

# An approach for Security Information and Event Management with Hadoop MapReduce

Anand Mehta<sup>1</sup> Manish Kumar Abhishek<sup>2</sup>

<sup>1</sup>PG Student <sup>2</sup>Manager (IT Infrastructure)

<sup>1,2</sup>Department of Computer Engineering

<sup>1</sup>GTU PG School, Gujarat, India <sup>2</sup>RailTel Corporation of India Ltd., Haryana, India

**Abstract**— Now a days, technology gives more profits in the municipal sector and private sector as well as the threats and influence of the coercions also high. This is a very problematic to promise a safety in a PC and IT systems because of the swiftly expansion of IT skills and except the Information Technology structure analysis of log is very significant. Infrastructure weaknesses is revealing openly due to lake of safety. This article contains an approach for design SIEM for handle large amount of log data Hadoop prospective is a best, with the help of the HDFS file structure and MapReduce mechanism processing of the logs is faster. So the user in a network operation room can visualize the attack or suspicious activity in real time.

**Key words:** SIEM, Security, Big Data, Hadoop, MapReduce, HDFS, Kibana

## I. INTRODUCTION

The quantity of fake-spams or cyber-bout are intensifying every months in a huge amount. The fabricator of antivirus tools, Kaspersky Lab notifies that its resolution is noticed 23,680,646 in 2008 to 51,887,400,554 in 2013 [1, 2].

As well as with the figures of the report given by the Verizon RISK group in 2012: 55% of malware was distinguished after the long time (month) from infection. Only 59% malware attack [6] was acknowledged in a single day [3].

Now a day's research concerning text based logs to visualizing has been consistently led since past may periods. In this period technologies and its ways constantly upgraded and user friendly and the size of logs has been quickly enlarged expanded through the through the advancement of Information Technology (IT). It is integral to conceptualize the log for efficient examination or mining of item sets [7]. This is too compelling in information or network safety field. Due to the alteration in security equipment and evolution in bulk packing, there is an edge to inspect security logs/data with limited social origin. So, enlargement of Security Information and Event Management (SIEM) [8] that interrogate and visualize a few security logs has guide. The SIEM description is operational prevention and resolution in consequence of the fact that classifying Advanced Persistent Threat (APT) attack [11].

APT is a set of quiet and constant PC hacking procedure, frequently arranged by humanoid aiming a precise unit [4]. In this article, we done the examination of large data sets of real time logs with the help of open source program Hadoop [13] Hadoop- MapReduce is appropriate and applied in many field for Big Data [14] analysis. As Log files [15] is also one of the type of huge data which growing fast so Hadoop is the finest and suitable platform for storing files and parallel execution of MapReduce program for examining them[9][10].

## II. ARCHITECTURE

### A. Security Information and Event Management (SIEM)

SIEM associate two dissimilar field, as per the fig. 1, right side is a Security Event Management (SEM) and left side is a Security Information Management (SIM). This fields are main attention because on the analysis and collecting of security significant data. Conversely, SEM accentuates the collection of log records in to adaptable volume of Information with the aid which Security event may be apportioned with almost though security info management (SIM) basically attentions on investigation of previous data in instruction to expand the extensive term usefulness and proficiency of infrastructure in information security structure. [12]

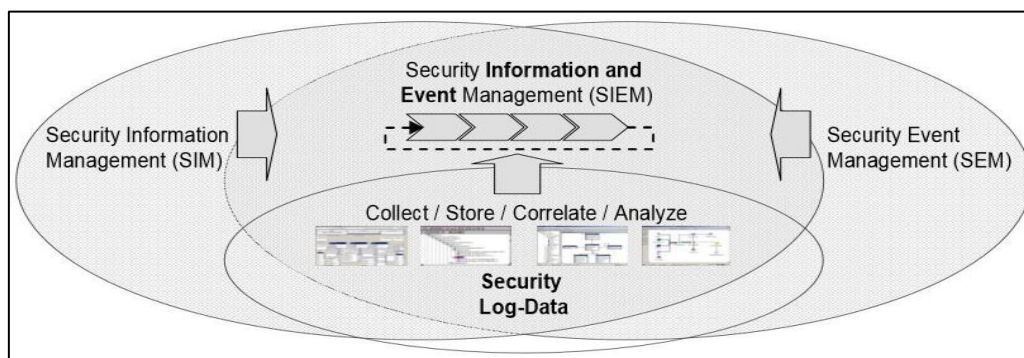


Fig. 1: Introspective Architecture of SIEM

The merger of SEM and SIM into a linked advancement of arrangement, controlling and auditing security applicable information on structure of data gathered from the various Information Security architecture is abbreviate in the term SIEM [5].

### B. Hadoop

Apache Hadoop is an open-source platform, which help in readying data and analogous administration in a circulated environment. Hadoop breach the Big Data base into piece of data and set aside over the association called clusters. For a handling the massive data, MapReduce is available and used for side by side processing on clusters, therefore it cut down the compilation time. The HDFS-Hadoop Distributed File System is basically a spread file system which is planned to perform on hardware. HDFS is tremendously fault-tolerant. HDFS also transports high output access to call data and is extremely fit for applications that have huge data sets.

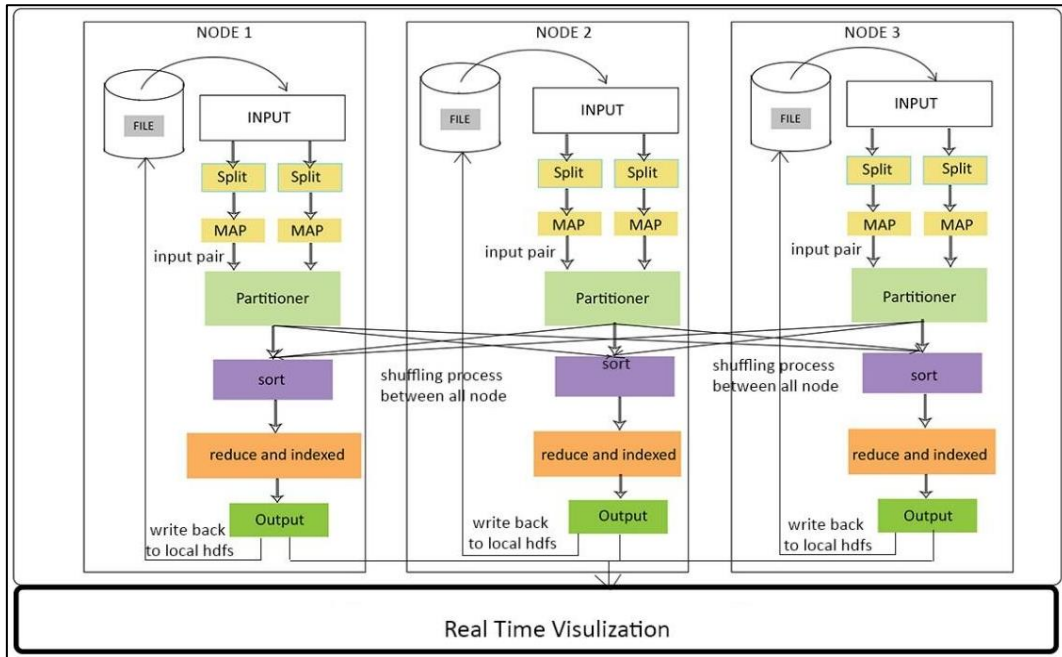


Fig. 2 Data Flow

### III. PROPOSED SOLUTION

In the proposed solution fig. 2 shows the map process and reduce process mechanisms of the request fit in a little nearer detail. The channel with more of its mechanism open. However only three nodes are described, the same channel and block can be simulated through a large number of nodes. A. Input files Input files is an initial log files or actual data for a MapReduce task, Input Files usually reside in Hadoop distributed File System (HDFS). This files format is random mostly the file format is text based. We normally handle multiline record. The log files are usually very large in size. B. Input Format Input format labels the input description for a map reduce job. An Input format is a session that delivers the subsequent functionality like authenticate the input requirement of the submitted task and breakdown the input file into logical Separations like small chunk, every of the chunk is then allocated to a separate mapper. Proposed algorithm:

- Step 1: Fetching entire file as an Input file from the Hadoop distributed file system.
- Step 2: Entire Input file is plain text file, then entire file splits in lines.
- Step 3: Now two input pairs generated called K and V.
- Step 4: Then, pairs K and V values submitted for map process.
- Step 5: Output of process mapped data aggregated in partitioned and the value will be shuffling between other nodes mapped data.
- Step 6: then value will be sorting as per the threshold value  $\leq$ .
- Step 7: then value will be reduce and indexed for visualization. (Indexing giving the benefit to organize the sorted data to use differently in dashboard as per user requirement)
- Step 8: Final output will be write back to the local hdfs and submitted parallel in to kibana visualization.

An Output format that inscribes basic text files is text output format. This is avoidance arrangement which inscribes (key, value) sets on separate positions of a text file. Text output set-up  $\langle K, V \rangle$  which is outspreads File output format  $\langle K, V \rangle$ . Categorization file output format is a middle arrangement amongst MapReduce works which used to serialize subjective data kinds for the file in Job.

#### IV. RESULT

Here we are present our proposed work, for the SIEM system and log analysis we observed live log data in Data center, in this log files it contains multiple field as IP addresses like source IP address, Destination IP address, URL, timestamps. We have installed Hadoop 2.7.2 on Cent OS machine with java 1.8. Log files are circulated consistently on these nodes on the cluster, The MapReduce job goes on these files and get examined consequences in the graphical presentations like graph charts using Kibana visualization.

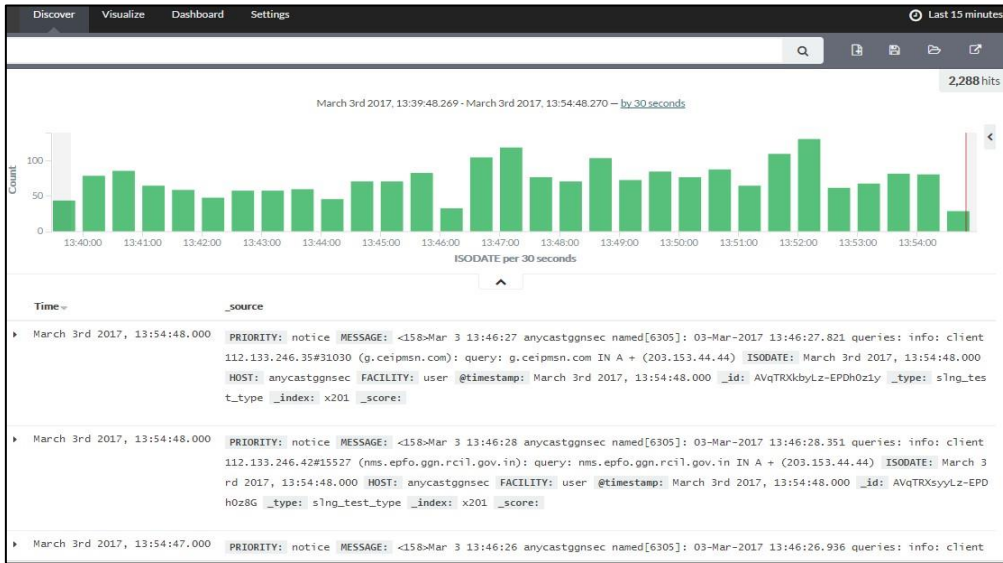


Fig. 2: Hits per 30 Seconds



Fig. 3: Timestamps count per minute

#### V. CONCLUSION

Log investigation supports to develop the business tactics and helps to make statistical reports. Hadoop MapReduce established log file investigation mechanism will provide us graphical reports with the help of Kibana Dashboard showing hits for pages, activity, in part of users are interested in sites, consistence login failures, traffic, attack etc. From these intelligences industry societies can estimate which portions of the site essential to be upgraded on behalf. Using of Hadoop MapReduce framework delivers parallel distributed and steadfast data storing by duplicating data for huge log files. At first stage, files get stored block wise in on a number of nodes in a cluster so that time obligatory can be reduced that save more execution time and give better presentation. Here hadoop main representative of improve response time. And Map reduce positively workings scattered for huge data set and providing the more well- organized results.

#### ACKNOWLEDGEMENT

The author1 is highly thankful to his respected guide Mr. Manish Kumar Abhishek for marvelous guidance and support to complete this paper. Author is also thankful to RailTel Corporation of India Limited (Ministry of Railways undertaking) and thankful to Mr. Srinivas Vasu (ADGM) of Railtel Corporation of India Ltd. For encouraging and support. The author1 would like to thank parents for financial and moral supports throughout their technical education.

#### REFERENCES

[1] Gostev, A. Kaspersky Security Bulletin: Statistics 2008, <https://securelist.com/analysis/kaspersky-security-bulletin36241/kaspersky-security-bulletin-statistics-2008/>

- [2] Funk, C.; Garnaeva, M. Kaspersky Security Bulletin. The Overall Statistics for 2013. Available online: <https://securelist.com/analysis/kaspersky-security-bulletin/58265/kaspersky-security-bulletin-2013-overall-statistics-for-2013/>.
- [3] Verizon RISK Team. Verizone 2012 Data Breach Investigations Report. Available online: <http://www.verizonenterprise.com/DBIR/2016/>. (Visited on Dec 2016)
- [4] Binde, Beth, Russ McRee, and Terrence J. O'Connor. "Assessing outbound traffic to uncover advanced persistent threat." SANS Institute. Whitepaper (2011).
- [5] Tobias Hoppe, Alexander Pastwa, Sebastian Sowa, "Business Intelligence Based Malware Log Data Analysis as an Instrument for Security Information and Event Management" International Journal on Advances in Security, vol 2 no 2&3, year 2009.
- [6] Roland Gabriel, Tobias Hoppe, Alexander Pastwa, Sebastian Sowa, "Analyzing Malware Log Data to Support Security Information and Event Management: Some Research Results" published in IEEE conference, year 2009
- [7] Tushar M. Chaure, Kavita R. Singh, "Frequent Itemset Mining Techniques – A Technical Review" published in IEEE WCFTR year 2016.
- [8] Sandeep Bhatt, Pratyusa K. Manadhata and Loai Zomlot, "The Operational Role of Security Information and Event Management Systems" published in IEEE Computer and Reliability Societies, year 2014.
- [9] Damian Hermanowski, "Open Source Security Information Management System Supporting IT Security Audit" published in IEEE, year 2015
- [10] Igor Anastasov, Danco Davcev, "SIEM Implementation for Global and Distributed Environments" published in IEEE year 2014.
- [11] Jaehee Lee, Changyeob Lee, Jaebin Cho, "A Study on Efficient Log Visualization Using D3 Component against APT How to visualize security logs efficiently?" published in IEEE year 2016.
- [12] Anand Mehta, Manish Kumar Abhishek, "A survey on log corelation in security information and event management with hadoop" published in IJARIE Vol-2 Issue-6 2016.
- [13] Apache Hadoop, <http://hadoop.apache.org/> (visited on 21 Dec 2016)
- [14] Big Data, [https://www.sas.com/en\\_us/insights/big-data/what-is-big-data.html](https://www.sas.com/en_us/insights/big-data/what-is-big-data.html) (visited on 19 Dec 2016)
- [15] Vangie Beal, "Log Files" by, [http://www.webopedia.com/TERM/L/log\\_file.html](http://www.webopedia.com/TERM/L/log_file.html) (visited on 15 Oct 2016)