

# Review on Big Data Security in Hadoop

Mr. Shrikant Rangrao Kadam<sup>1</sup> , Vijaykumar Patil<sup>2</sup>

<sup>1</sup>Department of Computer Science & Engineering, Engineering, JNTU University, Hyderabad Progressive College Of Engineering Hyderabad, Telangana, India,

<sup>2</sup>Department of Computer and Science, University of Pune, SKN'Sinhgad Institute of Technology, Lonavala, Pune, Maharashtra, India

\*\*\*

**Abstract** - Hadoop prominently utilized for handling vast measure of information on its dispersed programming structure with Hadoop disseminated document framework (HDFS), yet preparing touchy or individual information on circulated environment requests secure registering. Initially Hadoop was imagined with no security display. The encryption and unscrambling are utilized before composing and perusing information from Hadoop appropriated record framework (HDFS) separately. Propelled Encryption Standard (AES) empowers security to information at every bunch, it perform encryption/decoding before read/compose individually. Hadoop running on conveyed environment, it utilizes ware equipment which is a system show require a solid security instrument, also kerberos utilized for validation it offers get to control rundown and review components, to guarantee the information put away in the Hadoop record framework is secure.

**Key Words:**

Hadoop, DataNode, NameNode, TaskTracker, ASE, HDFS.

## 1. INTRODUCTION

Hadoop was produced from GFS (Google File System) [2, 3] and MapReduce papers distributed by Google in 2003 and 2004 individually. It has been prevalent as of late because of its very adaptable appropriated programming or registering system, it empowers handling huge information for information escalated applications and additionally numerous examination. Hadoop is a system of devices which underpins running application on enormous information and it is actualized in java. It furnish MapReduce programming engineering with a Hadoop circulated record system (HDFS), which has gigantic information handling capacity with a huge number of item equipment's by utilizing essentially its guide and diminish capacities.

Since Hadoop is typically executing in vast bunch or might be in an open cloud benefit. Like Amazon, Google, Yahoo, and so forth are such open cloud where various clients can run their employments utilizing Elastic MapReduce and distributed storage that is utilized as HDFS, it is fundamental to execute the security of client information on such stockpiling or bunch. Hadoop extend amid its initial outline organize the straightforward security instruments are utilized, for example, record authorizations and get to control list [4].

Encryption and unscrambling is key means for securing Hadoop record system (HDFS), where numerous DataNodes (or groups that is initially DataNodes) store document to HDFS, those are exchanged while executing MapReduce (client submitted program) work. It is accounted for that forthcoming Hadoop programming or form will incorporate encryption and unscrambling usefulness [5].

Web now creating vast measure of information consistently, IDC's distribute an insights in 2012 it incorporate the organized information on the web is around 32% and unstructured is 63%. Additionally the volume of computerized substance on web grows up to more than 2.7ZB in 2012 which is up 48% from 2011 and now soaring towards more than 8ZB by 2015. Each industry and business associations are presently a critical information about various item, creation and its market overview which is a major information gainful for profitability development. In business information examination application which is work on enormous information the Hadoop gets to be accepted stage, in forthcoming 5 year, over half of huge information applications are executing on Hadoop.

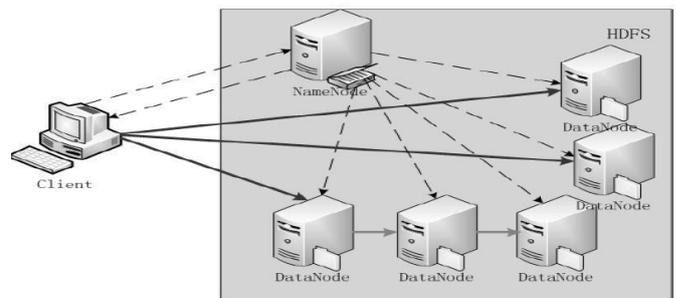


Figure 1: HDFS architecture [11]

Documents on Hadoop record framework (HDFS) are part into various squares and repeated with numerous DataNodes to guarantee high information accessibility and solidness to disappointment of execution of parallel application in Hadoop environment. Initially Hadoop bunches have two sorts of hub working as ace treatment or ace laborer design. NameNode as an ace and DataNodes are specialists hubs of HDFS. Where information documents are really situated in Hadoop is known as DataNode which just leads stockpiling. However NameNode contains data about where the diverse document squares are found yet it is not relentless, when framework begins piece may transforms one DataNode to another DataNode yet it answer to

NameNode or customer who present the MapReduce occupation or proprietor of Data occasionally [11]. The correspondence is in the middle of DataNode and customer NameNode just contains metadata.

## 2. Security Risks in HDFS

Hadoop utilizes 'whoami' and 'bash -c' groups utility of Unix for individual client and gatherings separately, this is the feeble point since which consents and record quantity are for customers. There are three sorts of security infringement in HDFS, unapproved get to, unapproved alteration of information and disavowal of administration or asset. Following are the areas where threat identify in Hadoop

**Hadoop does enforce authenticate any user or service:** unauthorized users may any HDFS cluster like owner via RPC of HTTP protocol.

- **DataNode can't have any access control mechanism to protect data block** : it is possible to write or modify existing data blocks to DataNode.
- **An attacker can presence as Hadoop service** : For example, code submitted by user register itself on MapReduce cluster as a new TaskTracker
- **Super-user of system does anything without checking**: User who takes control of NameNode is a super-user; it means somebody started the NameNode which have fully access on HDFS data.
- **An executing MapReduce may use the host operating system interfaces**: Some time execution of MapReduce demands access other tasks on the host OS, access local storage for instant Map output, but both executing on the same physical node.

## 3. Literature Review

Hadoop is initially a dispersed framework which permits us to store enormous information and backings for preparing it in parallel environment. Numerous associations utilizes huge information applications to foresee future extension, Hadoop bunch store the delicate data about such associations (data like profitability, money related information, client criticism and so on). As result Hadoop bunch require solid validation and approval with information security, for example, encryption.

The creators Seonyoung Park and Youngseok Lee, they show secure Hadoop design by including encryption and decoding capacities in HDFS in " Secure Hadoop with Encrypted HDFS ", J.J. Stop et al. (Eds.): GPC 2013, LNCS

7861, pp. 134–141, 2013 Springer-Verlag Berlin Heidelberg. They distribute secure HDFS by including the AES encode/unscramble class to CompressionCodec in Hadoop.

The creators Jason Cohen and Dr. Subatra Acharya, they show Trusted Computing Group (TCG, for example, the unavoidably accessible Trusted Platform Module (TPM) attentiveness toward accomplishing information secrecy and honesty in "Towards a Trusted Hadoop Storage Platform: Design Considerations of an AES Based Encryption Scheme with TPM Established Key Protections", IEEE tenth International Conference on Ubiquitous Intelligence and Computing and IEEE tenth International Conference on Autonomic and Trusted Computing in 2013. They distribute an encryption plot for Hadoop using equipment key securities and AES-NI for encryption speeding up.

The creators Hsiao-Ying Lin, Shiu-Tzuo Shen, Wen-Guey Tzeng and Bao-Shuh P. Lin, they display the information privacy issue by coordinating half and half encryption plans in the Hadoop dispersed record framework (HDFS) in " Toward Data Privacy by means of Integrating Hybrid Encryption Schemes and Hadoop Distributed File System", 26th IEEE International Conference on Advanced Data Networking and Applications, 2012 Springer-Verlag Berlin Heidelberg. They distribute two mixes, HDFS-RSA and HDFS-Pairing, as augmentations of HDFS.

The creators Thanh Cuong Nguyen, Wenfeng Shen, Jiwei Jiang and Weimin Xu, they introduce the plan to encode client's information locally before exchanging to HDFS on the off chance that he requires high privacy" A Novel Data Encryption in HDFS ", IEEE Worldwide Conference on Green Computing and Communications and IEEE Internet of Things and IEEE Cyber, Physical and Social Computing, 2013. They distribute novel technique to encode document while being transferred.

The creators Songchang Jin, Shuqiang Yang, Xiang Zhu, and Hong Yin, they introduce the most recent cryptography—completely homomorphic encryption innovation and confirmation operator innovation in " Design of a Trusted File System Based on Hadoop ", Y. Yuan, X. Wu, and Y. Lu (Eds.): ISCTCS 2012, CCIS 320, pp. 673–680, 2013. They distribute homomorphic encryption and validation operator innovation for securing HDFS.

The creators Monika Kumari and Dr.Sanjay Tyagi, they show three phase security model is displayed for Hadoop Environment in " A Three Layered Security Model for Data Management in Hadoop Environment", International Journal of Propelled Research in Computer Science and Software Engineering Volume 4, Issue 6, June 2014. They distribute secure record administration and dissemination over the protected Hadoop environment.

### 4. Secure Hadoop

Hadoop engineering comprises of an ace and all others are slaves. Ace contains NameNode that oversees metadata and get to control of record framework for mapping, DataNode and square of document, slaves are DataNode which store information. The HDFS contains information in piece of settled size, of course square size is 64 MB. Every square is recreated three circumstances in various DataNode, even in the wake of preparing or each time Hadoop keeps up replication figure three. Hadoop give MapReduce programming model which split employment into different assignments (guide or lessen) to process more than one HDFS information hinders in parallel. HDFS underpins a compose once-read-many model.

Secure Hadoop encode each record before written in HDFS. It is accounted for that each DataNode or slave is an item server which perform encryption or decoding at nearby site utilizing its CPUs [1]. Propelled Encryption Standard (AES) is most well known calculation that bolster square figure, henceforth it is appropriate for HDFS pieces. AES accessible with 128-piece AES, 192-piece AES and 256-piece AES, 128-piece AES is utilized the greater part of times in light of its straightforwardness. There are distinctive methods of operations of AES: ECB, OFB, CTR, XTS and CBC. It is accounted for that AES: ECB is great decision of encryption or unscrambling calculation since its simultaneously played out a calculation in a disseminated domain [1].

#### 4.1 ENCRYPTION IN HDFS

Figure 2 indicates operation that spare each piece into HDFS, customer split every record into settled size square and scrambles it before transfer to Hadoop document framework. It is accounted for that encryption and unscrambling can be actualized essentially by utilizing Java class [1]. Customers, itself perform encryption utilizing AES calculation on the CPU and exchange encoded piece to HDFS (DataNode). At that point collector DataNode (First DataNode where piece store) reproduce hinder into two different DataNodes.

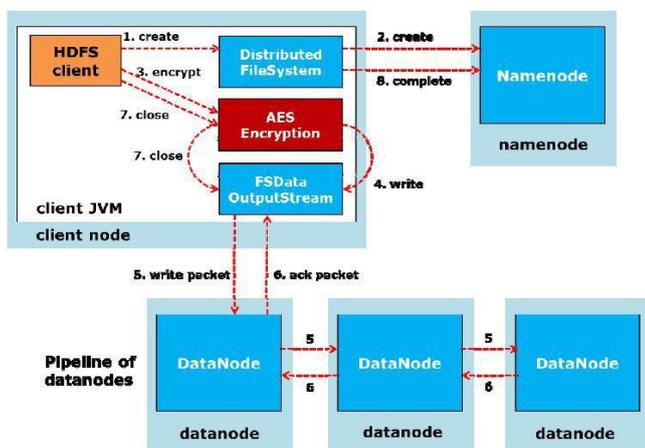


Figure 2: Writing a file by adding an encryption step[1]

#### 4.2 DECRYPTION IN HDFS

Information pieces are composed by customer to DataNode successively, however amid execution of MapReduce occupation numerous squares are perused (decoded) parallel at TaskTracker. Figure 3 demonstrates that MapTask read and encode information obstructs at TaskTracker utilizing AES encryption strategy. It is accounted for that various MapTasks are executing in Hadoop at specialist destinations. HDFS bolsters compose once-read-many model, it is accounted for that simultaneous decoding of HDFS square well reasonable for some MapReduce employments [1].

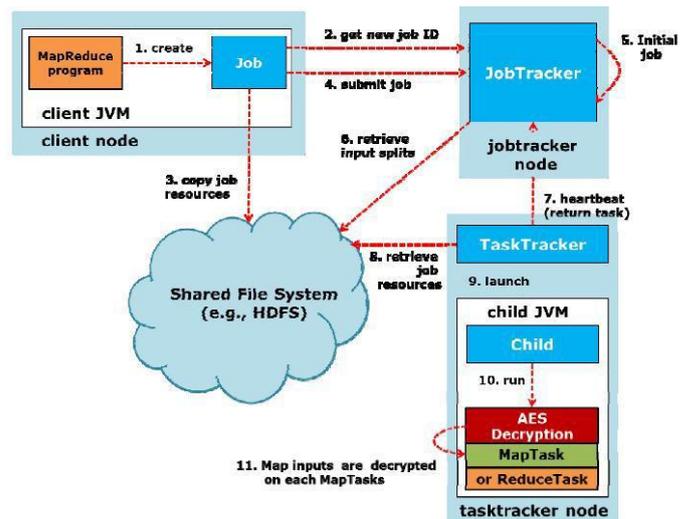


Figure 3: A MapReduce job that read an encrypted file[1]

### 5. Future Scope

Enormous information contains delicate and private data, so as to secure this huge volume that put away at various product equipment, important to actualize confirmation to check client or framework personality. Approval is valuable for giving access control benefits to client or framework; additionally the ACL's are aides for document consent. OAuth 2.0 is great decision for both validation and Authorization. Furthermore, review trails utilized for following every client action. OAuth 2.0 token effective component that bolster AES to give information classification and uprightness among various client

### 6. Conclusion

In the period of Big Data, where information is gathered from various sources, security is a measure issue, as there no any settled wellspring of information and no sort of security instrument. Hadoop received by different businesses to process such information, requests solid security arrangement. Consequently verification, approval and encryption or decoding strategies are much supportive to secure Hadoop record framework.

**REFERENCES**

[1] Seonyoung Park and Youngseok Lee –Secure Hadoop with Encrypted HDFS”

[2] Dean, J., Ghemawat, S.: MapReduce: Simplified Data Processing on Large Cluster. In:OSDI (2004)

[3] Ghemawat, S., Gobiuff, H., Leung, S.: The Google File System. In: ACM Symposium on Operating Systems Principles (October 2003)

[4] O'Malley, O., Zhang, K., Radia, S., Marti, R., Harrell, C.: Hadoop Security Design, Technical Report (October 2009)

[5] White, T.: Hadoop: The Definitive Guide, 1st edn. O'Reilly Media (2009)

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

[7] Jason Cohen and Dr. Subatra Acharya –Towards a Trusted Hadoop Storage Platform: Design Considerations of an AES Based Encryption Scheme with TPM Rooted Key Protections|| (2013)

[8] Lin, H., Seh, S., Tzeng, W., Lin, B.P. || Toward Data Confidentiality via Integrating sfsHybrid Encryption Schemes and Hadoop Distributed FileSystem|| (2012)

[9] Thanh Cuong Nguyen, Wenfeng Shen, Jiwei Jiang and Weimin Xu –A Novel Data Encryption in HDFS|| (2013)

[10] Devaraj Das, Owen O'Malley, Sanjay Radia and Kan Zhang –Adding Security to Apache Hadoop||

[11] Songchang Jin, Shuqiang Yang, Xiang Zhu, and Hong Yin –Design of a Trusted File System Based on Hadoop “ 2013

[12] Advanced Encryption Standard, [http://en.wikipedia.org/wiki/Advanced\\_Encryption\\_Standard](http://en.wikipedia.org/wiki/Advanced_Encryption_Standard)

[13] Sharma Y. ; Kumar S. and Pai R.M; –Formal Verification of OAuth 2.0 Using Alloy Framework ||