# Empowering Service Claims Fraudulent Using Big Data Analytics

[1] U. Kartheek Chandra Patnaik,[2]Ch. Navyatha, [3]B.V.Maliney, [4]B.Mounica, [5]J. Neelima

[1] *Associate Professor,* [2,3,4,5] *Students*

[1,2,3,4,5] *Department of Computer Science and Engineering.*

[1,2,3,4,5] *Lendi Institute of Engineering and Technology.*

[1,2,3,4,5] *Vizianagaram, Andhra Pradesh, India*

**Abstract-As more and more software moves to Data Analytics as a Service (DAaaS), the web application has become more ubiquitous and log file analysis is becoming a necessary task for analyzing the client's behaviour. Log files are getting generated very fast i.e., at the rate of 1-10 Mb/s per server. A single data centre can generate tens of terabytes of log data in a day which is very huge. In order to analyze such large datasets we need parallel processing system and reliable data storage mechanism. Virtual database system is an effective solution for integrating the data, but it becomes inefficient for large datasets. As log files are continuous stream data from distributed servers, an efficient way to handle such data is needed to store and send data into Storage servers. Our system uses Apache Flume to gather streams of log data from various servers and store them into HDFS, a Hadoop Distributed File System. MapReduce a parallel processing strategy breaks up input data and sends fractions of the original data to several machines in Hadoop cluster. This mechanism helps to process huge amounts of log data in parallel, using all the machines in the Hadoop cluster and computes result efficiently. This approach reduces the computation and response time as well as the load on to the end system. This paper proposes a server log analysis using Apache Flume for continuous streaming, Apache Hadoop through MapReduce and Pig for analyzing such logs, thereby providing accurate results for clients.**

**Keywords -DAaaS, Server Logs, Apache Flume, HDFS, MapReduce, Apache Pig**

## INTRODUCTION

Big Data and Cloud, two of the trends that are defining the emerging Enterprise Computing, show a lot of potential for a new era of combined applications. The provision of Big Data analytical capabilities using cloud delivery models could ease adoption for many companies, and in addition to important cost savings, it could simplify useful insights that could provide them with different kinds of competitive advantage.

Along these lines, Data Analytics as a Service (DAaaS) represents the approach to an extensible platform that can provide cloud-based analytical capabilities over a variety of industries and use cases. From a functional perspective, the platform covers the end-to-end capabilities of an analytical solution, from data acquisition to end-user visualization, reporting and interaction. Beyond this traditional functionality, it extends the usual approach with innovative concepts, like Analytical Apps and a related Analytical Appstore. In addition, the platform supports the needs of the different users who interact with it, including those of the emerging 'Data Scientist' role.

Analyzing is foremost scenario in the trends of DAaaS to make end user satisfaction. One such data analytics is Log File Analysis. Server logs are computer-generated log files that capture network and server operations data. They are useful for managing network operations, especially for security and regulatory compliance. Apache Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log streaming data into the Hadoop Distributed File System (HDFS). It has a simple and flexible architecture based on streaming data flows and is robust and fault tolerant with tuneable reliability mechanisms for failover and recovery.

Apache Hadoop is an open source software framework for storing and large scale processing of data-sets on clusters of commodity hardware. Hadoop is an Apache top-level project being built and used by a global community of contributors and users. It is licensed under the Apache License 2.0. All the modules in Hadoop are designed with a fundamental assumption that hardware failures (of individual machines or racks of machines) are common and thus should be automatically handled in software by the framework. Apache Hadoop's MapReduce and HDFS components originally derived respectively from Google's MapReduce and Google File System (GFS) papers.

## METHODS

A data is a collection of facts from the grids of web servers usually of unorganized form in the digital universe. Around 90% of data present in today's world are generated in last two years [1]. A large amount of the data available in the internet is generated either by individuals, groups or by the organization over a particular period of time. The volume of data becomes larger day by day as the usage of World Wide Web makes an interdisciplinary part of human activities. Rise of these data leads to a new technology such as big data that acts as a tool to process, manipulate and manage very large dataset along with the storage required.

Big Data is a high volume, high velocity and high variety information assets that demand cost-effective [2], innovative forums of information processing for enhanced insight and decision making. Big data, a buzz word in the business intelligence can handle petabytes or terabytes of data in a reasonable amount of time. Big data is distinct from large existing database which uses Hadoop framework for data intensive distributed applications.

Big Data analytics applies advanced analytical techniques of large datasets to discover hidden patterns and other useful information. It is performed using software tools mainly for predictive analysis and data mining. The growing number of technologies is used to aggregate, manipulate, manage and analyze big data. Some of the most prominent technologies are [3] NoSQL databases that include Cassandra, MongoDb, redis, Hbase and Hadoop framework includes Hadoop, HDFS, Hive, Pig.

Numerous research works are carried out in web log mining, hadoop and some of them are reviewed below. Murat et al., [4] proposed the smart miner framework that extracts the user behavior from web logs. The framework used the smart session construction to trace the frequent user access paths. Sayalee Narkhede et al., [5] introduced the Hadoop-MR log file analysis tool that provides a statistical report on total hits of a web page, user activity, traffic sources. This work was performed in two machines with three instances of hadoop by distributing the log files evenly to all nodes.

Milind Bhandare et al., [6] put forth a generic log analyzer framework for different kinds of log files such as a database or file system. The work was implemented as a distributed query processing to minimize the response time for the users which can be extendable for some format of logs. Parallelization of Genetic Algorithm (PGA) was suggested by Kanchan Sharadchandra Rahate et al., [7]. PGA uses OlexGA package for classifying the document. The train model data is stored in HDFS and the test model categories the text document.

A framework for unstructured data analysis was proposed by Das et al., [8] using big data of public tweets from twitter. The tweets are stored in Hbase using Hadoop cluster through Rest Calls and text mining algorithms are processed for data analysis. The semi structured log files are large datasets which are challenging to store, search, share, visualize and analyze. Almost 26% of web log types of data require big data technology to perform an analysis [9]. In order to improve the usage of a website and to track the user behavior, in online advertising and E-commerce the web log mining is performed using Hadoop.

The related works so far stated above performs the work with good scalability but fails to experiment the time efficiency between the different modes of hadoop and necessity of the scalability. The proposed work analyses the working of both hadoop modes and the time efficiency in each mode specifically for semi structure log data, along which a statistical report is made.

*Apache Hadoop* is an open source framework for distributed storage. It process large amounts of datasets on a commodity hardware. Hadoop enables is to take instant decisions in businesses from massive amounts of structured and unstructured data.

*Hadoop Distributed File System (HDFS)* is the core technology for the efficient large storage layer, and is also designed to run across low-cost commodity hardware. It stores huge amounts of datasets in it and provides the datasets for processing when required.

*Apache Flume* is a reliable, distributed and available service for efficiently collecting, aggregating, and moving large amounts of streaming data into the Hadoop Distributed File System (HDFS). It has a simple and flexible architecture for streaming data flows. It is also robust and fault tolerant with reliability mechanisms for failover and recovery.
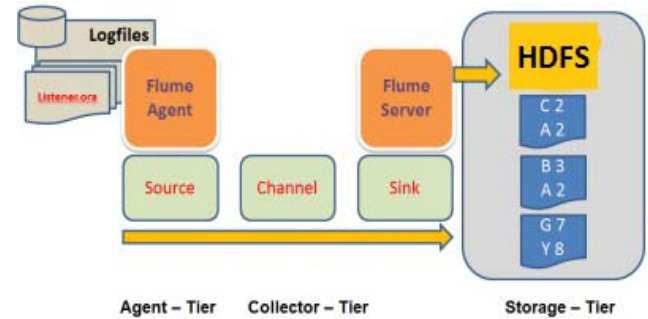


**Fig1. Flume Architecture**

*MapReduce* is a framework used for writing applications that process large amounts of structured and unstructured data in parallel across a cluster of hundreds of machines, in a fault-tolerant and reliable manner. *MapReduce* splits a large dataset into independent chunks of data and organizes that data in key, value pairs for parallel processing. This parallel processing of data in clusters improves the speed and reliability of the cluster and returns the solutions more quickly and with more reliability.

The *Map* function divides the input data into chunks by the Input Format and creates a map task for each chunk in the input. The JobTracker distributes tasks to the worker(slave) nodes. The output of each map task is taken and partitioned into a group of key-value pairs for each reduce.

The *Reduce* function then collects the various results from worker nodes and combines them for an efficient answer for the problem, which the master node was trying to solve. Each reduce pulls the appropriate partition from the machines where the maps have executed, and then writes its result back into HDFS. Thus, the reduce is able to collect the data from all of the maps for the keys and is responsible for combining the results in order to solve the problem.
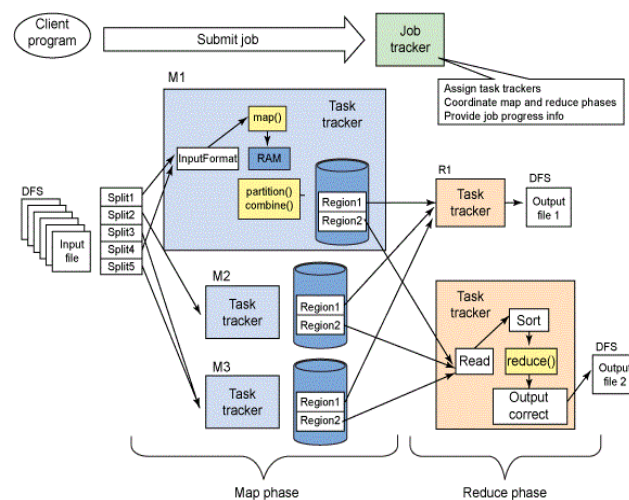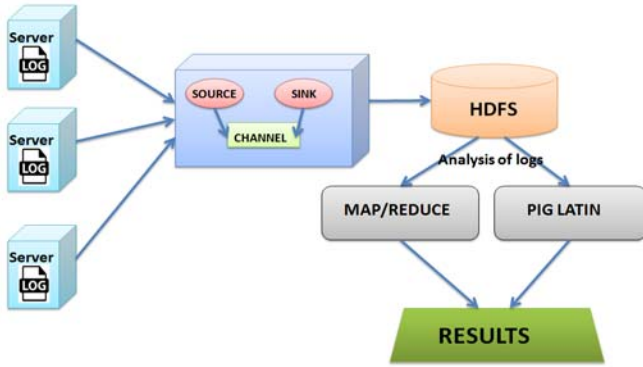


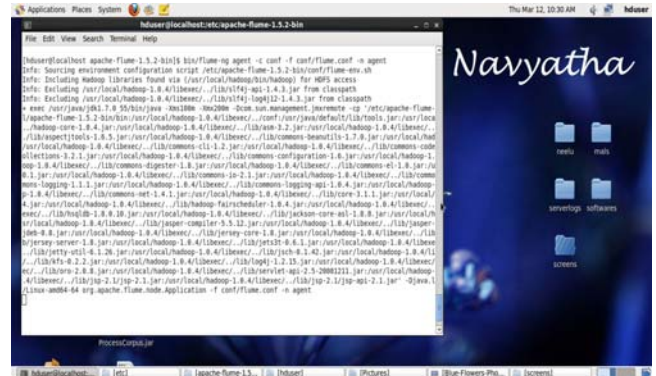**Fig2. MapReduce Architecture**

Fig 3.System Architecture



Fig 4.Apache Flume Agent

*Apache Pig* allows us to write complex MapReduce transformations using a very simple and easy scripting language. Pig Latin (the language) defines a set of transformations on a data set such as aggregate, join and sort. Pig translates the Pig Latin script into MapReduce so that it can be executed within Hadoop®. Pig Latin is sometimes extended using UDFs (User Defined Functions), which the user can write in Java or a scripting language and then call directly from the Pig Latin.

Pig runs on Hadoop and makes use of MapReduce and the Hadoop Distributed File System (HDFS). The language for the platform is called Pig Latin, which abstracts from the Java MapReduce idiom into a form similar to SQL. Pig Latin is a flow language whereas SQL is a declarative language. SQL is great for asking a question of your data, while Pig Latin allows you to write a data flow that describes how your data will be transformed. Since Pig Latin scripts can be graphs (instead of requiring a single output) it is possible to build complex data flows involving multiple inputs, transforms, and outputs. Users can extend Pig Latin by writing their own functions, using Java, Python, Ruby, or other scripting languages.

The user can run Pig in two modes:
- Local Mode. With access to a single machine, all files are installed and run using a local host and file system.
- MapReduce Mode. This is the default mode, which requires access to a Hadoop cluster.

The user can run Pig in either mode using the "pig" command or the "java" command.

## RESULTS & DISCUSSIONS

The data from multiple servers leads to Big Data which is meant as semi-structured, and the need to analyze such data is daunting task. We analyzed the log data using both Apache Pig and MapReduce, there by comparing the performance of analyzing and time take to process the map and reduce functions, showed in Fig.10.

Apache Flume, a framework which is used by the Web server to collect different server log data by different agents, there by collected by collector nodes and finally stored in HDFS Sink, as represented in Fig. 4 and Fig.5
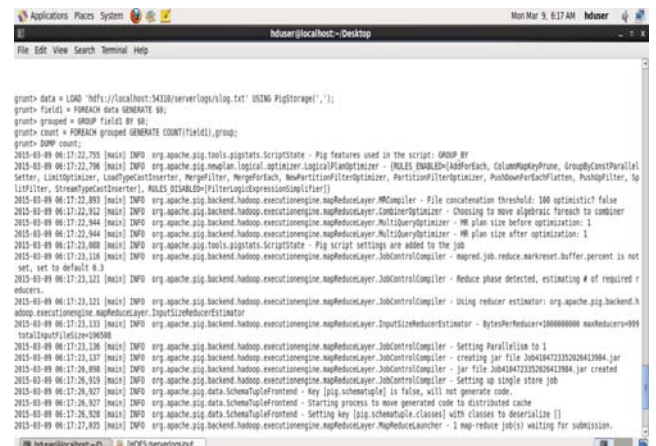


Fig 5.Data stored in HDFS Sink



Fig 6.Pig Analysis in Cent OS
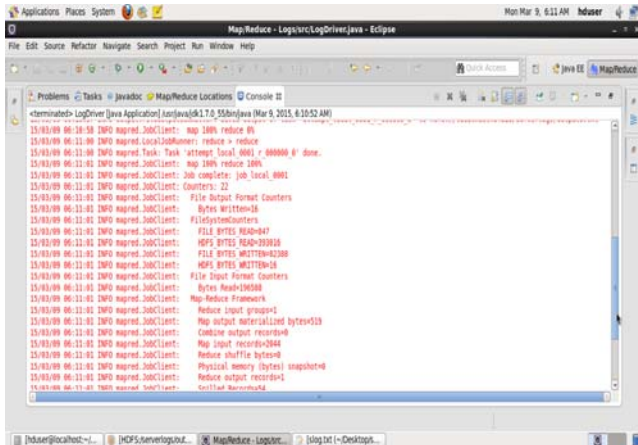


Fig 7.IPAddress Retrieval in Hadoop Using Pig

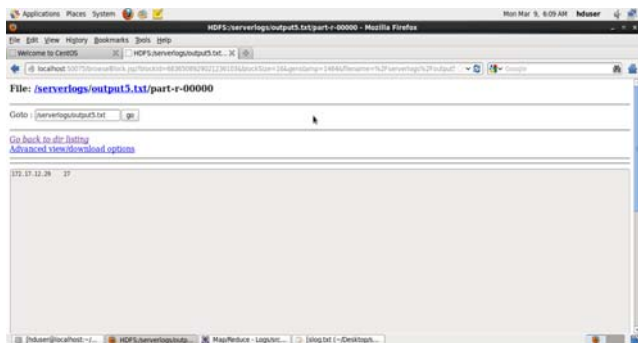**Fig 8.MapReduce functions for Logs**



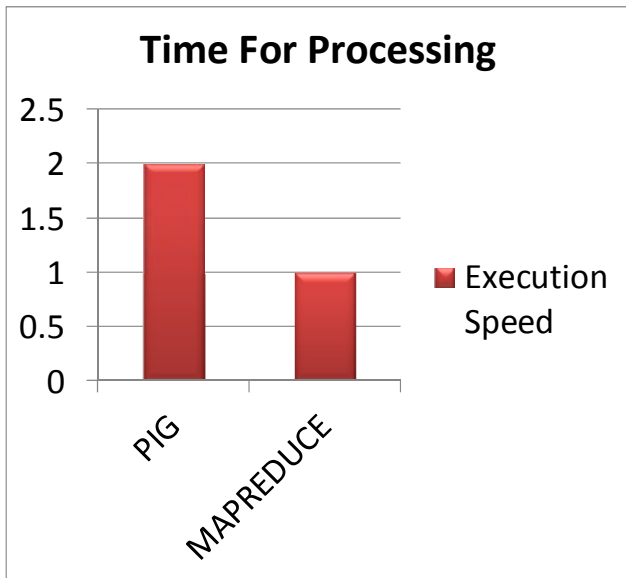**Fig 9.IPAddress Count Retrieval in Hadoop Using Mapreduce**



**Fig.10 Comparative analysis between Pig and MapReduce**

## CONCLUSIONS

As the amount of data in real-world increases, the smoother way to analyze such data plays a vital role. Hence our paper enables the easier way to analyze such huge amount by Using Apache Pig and Map Reduce and their by finding the better way to reduce communication cost and time complexity.

## REFERENCES

[1] Ruchi Verma, Sathyan R Mani, "Use of Big Data Tehnologies in Capital Markets," 2012 Infosys Limited, Bangalore, India.

[2] James Manyika, Brad Brown et.al, "Big Data: The next frontier for Innovation, Competition, and Productivity," McKinsey Global Institute, June 2011.

[3] Murat Ali , Ismail Hakki Toroslu, "Smart Miner: A New Framework for mining Large Scale Web Usage Data," WWW 2009, April 20-24. 2009 Madrid, Spain. ACM 978-1-60558-487-4/09/04.

[4] Sayalee Narkhede and Tripti Baraskar, "HMR Log Analyzer: Analyze Web Application Logs Over Hadoop MapReduce," International Journal of UbiComp (IJU) vol.4, No.3, July 2013.

[5] Milind Bhandare, Vikas Nagare et al., "Generic Log Analyzer Using Hadoop Mapreduce Framework," International Journal of Emerging Technology and Advanced Engineering (IJETAE), vol.3, issue 9, September 2013.

[6] Kanchan Sharadchandra Rahate et al., "A Novel Technique for Parallelization of Genetic Algorithm using Hadoop," International Journal of Engineering Trends and Technology (IJETT), vol.4, issue 8, August 2013.

[7] T. K. Das et al., "BIG Data Analytics: A Framework for Unstructured Data Analysis," International Journal of Engineering and Technology (IJET) vol 5, No 1, Feb-Mar 2013.

[8] Joseph McKendrick, "Big Data, Big Challeneges, Big Opportunities: 2012 IOUG Big Data strategies survey," September 2012.

[9] Jeffrey Dean and Sanjay Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters," Google, Inc

[10] Tom White, "Hadoop: The definitive Guide," Third Edition, ISBN: 978- 1-449-31152-0-1327616795.

[11] Ian Mitchell, Mark Locke and Aundy Fuller, "The White Book of Big Data"