

A Novel Approach for Identification of Hadoop Cloud Environment

Chella Esther Varma, Associate Professor,

Vuppala Bhavana Eswar, Assistant Professor,

M.Sunil Kumar, Assistant Professor,

Geethanjali College Of Engineering And Technology, Cheeryal, R.R. Dist.

ABSTRACT: Due to the modern-day tendencies within the discipline of science and technology resulted in the trends of effective information switch, ability of handling colossal data and the retrieval of data efficiently. On account that the data that is saved is increasing voluminously, methods to retrieve relative understanding and protection associated concerns are to be addressed effectively to at ease this bulk data. Additionally with emerging standards of giant information, these security issues are a challenging undertaking. This paper addresses the limitation of comfortable data switch using the principles of data mining in cloud environment making use of hadoopmapreduce. Based on the experimentation achieved outcome are analyzed and represented with recognize to time and space complexity when compared hadoop with non hadoop approach.

KEYWORDS- Big Data, Hadoop, Mapreduce, Cloud Computing, Temporal Patterns

I. INTRODUCTION

The contemporary technological traits witnessed the storage of enormous data and methodologies special towards efficient retrievals. For the reason that this information is on hand are surmounting, safety breaches and upholding the privacy is a foremost trouble. These security issues are way more difficult whilst on the grounds that the data transfers in cloud atmosphere or parallel processing architectures [1]. So as to manage this data effectually ideas of mapreduce [2] is concentrated in the literature. This is as a result of its capabilities of faulttolerance and scalability in conjunction with simplicity.

One other fundamental advantage of highlighting the mapreduce notion is it allows the parallel processing environments which aid not directly in the direction

of colossal data storage [3]. The notion of mapreduce may also be simply implemented utilisinghadoop environment [4]. Many methodologies had been discussed in literature [5, 6, 7,8] to deal with the disorders of safety in client server environment. However among the many restrained algorithms used for safety in dispensed environments Symmetric Encryption is most of the time projected as a result of its robustness and capability [12] of utilization in each 64-bit and 128-bit key format. Within the trendy scenario as a result of the broaden in the rate of application, upkeep of application, storage of application, pressured the clients or manufacturers to undertake cloud computing atmosphere. In this atmosphere the software or data is stored principally within the form of clusters. These clusters will likely be transmitted over a cloud based on the clients request varieties which will also be SAS, PAS and IAS [8]. Among the many different offerings offered by way of the cloud atmosphere, the ordinarily used services comprise providing occasions on demand and offering computational capabilities on demand.

The map shrink concept addressed in this paper helps the dispensed computing for large data units on clusters of computer systems for offering computing ability on demand. To facilitate this carrier hadoop is ordinarily used due to its capacity of handling HDFS records in which data related to exclusive machines alongside the globe may also be saved. Mapreduce is a performance of hadoop which helps in data preprocessing. This preprocessed data can be worthwhile for the effective evaluation of bigdata. Data mining is the exploration of data with the purpose of discovering hidden structure.

In many real-world functions, it is fundamental to study the change of temporal sides of a nonstationary time sequence, and identify the ones which can be representing the value of time situations. For example, it is relevant in data leakage functions from where the data has been leaked or it is elaborate to identify IP of an unauthorized user who logged at any time or irregular interval of time in a cloud atmosphere as a rule such time series are viewed non-stationary. Traditional time sequence evaluation employs statistical approaches to model and provide an explanation for the data and predict future values of the time series. It's not convenient, nonetheless, to determine the primary temporal patterns of the time series using these normal methods. Utilizing a suite of observations, in this paper, we gift a brand new process for time series data mining. By incorporating symmetric key encryption with the use of hadoop, temporal patterns (user's log history at regular or irregular time interval) will also be effortlessly published in non-stationary (cloud) atmosphere. In order to handle the colossal data and transmit the data throughout the globe effective data switch methodologies are to be adopted by means of making use of symmetric encryption.

II. PREVIOUS WORK

P. Srinivasa Rao et al [18] proposed an approach to look after web usage from unauthorized clients with the aid of utilizing hadoop mapreduce where a namenode log file technique is proposed wherein identification of user's temporal patterns process is experimented.

Elisa Bertino et al [10] proposed a process of Digital identification administration for a cloud using Multifactor Authentication manner. S. Fischer-Hubnar et al [11, 15, 16] proposed a privacy and identification administration for Europe where it presents privacy maintenance Authentication using Errorneous Credentials.

Basker Prasad Rimal et al [8]. proposed a process in working out of taxonomy and survey of cloud computing methods. Kumar Gunjan et al [13] gave an outline thought of identity administration in cloud computing. Mark D. Ryan et al [14, 17] explained cloud computing protection: the scientific assignment

and survey of options. Taking into consideration all of the above disorders, on this paper we are going to deal with defense of colossal data that has been transmitted in cloud through Hadoop dispersed process by means of making use of DES Algorithm.

To reduce the delay because of decryption process on the receiving finish, an alternative procedure can also be adopted for authorized users by sending raw data. In this paper we propose a novel methodology the place in the safety will be offered in two phases in Hadoop Cloud atmosphere.

III. SYSTEM MODEL

In the Figure 1 illustrated above, at any slave (datanode) user may send the data to any other datanode in the cloud. The process of sending the data securely to other destination node is clearly visualized in the below Fig. 1

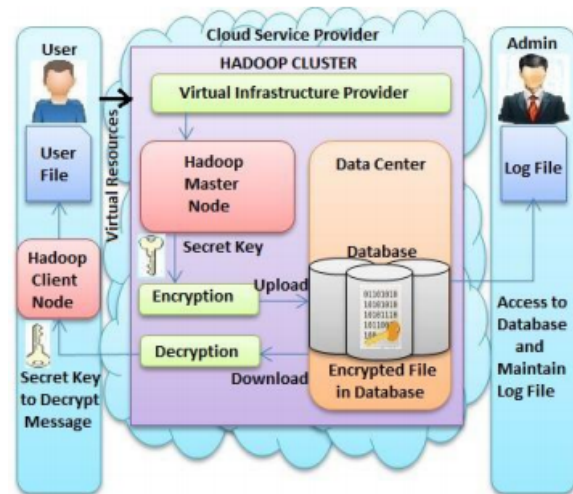


Fig. 1: Architecture of Hadoop Cluster in Cloud

At any datanode if the user is getting authorization to enter into cloud, he can be allowed to send or receive data that can be processed is shown in below Fig. 2.

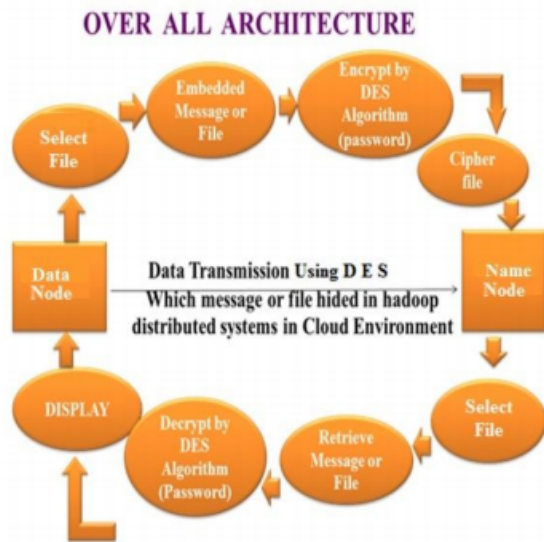


Fig. 2: Encrypted Data Transfer

A. Methodology

This paper addresses the method of mighty dataswitch to the authenticated clients and it comprises a mechanism where in the unauthorized men and women may also be blocked from receiving the data, and among the many approved users a security tag is attached so that we could determine the source of data leakages. The users for whom the data is transmitted assuming to be authorized, if an data is leaked the security tag is ready to 1 else this security tag is ready 0. For all these security tags where the flag is one, the corresponding IP of the clients might be scrutinized and an error warning might be notified. If the method is repeated the user or a individual corresponding to detailed data node with the particular IP will be blocked from receiving additional data. In an effort to avert unauthorized users to view the content material the data is encrypted making use of DES and File key symmetric encryption algorithms. Symmetric encryption which is used on this paper is extra nice than the uneven encryption which requires extra CPU cycles and CPU reminiscence moreover to a couple of obstacles explained by package File white et al [12].

B. Mapper and Reducer:

The mapper that includes my TMap algorithm is applied to each input data that has been transmitted in

hadoop allotted atmosphere. The data that is transmitted will have to be encrypted at each and every node with DES and Mapkey. The encrypted bundle of data will be stored at a customary memory of hadoop called HDFS (Hadoop allotted File system). The patron called datanode is allowed to learn the understanding blocks if he is having authorization, otherwise tag worth will probably be incremented through quantity of times he attempted to grab the understanding. When the user at a distinctive data node is having a tag price more than 0 can be recorded at log file of namenode so that the writer will not be allowed to access any extra understanding in that cloud. That's the IP associated with that unique user shall be blocked. This processing shall be carried out within the Mapper and reducer whose job is to segregate all url's of a distinctive user so that temporal patterns of the user can be located. A common algorithm with Map and cut back functions for determining such temporal patterns is listed within the table. 1.

Table 1: TMap and TReduce Algorithm

Set the input path and the output path
Step 1: Client selects the file at any datanode.
Step 2: Authentication using Mapkey then goto step 3.
Step 3: Check for the captcha, if captcha is not matched and tag value is greater than 0 then goto Step 10 else goto step 4. // Start of Map Function.
Step 4: Map (key, value)
Step 5: Client can send or receive data
Step 6: Client encrypts file by using Encryption algorithm (DES) and Map key by giving Authentication (Password).
Step 7: The cipher file is transmitted over the cloud through hadoop HDFS.
// Start of Reduce Function
Step 8: Reduce (key, value).
Step 9: If at any IP, any unauthorized user is attempted to access the data, tag value will be set and

the log file of name node will be updated so that the IP

will be blocked.

Step 10: If the IP is with authorized user then decrypt data by giving password.

Step 11: The message or file is retrieved at any data node of the respective cloud.

Step 12: Logout from datanode.

C. Mapkey Algorithm:

The Mapkey algorithm is one of the security algorithms used to furnish safety for the user data and store them in an encrypted format. A random quantity is generated utilising Password based Key Derivation function (PBKDF2) Algorithm that derives the important thing with the aid of making use of SHA1, SHA256, MD5 etc., algorithms for the random generated number and again different algorithms like AES, DES and so on., are applied through settling on the random quantity as a secret key. When a user uploads the file in a cloud, the protection algorithms are applied over the user file and encrypt the file and then the cloud supplies an encrypted output file. These files are saved in a cloud storage database, which also provide high degree safety to the cloud computing environment. A secret key is supplied to the licensed user to access his records in the cloud environment.

MapKey Algorithm:

1. Start
2. Read user file from at any datanode
3. Generate Random Number (n) // e.g.: 12345
4. Perform PBKDF2 Algorithm to derive the Key (k)
5. Return MapKey
6. Stop

PBKDF2 Algorithm:

Input: Pwd Password

S salt Function

Ic Iteration Count

Kl Key length in bits $(2^{32} - 1) * Hl$

Parameters: Prf \rightarrow HMAC Function

Hl \rightarrow Hash Function Digest System

Output: Mk \rightarrow Master Key (Mk)

Algorithm: if $(Kl > (2^{32} - 1) * Hl)$

Return Error and Stop

Initialize $L \rightarrow [Kl / Hl]$

$Q = Kl - (L - 1) * Hl;$

For $(i = 1$ to $l)$

$Xi = 0;$

$V0 = S \parallel int(i);$

For $(j = 1$ to $l)$

$Vj = HMAC(Pwd, Vj-1);$

$Xi = Xi XOR Vj$

Return $Mk = Xi \parallel X2 \parallel \dots \parallel Xi // < 0 \dots Q - 1 >$

As shown in the above TMAP algorithm, the data encryption standards are one of the most protection algorithms used to furnish security to the person data and retailer them in an encrypted format. When a consumer uploads the file in a cloud, the protection algorithms are applied over the user file. These files are stored in a HDFS which is a normal hadoop Storage subject for all data nodes in a hadoop distributed atmosphere. A secret key is offered to the licensed users to access this data within the cloud environment.

IV. CONCLUSION

In this paper, a process is provided to explain the safety framework for cloud atmosphere. This framework helps in offering the protection to the user data in an encrypted layout that are uploaded by using the provider user right into a cloud, by incorporating the key facets of distinct algorithms like DES, FileKey methods, that are placed in a Hadoop cluster. Key ideas of this structure are the definition of unique safety parameters for expressing security requirements and security performance, we also provide a safety approach to the cloud environment with the aid of fusing the safety capabilities and following the protection parameters and safety insurance policies on the time of user login to provide authentication to the user and to search out temporal patterns within the cloud.

REFERENCES

- [1] KejaingYe et al., vHadoop: A Scalable hadoop virtual cluster platform for MapReduce-

Based Parallel Machine learning with Performance Consideration, 2012 IEEE International Conference on Cluster Computing Workshops, PP: 152-160, 2012

[2] J. Dean and S. Ghemawat, "Map Reduce: Simplified Data processing on large clusters", *communications of the ACM*. Vol. 51, no. 1, pp 107-113, 2008.

[3] T. White, *Hadoop: The Definitive guide*. yahoopress, 2010.

[4] Mohamed H. Almeer et al "cloud hadoop mapreduce for remote sensing image analysis" *Journal of Emerging Trends in Computing and Information Sciences*, VOL. 3, NO. 4, April 2012.

[5] Abhipal Singh et al "File Transfer Using Secure Sockets in Linux Environment", *Proceedings of the 4th National Conference; INDIACOM-2010*.

[6] Matthieu Bloch et al "Network Security for Client Server Architecture Using Wiretap Codes", *IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY*, VOL. 3, NO. 3, SEPTEMBER 2008

[7] Xinyi Huang et al "Further Observations on Smart-Card-Based Password-Authenticated Key Agreement in Distributed Systems", *IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS*, 2013.

[8] Bhaskar Prasad Rimal et al "A Taxonomy and survey of cloud computing system", *Fifth International Joint Conference on Information Systems, IDC*, 2009.

[9] Poulami Dutta et al "Data Hiding in Audio Signal: A Review" *International Journal of Database Theory and Application* Vol. 2, No. 2, June 2009

[10] Elisa Bertino et al, "privacy-preserving digital identity management for cloud computing" *IEEE computer society technical committee on data Engineering*, 2009.

[11] S. Fischer-Hubner and H. Hebdon, "PRIME privacy and identity management for Europe" August 2010.

[12] kitu File white et al "symmetric vs Asymmetric encryption" 2010.

[13] Kumar Gunjan et al, *international journal of Engineering Research and Technology (IJERT)*, ISSN 2278-0181 vol 1 issue 4 june 2012.

[14] Make D. Ryan et al "Cloud computing Security: The Scientific journal of systems and Software challenge, and a survey of solutions", Elsevier, 2013.

[15] Jittin et al. "An analysis on privacy preserving in cloud computing", *International journal of computer trends and Technology (IJCTT)*, volume 4, issue 6 june 2013.

[16] Man Qi* et al, "Social Networking searching and privacy issues" *information security technical report*, Elsevier 2011.