



## Internet Traffic Flow Analysis using Hadoop

By Sudipa Biswas & Subhradeep Biswas

**Abstract-** The internet traffic analysis elucidates the network administrator for monitoring the ongoing operation in the network and to understand the network so that the behavior could be examined and large problem can be examined. Flow analysis assists in traffic management, allocation of resources and fault tolerance. Due to the fast increase in internet user simultaneously the network usage has also escalated rapidly. The major problem of this fast growth in network is the traffic management, storing of traffic data and analysis this enormous amount of data in a single machine. To resolve this issue hadoop has been implemented to scan multiple input data and produce output for traffic identification and clustering flow. In this paper internet traffic flow analysis has been done using hadoop. In this proposed method system accepts packet data as input from network and this input is appended to hadoop distributed file system (HDFS) and at last processing is done through MapReduce. Once the output has been generated the network administrator analyses the internet traffic and troubleshoot any problem if necessary.

**Keywords:** HDFS, traffic analysis, traffic identification, traffic clustering, mapreduce, and hadoop framework.

**GJCST-B Classification:** C.2.5



*Strictly as per the compliance and regulations of:*



# Internet Traffic Flow Analysis using Hadoop

Sudipa Biswas <sup>α</sup> & Subhradeep Biswas <sup>σ</sup>

**Abstract-** The internet traffic analysis elucidates the network administrator for monitoring the ongoing operation in the network and to understand the network so that the behavior could be examined and large problem can be examined. Flow analysis assists in traffic management, allocation of resources and fault tolerance. Due to the fast increase in internet user simultaneously the network usage has also escalated rapidly. The major problem of this fast growth in network is the traffic management, storing of traffic data and analysis this enormous amount of data in a single machine. To resolve this issue hadoop has been implemented to scan multiple input data and produce output for traffic identification and clustering flow. In this paper internet traffic flow analysis has been done using hadoop. In this proposed method system accepts packet data as input from network and this input is appended to hadoop distributed file system (HDFS) and at last processing is done through MapReduce. Once the output has been generated the network administrator analyses the internet traffic and troubleshoot any problem if necessary.

**Keywords:** HDFS, traffic analysis, traffic identification, traffic clustering, mapreduce, and hadoop framework.

## I. INTRODUCTION

Internet is inclusive system which connects numbers of computer with each other. It uses TCP/IP to get connect this devices which contains packets getting from source to destination computer. This computer network is usually administered by [1] software defined network (SDN). This assist the network in by decoupling that drive the outcome about the traffic is remitted. To implement the traffic analysis of the big data the collection needs to be done in order to measure the data for Varsity sources. The huge amount of data which can be of any form like 3d data, audio, video, text and many more which cannot be processed by any traditional way but by the big data approach we can measure and analyze the data to be further analyzed for resolving network related issues. Hadoop[2] is implemented that uses basic programming to process large amount of data sets. The main intention of this paper is to design and implement a system to for network traffic analyze utilizing hadoop clusters [3]. Once the input is given the detailed measuring and analyze is done on the input and output is derived on the basis of the input given to the system.

**Author α:** Department of software engineering, BITS pilani, India.

**e-mail:** sudipabiswascs@gmail.com

**Author σ:** TCS, India.

## II. PROPOSED SYSTEM

### a) Hadoop Framework

Hadoop is literally open source which is java based programming that will support analyzing processing and storage for large data set in a computing environment. Hadoop makes possible on the application on system with thousands of commodity hardware nodes and handling terabytes of data. Its file system which is distributed and facilitates rapid data transfer rates among nodes and continue operating a nodes failure. Hadoop has emerged as big data foundation and scientific analytics, business and planning.

### b) System Description

The process of flow analysis consists of mainly three main factors. Firstly the data exchange, secondly the analysis and thirdly the user interface. In the data exchange process the HDFS [4] is implemented to store the information or the data to be used as an input in analyzing. In the analysis process the data is analyzed and managed and other factors such as node, link, and flow analysis are also implemented in this process. In the user interface the user interface the system and graphical display is displayed to enhance the user understands the analysis flow. The API [5] and GUI [6] tool is also implemented which enhances the communication. The network gives the input to the system. Below in Fig1 component analysis flow architecture is given.

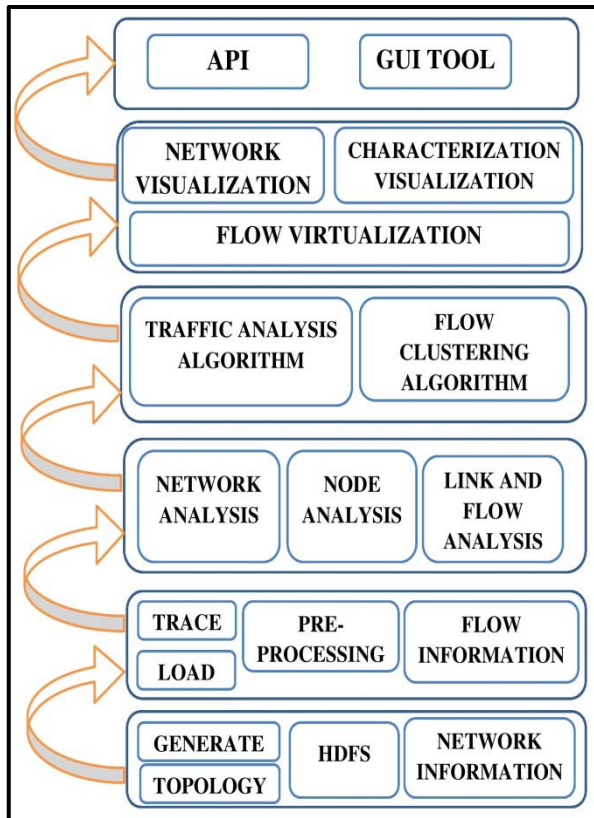


Fig. 1: Component analysis flow architecture

The input is actually data packet following from network. All the information is actually stored in HBASE. After that the system executes HDFS and MapReduce[7] function on the input data packet. After the function gets implemented and packets are stored the flow mechanism is implemented. The statistical view of the data can be viewed by the user. The different platform like Windows, Linux, and Mac OS/X are all on written in java platform. The specified format of the input helps in parsing of the input data. The sorting of the input file becomes necessary if the required format is done available. Parsing of the input is based on source IP address [8], types of packet, destination IP address, source port address and destination port address. The input which has same source port and destination port will form and stored together and input from the same source IP address and same destination IP address is clustered and stored as unstructured data in the database. In Fig2 block diagram of flow system analysis.

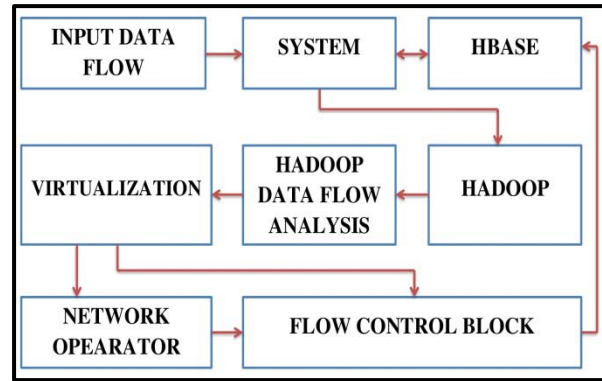


Fig. 2: Block diagram of flow system analysis

### III. CONCLUSION

In this paper we have done the flow analysis and identification on hadoop platform. We have provided detailed analyses of the input data packet and classification of the type. This paper shows the methodology for analyzing packet file and statistical analysis of the original input packet and flow. The future work of this paper can be worked on the various problems like more networks makes more congestion troubleshooting the problem using the hadoop technology.

### REFERENCES RÉFÉRENCES REFERENCIAS

1. Ta-Yuan Chou; Wun-Yuan Huang; Hui-Lan Lee; Te-Lung Liu; Joe Mambretti; Jim Hao Chen; FeiYeh "Heterogeneous Interconnection between SDN and Layer2 Networks Based on NSI" : 2017 31st International Conference on Advanced Information Networking and Applications Workshops (WAINA).
2. Sunita Choudhary; Preeti Narooka "Hugepage & Swappiness functions for optimization of the search graph algorithm using Hadoop framework": 2016 IEEE International Conference on Computational Intelligence and Computing Research (ICIC).
3. Yuanqi Chen; Yi Zhou; Shubbhi Taneja; Xiao Qin; Jianzhong Huang "aHDFS: An Erasure-Coded Data Archival System for Hadoop Clusters": IEEE Transactions on Parallel and Distributed Systems.
4. Youngho Song; Young-Sung Shin; Miyoung Jang; Jae-Woo Chang "Design and implementation of HDFS data encryption scheme using ARIA algorithm on Hadoop" : 2017 IEEE International Conference on Big Data and Smart Computing (BigComp).
5. Yang Lai; Shi Zhong Zhi "An Efficient Data Mining Framework on Hadoop using Java Persistence API": 2010 10th IEEE International Conference on Computer and Information Technology.
6. Nikolay Laptev; Kai Zeng; Carlo Zaniolo "Very fast estimation for result and accuracy of big data analytics: The EARL system" : 2013 IEEE 29th

International Conference on Data Engineering (ICDE)

7. Yasser Altowim; Sharad Mehrotra "Parallel Progressive Approach to Entity Resolution Using MapReduce": 2017 IEEE 33rd International Conference on Data Engineering (ICDE).
8. Elias Bou-Harb; Mourad Debbabi; Chadi Assi "BigData Behavioral Analytics Meet Graph Theory: On Effective Botnet Takedowns" : IEEE Network.
9. <http://searchcloudcomputing.techtarget.com/definition/Hadoop>





This page is intentionally left blank