

Analysis of Network Traffic Data Using Hadoop and Apache Hive

Swathi Prabhu^[1], Pooja Parashuram Bajantri^[2], Prajna B Nayak^[3], Priya Shetty^[4], Rachana N^[5]

Abstract— Computer network traffic can be classified into various classifications based on parameters such as port number, IP address or protocol. The sequence and pattern of the traffic in the network is illuminated using flow analysis. This helps network administrator to understand the network usage and to examine the behavior of the user using the network and to monitor the operations going on in the network. Flow analysis helps in fault tolerance, resource allocation and network capacity planning. Since this is data age, due to the fast growing network day by day the volume of the traffic is getting huge. So it is very difficult to collect, store and analyze this large data on a single machine. It requires scalable tools to measure, classify and analyze this traffic data. Due to the limited computational capability and storage capacity traditional tools are failed. Hadoop is a leading distributed framework which is designed to execute tremendous data that can be of hundreds of terabytes and even petabytes of data. In this paper a Hadoop framework based network traffic data analysis is done. Here the input is appended to the Hadoop Distributed File System (HDFS) and process the data using Map-Reduce, Hive. The output is displayed using Tableau and RStudio.

Index Terms— Hadoop, HDFS, Hive, Internet Traffic, Map-Reduce, Network flow analysis.

1 INTRODUCTION

A network consists of two or more computers that are linked in order to share resources, exchange files or allow electronic communications. The computers on a network may be linked through cables, telephone lines, radio waves, satellites or infrared light beams.

Data analytics has become a key element of the business decision process over the last decade. Classic systems based on relational databases and expensive hardware, while still useful for some applications, are increasingly unattractive compared to the scalability, economics, processing power, availability offered by today's network driven distributed solutions.

Network traffic analysis [1] is nothing but analyzing the amount of data moving across a network at a given point of time. Network packets which are mostly the network data are the loads in the network. The main objective of monitoring network traffic is to check the availability and smooth operations of the network. Network traffic requires reviewing of each incoming and outgoing packets. It can be classified using HTTP, IP, ICMP, TCP and UDP. To address the challenges in the traditional systems we use Apache Hadoop. It is an open source software framework for storage and large scale processing of datasets.

1. Scalability: Ability to store data in a distributed fashion in a large set of servers
2. Cost-effective Solution: Hadoop's scale out architecture with Map-Reduce programming model, allows the storage and processing of BigData in affordable manner.
3. Flexibility: The Map-Reduce model of Hadoop Framework can deal with both structured and unstructured datasets.
4. Fast: Parallel Processing allows multiple processes to divide the task and hence take fewer-time.

The data analysis is carried out through Hadoop Map-Reduce and Apache Hive.

Map-Reduce[7] is a programming model used in the Hadoop ecosystem to process large amount of data in a distributed manner. Map-Reduce programs are inherently parallel, thus able to perform very large-scale data analysis. Although it is a powerful tool for processing, it is harder to understand and to develop. Hence we can use Apache Hive, a project developed on top of Hadoop framework which is easy to understand and develop.

Apache Hive is the SQL-On-Hadoop technology to query heterogeneous data stored in different databases and filesystems that are fused with Hadoop. This approach is significantly faster and has new features that will support performing inserts and updates to tables. To visualize the data, we made use of Tableau Public and RStudio.

Tableau public is an open source tool for pattern discovery using data visualization. Learning and using Tableau is a very low time consuming activity. But there are limitations, such as we can only read from txt, Excel or Odata sources and there is no security for the data. Hence we used RStudio for analysis.

RStudio provides enhanced security and authentication, advanced resource management and priority email support.

[1] Assistant Professor, Dept. of CSE, SMVITM, Bantagal.
prabhuswathi2@gmail.com

[2][3][4][5] Student, Dept. of CSE, SMVITM, Bantagal.
poojauapadhya96@gmail.com [2], nayak.prajna97@gmail.com [3],
shettypriya495@gmail.com [4], rachanasmoitm@gmail.com [5]

2 ARCHITECTURE

2.1 Apache Hadoop

Apache Hadoop 1.0 and apache Hadoop 2.0 are the two versions of Hadoop. The 1.0 version was the initial release, where there were only two modules: Map-Reduce and HDFS (Hadoop Distributed File System). The Map-Reduce was the programming model used for large scale processing and was also taking care of cluster resource management. The HDFS is used to manage the storage across the network of machines.

Apache Hadoop 2.0 represents a generational shift of architecture in the apache Hadoop. With the introduction of YARN (Yet another Resource Negotiator), Apache Hadoop became a significantly strong platform. YARN module takes care of cluster resource management in Hadoop 2.0, which allows the multiple applications to run in Hadoop. Only the data processing is done by Map-Reduce module. The Fig 1. depicts the Hadoop Architecture.

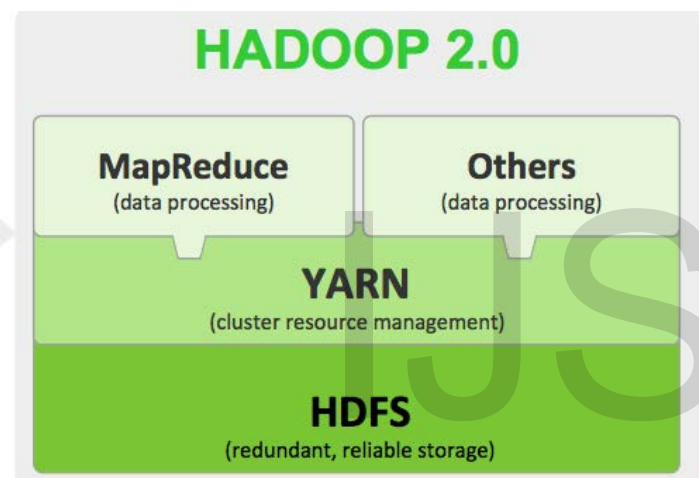


Fig 1. Hadoop Architecture

2.2 Map-Reduce

Google[3] was the first to propose the Map-Reduce programming model for page ranking or web log analysis. Map-Reduce is a programming model and an associated implementation for processing and generating big datasets with the parallel, distributed algorithm on a cluster. A Map-Reduce program is composed of two phases: a Map phase and a Reduce phase. Each phase has key-value[7] pairs as input and output.

Map phase: Master node takes large problem input and slices it into smaller sub problems and attributes these to worker nodes.

Reduce Phase: Master node collects the answers to the sub problems from the worker nodes and combines them in a pre-defined way to get the output.

2.2 Apache Hive

Apache Hive is a data warehouse software project built on top of Apache Hadoop for providing data summarization, query and analysis. Initially Hive was developed by Facebook, later the Apache software foundation took it up and further developed as an opensource under the name Apache Hive. It is a

platform used to develop SQL like scripts to do Map-Reduce operations.

Some of the features are as follows:

- It stores schema in the database and process the data into HDFS.
- It provides SQL like query language called HiveQL or HQL.
- It is designed for OLAP.

2.3 Comparison between Map-Reduce and Hive

The data processing is done using both Hadoop's Map-Reduce programming model and Apache Hive. The comparison between the two programming model is given in Table 1.

TABLE 1
COMPARISON BETWEEN MAP-REDUCE AND HIVE

Hadoop Map-Reduce	Apache Hive
Compiled language	SQL-Like query language
Low level of abstraction	Higher level of abstraction
More lines of code	Lesser lines of code
More development effort is required	Development effort is less
Code efficiency is high	Code efficiency is less

2.4 Tableau Public

Tableau is a Business Intelligence tool for visually analyzing the data. Users can create and distribute an interactive and shareable dashboard, which depicts the Trends, Variations and Density of the data in the form of Graphs and Charts. Tableau can connect to files, relational and Big Data sources to acquire and process the data. The software allows data blending and real time collaborations which makes it unique. It is used by businesses, academic researchers and many government organizations for visualize data analysis. Tableau provides solutions for all kinds of industries, departments and data environments. Speed of analysis, Self-reliant, Visual discovery, Blend diverse data sets, Real-time collaborations and Centralized data are some of the features of the Tableau Public.

2.5 RStudio

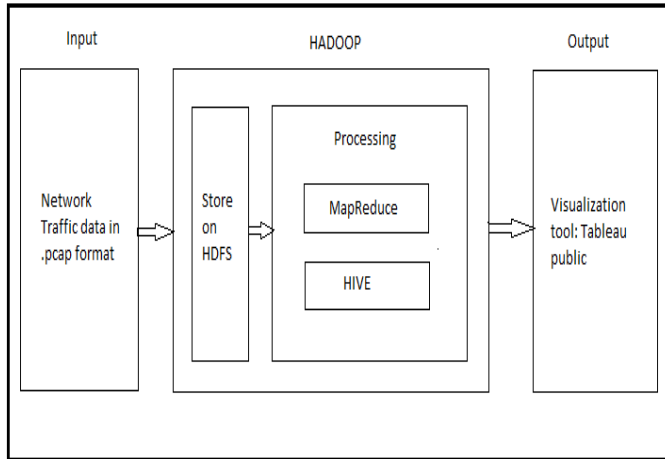
RStudio is a free and open-source integrated development environment for R, a programming language for statistical computing and graphics. RStudio was founded by JJ Allaire. RStudio integrates with R as an IDE (Integrated Development Environment) to facilitate the user with further functionality. RStudio is a development environment facilitates with source code editor, build automation tools and the debugger. Interactive, Easy debugging, Friendly environment for package development are some of the advantages of using RStudio.

3 SYSTEM OVERVIEW

The proposed system involves four main phases: Collection

and Conversion of input data set, storing in HDFS, processing of data using Map-Reduce and Hive, visualization using Tableau Public and RStudio. The working procedure of the system is represented by Fig 2.

Fig. 2 Proposed System Architecture



Four phases are as follows:

A. First Phase:

- a) Collection of dataset: We are using pcap file as a dataset, collected from [6]. And also we used Wireshark tool to capture the packets.
- b) Conversion of data set from .pcap format to .csv file format: Wireshark tool is used for conversion of file format.

B. Second Phase:

- a) Setting up of Hadoop platform: Installed Hadoop on Ubuntu 16.04 LTS version.
- b) Load dataset into HDFS

C. Third Phase:

- a) Processing by Map-Reduce: We used Map-Reduce programming model in which we defined separate mapper and reducer class to process the pcap data.
- b) Processing by Hive: Initially we loaded the data into External hive table and then the data queried using HiveQL, a SQL like querying language.

D. Fourth Phase:

Analysis using visualization tool: Tableau Public, an open source visualization tool is used to display the analyzed result. We have also made use of RStudio for analysis.

Hadoop: we used Hadoop 2.7.4 version with Multi-node Cluster containing one master and a slave.

Apache Hive: we used Hive 2.2.0 version which runs on top of Hadoop Distributed File System (HDFS).

Tableau Public: A visualization tool of version 10.5 to display results graphically.

RStudio: A statistics Computational tool of version 1.1 to display the results.

4.2 Overall Results

1. PCAP dataset

The dataset employed here is basically a pcap file which contains fields such as time of the packet capture, source IP address, destination IP address, protocol, length and other information about the packet. A screenshot of dataset is shown in Fig 3.

```

@pcap-Network
File Edit Format View Help
*No. *Time *Source *Destination *Protocol *Length *Info
*1 *0.000000 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*2 *0.033094 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*3 *0.033350 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*4 *0.066439 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*5 *0.066679 *130.38.180.250 *130.38.180.254 *CIGI *60 *130.38.180.250 => 130.38.180.254 (16 bytes)
*6 *0.069735 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*7 *0.100061 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*8 *0.133016 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*9 *0.133279 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*10 *0.369418 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*11 *0.400131 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*12 *0.432785 *130.38.180.254 *130.38.180.250 *CIGI *186 *130.38.180.254 => 130.38.180.250 (144 bytes)
*13 *0.433069 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*14 *0.466020 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*15 *0.466330 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*16 *0.499323 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*17 *0.499629 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*18 *0.532827 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*19 *0.532857 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*20 *0.569936 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*21 *0.566255 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*22 *0.569234 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*23 *0.599550 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*24 *0.632537 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
*25 *0.632783 *130.38.180.250 *130.38.180.254 *CIGI *130 *130.38.180.250 => 130.38.180.254 (88 bytes)
*26 *0.665840 *130.38.180.254 *130.38.180.250 *CIGI *122 *130.38.180.254 => 130.38.180.250 (80 bytes)
    
```

Example Dataset

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
rw-r--r--	hadoop	supergroup	54 B	20/02/2018 9:28:52 AM	1	128 MB	inputs
rw-r--r--	hadoop	supergroup	61 B	20/02/2018 9:24:04 AM	2	128 MB	inputs
rw-r--r--	hadoop	supergroup	1.51 MB	20/02/2018 9:25:19 AM	2	128 MB	input2
drwxr-xr-x	hadoop	supergroup	0 B	18/02/2018 9:31:34 AM	0	0 B	output
drwxr-xr-x	hadoop	supergroup	0 B	20/02/2018 9:28:30 AM	0	0 B	output3
drwxr-xr-x	hadoop	supergroup	0 B	20/02/2018 9:28:31 AM	0	0 B	output2
rw-r--r--	hadoop	supergroup	702.59 KB	20/02/2018 9:32:08 AM	2	128 MB	p1
rw-r--r--	hadoop	supergroup	163.5 KB	20/02/2018 9:32:36 AM	2	128 MB	p2
rw-r--r--	hadoop	supergroup	513.54 KB	20/02/2018 9:32:53 AM	2	128 MB	p3
rw-r--r--	hadoop	supergroup	2.21 KB	20/02/2018 9:33:00 AM	2	128 MB	p4
rw-r--r--	hadoop	supergroup	43.69 MB	20/02/2018 9:33:11 AM	2	128 MB	p5
rw-r--r--	hadoop	supergroup	1.31 KB	20/02/2018 9:33:25 AM	2	128 MB	p6
rw-r--r--	hadoop	supergroup	157.12 MB	20/02/2018 9:33:47 AM	1	128 MB	p7
rw-r--r--	hadoop	supergroup	280.5 MB	20/02/2018 9:34:29 AM	2	128 MB	p8
rw-r--r--	hadoop	supergroup	73.43 KB	20/02/2018 9:34:37 AM	2	128 MB	p9

Fig 4. File Storage in HDFS

2. Copying the Network Traffic data on HDFS is shown in the Fig 4.

3. Processing the dataset using Map-Reduce

4 EXPERIMENTS

4.1 Experimental Setup

Pre-requisites:

- RAM: Minimum 4GB
- Operating System: Ubuntu 16.04 LTS
- Java Version: JDK 1.8

Priyanka Lokhande, "Network Traffic Analysis Measurement and Classification using Hadoop", IJARCCCE, vol. 5, Issue 3, March 2016.

- [2] Rakshitha Kiran P, "Hadoop Technology for Flow Analysis of the Internet Traffic", International Journal of Innovative Research in Computer and Communication Engineering, vol. 3, Issue 2, February 2015.
- [3] Youngseok Lee, Wonchul Kang, Hyeongu Son, "An Internet Traffic Analysis method with MapReduce", IEEE, 2010.
- [4] Hadoop, <http://hadoop.apache.org>.
- [5] Wireshark, <http://wireshark.org>.
- [6] The training dataset collection - <http://www.netresec.com/>.
- [7] Tom White, Hadoop-The Definitive Guide, 4th Edition, 2015.

IJSER