested protocol cannot be falsified, can be found, is accurate, and is detectable. In conclusion, we provide both theoretical and experimental data to validate the recommended technique. [32] CDS is a popular cloud-based solution because of its cheap cost and excellent efficiency in disseminating and sharing content with end-users, partners, and insiders. By duplicating commonly used resources, this service increases both accessibility and I/O speed. This adds more work for the network and storage facilities. This study proposes a threepronged strategy for enhancing the effectiveness of I/O operations in the cloud and increasing CDS use. [33] As cloud storage services grew more common, there was an increase in worry over the integrity of data that was outsourced and stored on servers that could not be trusted. Provable data possession (PDP) protects the integrity of cloud-stored data by demanding the cloud server to prove the data hasn't been updated or deleted without giving users access to the actual data. ID-P33DP is an identity-based, privacy-preserving data possession mechanism for secure cloud storage. A cloud user must utilise an outsourced file and a global parameter to create ID-P33DP authenticators. Any third-

party auditor (TPA) may ascertain whether an outsourced file was retained in its original state by validating homomorphic authenticators. ID-P33DP supports the aggregation of user-generated identity-based homomorphic authenticators using RSA. The cloud may aggregate identitybased homomorphic authenticators to create a data possession proof that verifies data integrity. A zero-knowledge proof may prevent TPA data leaks. RSA proves ID-soundness, P33DPs' privacy, and TPA's privacy. Cross-user aggregate verification reduces TPA's computational and communication cost [34].

CONCLUSION The storage of a vast quantity of data in the cloud, which is then delivered via the internet, enables Cloud Computing to make doing business easier by providing smooth access to the data and eliminating device compatibility limits. Data that is in transit, on the other hand, may be intercepted by a man-in-the-middle attack, a known plain text assault, a selected cypher text attack, a related key attack, or a pollution attack. Uploading data to a single cloud might, as a result, increase the likelihood that the secret data would be damaged. The Hadoop Distributed File System (HDFS) is a popular choice among distributed file systems for large-scale data processing on Hadoop and similar frameworks. We plan to enhance the safety

of HDFS by using a method enabled by blockchain technology (hereafter referred to as BlockHDFS). More specifically, the proposed BlockHDFS uses the Hyperledger Fabric platform, which is designed for corporate usage, to make the most of the information inside files to provide trustworthy data security and traceability in HDFS. In the results section, the performance of AES is superior to that of other encryption algorithms because it ranges from 1.2 milliseconds to 1.9 milliseconds. In contrast, DES ranges from 1.3 milliseconds to 3.1 milliseconds, RC2 milliseconds to 3.6 milliseconds, and RSA milliseconds to 1.4 milliseconds, with data sizes ranging from 910 kilobits.

REFERENCE

[1] G. Kumar et al., "A Novel Framework for Fog Computing: LatticeBased Secured Framework for Cloud Interface," in *IEEE Internet of Things Journal*, vol. 7, no. 8, pp. 7783-7794, Aug. 2020, doi: 10.1109/IJOT.2020.2991105.

[2] X. Liu, G. Yang, Y. Mu and R. H. Deng, "Multi-User Verifiable Searchable Symmetric Encryption for Cloud Storage," in *IEEE Transactions on Dependable and Secure Computing*, vol. 17, no. 6, pp. 1322-1332, 1 Nov.-Dec. 2020, doi: 10.1109/TDSC.2018.2876831.

[3] Benet, J. Ipfs-content addressed, versioned, p2p file system. arXiv 2014, arXiv:1407.3561.

[4] J. Wei, X. Chen, X. Huang, X. Hu and W. Susilo, "RS-HABE: Revocable-Storage and Hierarchical Attribute-Based Access Scheme for Secure Sharing of e-Health Records in Public Cloud," in *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 5, pp. 2301-2315, 1 Sept.-Oct. 2021, doi: 10.1109/TDSC.2019.2947920.

[5] H. Wang, L. Feng, Y. Ji, B. Shao and R. Xue, "Toward Usable Cloud Storage Auditing, Revisited," in *IEEE Systems Journal*, vol. 16, no. 1, pp. 693-700, March 2022, doi: 10.1109/JSYST.2021.3055021.

[6] Apache Hadoop, URL, <http://hadoop.apache.org>, 2006.

[7] K. Shvachko, H. Kuang, S. Radia, R. Chansler, The hadoop distributed file system, in: 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies (MSST); 3-7 May 2010; Incline Village, NV, USA, IEEE, Piscataway, NJ, USA, 2010, pp. 1-10.

[8] S.A. Weil, S.A. Brandt, E.L. Miller, D.D.E. Long, C. Maltzahn, Ceph: a scalable, high performance distributed file system, in: *Proceedings of the 7th Symposium on Operating Systems Design and Implementation, OSDI '06*; 6-8 Nov

2006; Seattle, WA, USA, USENIX Association, Berkeley, CA, USA, 2006, pp. 307–320.

[9] F. Schmuck, R. Haskin, Gpfs: a shared-disk file system for large computing clusters, in: Proceedings of the 1st USENIX Conference on File and Storage Technologies, FAST '02; 28–30 Jan 2002; Monterey, CA, USA, USENIX Association, Berkeley, CA, USA, 2002.

[10] C. Ungureanu, B. Atkin, A. Aranya, et al., HydraFS: A highthroughput file system for the hydrastor content-addressable storage system, in: Proceedings of the 8th USENIX Conference on File and Storage Technologies, FAST'10; 23–26 Feb 2010; San Jose, CA, USA, USENIX Association, Berkeley, CA, USA, 2010, pp. 225–239.

[11] J. Dean, S. Ghemawat, Mapreduce: simplified data processing on large clusters, *Commun. ACM* 51 (1) (2008) 107–113.

[12] M. Zaharia, R.S. Xin, P. Wendell, T. Das, M. Armbrust, A. Dave, X. Meng, J. Rosen, S. Venkataraman, M.J. Franklin, et al., Apache spark: a unified engine for big data processing, *Commun. ACM* 59 (11) (2016) 56–65.

[13] X. Chen, L. Hu, L. Liu, J. Chang and D. L. Bone, "Breaking Down Hadoop Distributed File Systems Data Analytics Tools: Apache Hive vs. Apache Pig vs. Pivotal HWAQ," 2017 IEEE 10th

International Conference on Cloud Computing (CLOUD), 2017, pp. 794–797, doi: 10.1109/CLOUD.2017.117.

[14] M.N. Vora, Hadoop-hbase for large-scale data, in: Proceedings of 2011 International Conference on Computer Science and Network Technology; 24–26 Dec 2011; Harbin, China, IEEE, Piscataway, NJ, USA, 2011, pp. 601–605.

[15] F. Shahid, H. Ashraf, A. Ghani, S. A. K. Ghayyur, S. Shamshirband and E. Salwana, "PSDS–Proficient Security Over Distributed Storage: A Method for Data Transmission in Cloud," in *IEEE Access*, vol. 8, pp. 118285–118298, 2020, doi: 10.1109/ACCESS.2020.3004433.