# A New HadoopBased Network Management System with Policy Approach

Sahar Namvarasl, Vajihe Abdi, Reza Javidan

*Department of Computer Engineering and IT, Shiraz University of Technology, Modarres Blvd, Shiraz, Iran*
*sahar.namvarasl@gmail.com, abdivajihe@gmail.com, reza.javidan@gmail.com*

## ABSTRACT

In recent years with the improvement in the field of network technology and decreasing of technology cost, lots of data are produced. This massive amount of data needs mechanism for processing and mining information rapidly. In this paper a new Hadoop based network management system with policy approach which is considered hierarchical manager is presented. Storing and processing massive data efficiently are two capability of Hadoop technology by using HDFS and MapReduce. In this paper, processing time is considered as a main factor. As a result it is proved that this management system using policy approach increases the performance of entire system without putting on extra cost for implementation. This system in contrast with pure Hadoop and centralized system is several times more rapid.

**Keywords**: Network Measurement, Policy Based System, Hadoop, Big data

## 1. INTRODUCTION

In view point of scale, there are different types of computer networks as follows: PAN, LAN, CAN, MAN, WAN and GAN. A metropolitan area network which is called MAN connects some local area network (LAN). Also this type of computer network is limited in a size of a town/city [1].

Network management is a process which aims to control huge amount of data produced in networks in order to increase its efficiency and productivity. Network management is trying to afford to the large amount of data, produced in various networks and to be convinced with its scalability, reliability, performance, efficiency and transparency presented to users [2].

Policy based approach is a kind of system that swap between various system rules, defined by programmer, in different situation.

There are three architectures for network management: centralized, distributed and hierarchical. The first architecture, centralized management, is used when data are collects from networks via a single management system and send back controlling data to each node. In this case the management system acts like client-server. Distributed management, the second one, is used when several management systems are situated in all network domains, collect data from each local zone in order to complete its role. Hierarchical management is used when monitoring, displaying, storing and processing are stand on various devices aside [2]. In this paper different levels of distributed management system where monitoring is distributed are

considered. The first one is to locate each manager system on each MAN. In this case every manager system would gather its essential information from its local MAN in order to control its local area. In order to manage extensive area, a manager system is located to control different manager systems located in different MANs. These steps can go up based on the area of management task. Besides common advantages of distributed management, one of the advantages of using this kind of architecture in this paper is that management system could place everywhere without puts on extra cost.

Project of Apache Hadoop is aimed to improve extensibility, reliability and distributed computing. Hadoop is an open source and distributed infrastructure which is Google's Cloud computing platform idea and Apache Software Foundation development. Good extensibility to store, having storage mechanism, processing large amount of data and Fault-tolerant in periods of computing are some of advantages by using Hadoop [3]. Hadoop has two main components: HDFS and MapReduce. HDFS is Hadoop distributed file system that is used to store large data sets and access data in applications with high efficiency [4]. MapReduce is a parallel programming for computing large data sets stored in HDFS [5]. In this paper a measurement system for wide network based on Hadoop and policy approach using the common architecture of network management (distributed) is designed in order to achieve their advantages in a wide network environment.

The architecture of Hadoop is deployed in Cluster form, which contains one master (which called Namenode too), some slaves (which also named Datanode) and some probes (client) having data to store and processing function in this system is shown in Figure 1.
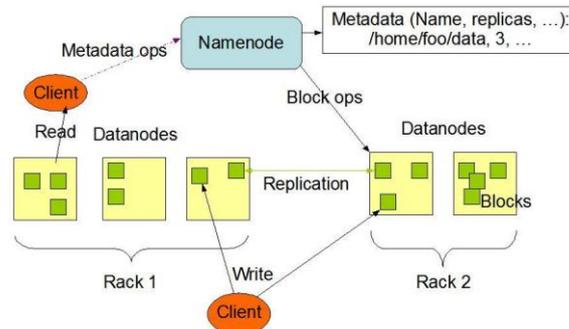


FIGURE 1. Architecture of HDFS [6]

Master does some functions such as managing file system and adjusting client's (manager) access and slave is responsible for managing storage. In cluster architecture of Hadoop deployment, each node (master or slaves) stores probe's data in Hadoop DFS. In order to store data, at first a connection establishes between client and master. Second the master decides to store data in which slaves. If there is massive data to store, data will break into blocks and each block stores in determined slaves. Also each block is replicated based on number of replicas specified in Hadoop adjustment.

Using MapReduce programming which is implemented in each node stands in the cluster, cause parallel computing and hence reduces total processing time. When a manager asks for a process on his data with calling methods of MapReduce, master determines slaves containing related data and sends them map function; by passing a key. The operation (map) will be done based on received key and After Map

operation, results are reduced and send to client [6]. According to declared explanation MapReduce function is shown in Equations (1) and (2).

$$Map : (K1, V1) \rightarrow list(K2, v1) \qquad (1)$$

$$Reduce: list (K2, list(V2)) \rightarrow list (k3, v3) \qquad (2)$$

In this paper a new Hadoop based network management system with policy approach is suggested. These components are applied in this system: HDFS and MapReduce. Also this system has good flexibility, reliability and scalability in a wide environment.

The rest of this paper is organized as follows: in Section 2 some related works are explained, in Section 3 a new Hadoop based network management system with policy approachis proposed and in Section 4 experiment and analysis presented. Finally in Section 5 a conclusion is made.

## 2. RELATED WORKS

In order to manage network in a wide environment, there is some suggested architecture. The first idea for achieving this goal is to locate one central management system in network, collecting data from entire components. But this management system is prone to become a bottleneck if network is specified in a wide environment with massive data.

D. Zhang et al. introduce a framework model for next generation network management in [7]. In this framework a cloud is created to present management services so that network management systems use these services to tackle their underlying physical network. General characteristics of the cloud to cloud, information and communication and functional model of this framework is observable in [7].

H. Wang et al. in [8] use a cloud pattern to introduce their network traffic analysis platform. That Architecture consists of three layers; resource collection layer which responsible for administrative tasks on cloud environment, resource management layer which routinely takes information from the lower layers and performs necessary calculations and stores result in computing resource table, finally the task of open service layer is analysis and distribution of traffic. Toward implementation, IPTAS (IP Trace Analysis System) is used to provide service for traffic distribution and analysis.

Y. Bao et al. in [9] also introduce layer cloud based network management architecture; the difference is that to implement the cloud part has used specifically Hadoop. In the architecture that is presented in this paper, Hadoop cluster is implemented in cloud data center layer and uses separate PCs for cluster. The main problem of these two architectures is that both of them consider separate system for cloud; in addition it imposes costs, it also needs more time to transmit data to cloud system.

There is also some related works about Hadoop main system, HDFS and MapReduce, and other Hadoop technologies such as Flume and Cukwa for data collection, Cassandra and CouchDB for storage and Pig, Hive and Mahout for

processing. One of these related works is distributed storage of network measurement data on Hbase [10]. According to [10] traditional storage system that uses relational database has limited capacity, no extensibility and elasticity. More over this kind of storage does not tackle processing and storage of big data. Potential problem of offline data storage and its basic feature, MySQL problem and Using Hadoop and Hbase, non-relational and column-oriented distributed storage, to store measurement data in real-time and non-real-time applications have been considered in [10]. There is also a read and write comparison performance of Hbase and MySQL in [10].

Also because of Hadoop tools and its capability to process big data, using its component in cloud platform is considered noteworthy hence some papers have been proposed in this field. One of these papers named an internet traffic analysis method with MapReduce [11] in cloud platform ascertains flow statistical computation time improvement in MapReduce-based flow analysis method. Y. Lee et al. in [11] propose simply architecture consists of one cluster master and some cluster nodes which are standing physically in different networks. Cluster nodes gather data from their network and submit it to cluster master. But it is not a layered architecture like [8, 9] and so reduce the speed and complexity. More accurate system of this architecture is introduced in [12].

According to [12], Z. Quing et al. present a distributed network measurement system based on Hadoop which contain two parts: one for measuring network by using nodes called probes; collecting raw data from network and the other is distributed network management consists of three parts: centralized management platform, data analyzer that performs Hadoop functions in a cloud environment and a load balancer. In this architecture, after each probe collects data, send them to one of the Hadoop nodes. To establish a connection with the best data analyzer depends on work load and processing, at first instant probe asks IP and port of best node exists in Hadoop platform from load balancer and after which, probe transmits its data to determined node. Because of using Hadoop in this system, good storage and computing, flexibility, reliability and scalability is provided. Although this architecture solve problem happening in centralized model, it has some problems. Having a single point of failure problem in load balancer operation, deploying one management platform and its entire related components in all networks and operate Hadoop system in a cloud platform apart from measured network, causing additional cost are some major obstacles in this system.

The other architecture for internet traffic measurement and analysis with Hadoop is [13]. Y. Lee et al. present a system which works likely based on Hadoop measurement with more details. In this system there are three format of storage: text file, binary files and a format called libpcap. The third format is similar to the type of network packet. In [13] paper they use libpcap format with heuristic algorithm. Also this system uses apache Hive which is compatible with HDFS to deal with the fast query coming from managers and finally IP, TCP, and HTTP traffic analysis MapReduce algorithms are performed in their system. Despite Z. Quing et al., Y. Lee et al. uses an especial format, libpcap, so that converting network packet format and libpcap format performs more easily but this work aggravates overload. Using Hive in this system whereas does fast queries, needs more configuration versus running pure Hadoop.

In this presented paper, a new Hadoop based network management system with policy approach is defined in order to solve declared problems whereas having their advantages. In this offering suggested architecture, Hadoop configuration, using

HDFS and MapReduce, does no utilize an unknown packet format; hence no overload is sustained and also managers can retrieve information from HDFS at any time and in any place they are. The other advantage of this suggested management system is that no bottleneck exists in that; this is because of Hadoop distributed system and eke there is no single point of failure.

## 3. THE PROPOSED SYSTEM

In wide area networks large volumes of data are continuously generated. First goal of offering design is to implement network management system for widespread network, so no mechanism which able to process big data easily and rapidly is needed. Indisputable centralized management system is not able to perform this heavy process. Therefore a combination of Hadoop processing with traditional network management is used and a novel distributed network management system is introduced. Another important point is that one centralized manager is not suitable for widespread network. A better option is to use hierarchical manager. Widespread network is a collection of MAN and each MAN has at least one mid-level manager and one or more low-level manager. Therefore Hadoop nodes are expanded such that each MAN has one or more management system.

Despite all features have mention for Hadoop, due to complexity of the communication between its components; centralized system will perform better in some cases. As a result a policy based approach which depends on opportunity uses Hadoop based or centralized system is implemented.

Figure 2 shows suggested architecture with details. In this system probes are components of network and they are responsible for management data collection. Slaves and master organize Hadoop cluster and do management task consist of converting raw network data to management information usable for managers. Managers by using this information can take management decisions and apply on network.
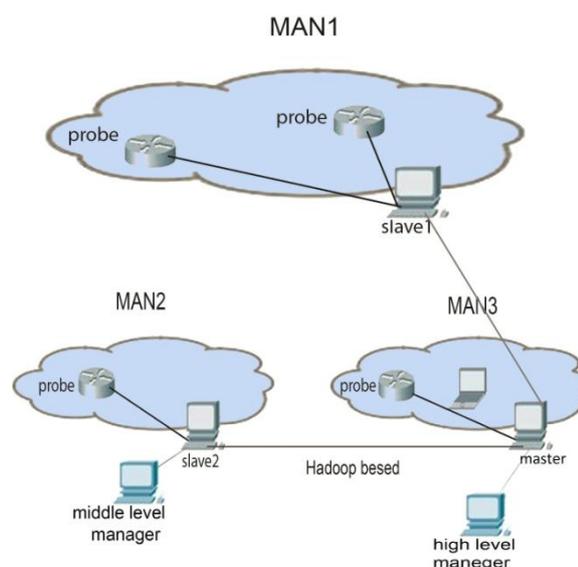


FIGURE 2. The architecture of Hadoop policy based network management

In the following more about different parts of the system is described.

## 3.1 POLICY AND RULES

This system contains different rules for each policy. The first one is according to network analyses that determine in which time the network traffic is high. The second one is according to amount of data which is gained from probes. And the last one is according to type of process which manager request and load level of these processes. Eventually based on these rules, system decides to work centralized or Hadoop based. Each time load and traffic level is high, system switch to Hadoop storing and processing. Therefore the architecture contains two main parts: local and traditional centralized system and Hadoop based system.

## 3.2 SLAVE NODES

As mentioned wide network consists of several MAN. Each MAN contains one or more slave nodes, which formed whole Hadoop system. Each slave can be one of the network components (master or slave), so network developer does not need to cost a lot to develop this network management system.

After probes collect data, send them to slave node in their MAN. As for short distance between slave node and probe locating in the same MAN, RTT will not increase incredibly.

Now slave depends on policy rules, decides to store and process locally or based on Hadoop.

## 3.3 MASTER NODE

Another part of Hadoop system is master node. Master can be one component of network with high availability. While slave decides to do management tasks by using Hadoop, does it with help of master. A master node is also responsible for segmentation or clustering. In segmentation, input file which is consisting of huge amount of data is segmented with one effective field of input. For example master node is programmed to segment the whole input data with its IP address. In order to cluster, clustering algorithms are acceptable. For example K-means clustering is used to cluster the given input file based on a centroid.

## 3.4 MANAGERS

The manager of system can be located in each position. In this architecture, hierarchical manager system is considered. At the middle level of the network management system, the manager is a part of MAN which needs management information from its own MAN, so with requesting from slave, analysis and measurement information is gained. At higher level, manager wants to manage the whole network. In this case manager should request information from master. Because of Hadoop filtering feature, it is possible that low level managers, monitor some part of network.

## 4. EXPERIMENTAL AND DISCUSSION RESULTS

Experimental results prove that this architecture work efficiently. For implementation the Hadoop cluster system on 4 PCs (CPU: Core i32.93 GHZ; memory 4G; Hard disk 500G) are considered; three nodes as slave and one node as master.

For generating network data NS2 network simulator is used. Using network simulator is a best option to consider behavior of the real network with low cost. NS2 is an open source network simulator, design specifically for researcher in the field of networking, contains different modules so designer can easily design any network, with any size [14]. For this project about 50 components for each MAN are considered and one MAN for each Hadoop node (master and slave).

The main factor that is considered for this proposed system is processing time. From the time which Hadoop node gathers data from probes until the final results are showed to manager is considered as the processing time.

Finally storing and processing operations of data is performed based on policy rules. When system prefer to operate centralized, Hadoop node in the MAN is responsible for management task. Also results are compared with centralized and fully Hadoop based system. Final result shows that Hadoop based system work more efficiency than centralized system in some situation.

According to experimental result at Table 1 and Table 2 measuring throughput parameter on centralized management with 150MB data and 10 reputations, is more rapidly than Hadoop based, but on delay measurement it is vice versa. According to experimental results, centralized method is more relevant on short period of time computing and Hadoop based method is more relevant on computing which needs to preserve log of system and long period of time processing.

TABLE1.
Average Computing Time (ms) of delay measurement on 150MB data

| Centralize | Hadoop base |
|------------|-------------|
| 106422 | 16052 |

TABLE 2.
Average Computing Time (ms) of throughput measurement on 150MB data

| Centralize | Hadoop base |
|------------|-------------|
| 3282 | 15385 |

Note that Hadoop based computing has more overload and it is not suitable on all situations. Therefore a system with policy based management approach is designed to utilize the benefit of both of centralized and Hadoop based management.

Also Hadoop nodes determine volume of generated data by network and switch to Hadoop based system when the volume of data is more than the threshold. This threshold depends on the features of node that is supposed to act as the central management system and during the execution of system this amount may be changed based on a learning algorithm. NS2 component is planned to generate more data on specific time. Thereupon experimental proves that this system is more

efficient than the system which always works on Hadoop, even with a low volume of data and response time is less.

As mentioned, one of presented system feature is reliability. This system is implemented to operate locally against failure of master; this means that when the connection between master and slave is interrupted, slave act as a centralized system.

In order to implement MapReduce, main part of processing task, has used java programming. The Java programming language is concurrent and object-oriented and it is a simple enough that programmers can expand their program easily [15], thus to adapt the system with the specific requirement, will not encounter many difficulties. Also experience with Java and Hadoop shows that this system can support any type of network measurement data. Therefore the system is applicable on any network environment.

To better evaluate the system performance, the required time for processing different amount of data for 10 repetitions is measured and mean of result is shown in figure 3. As mentioned before processing time is period of time that network data is gathered, analyzed and delivered to managers. This process involves the data storage and calculation of different parameters such as delay and throughput. The amount of data is produced by the network is greater; the size of the input data into a text file is also more. The final result is compared with centralized and pure Hadoop based system. It is clear that policy based approach in contrast with traditional system and Hadoop based system that is introduced in recent papers, has more rapid performance.

In more detailed comparison of the pure Hadoop based system and policy approach system is shown in Figure 4.
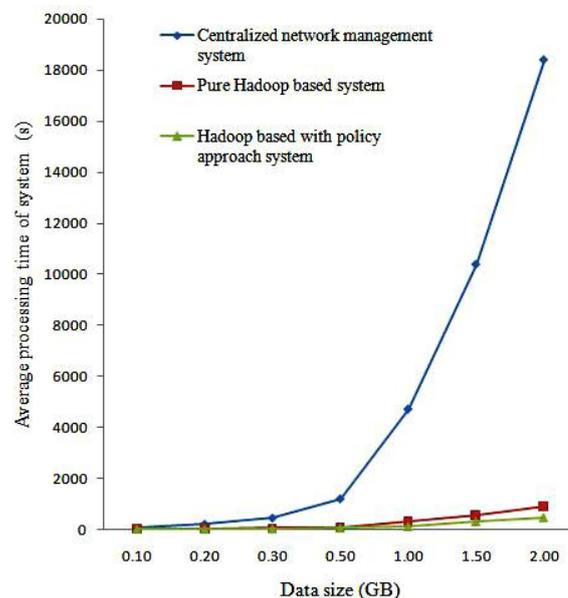


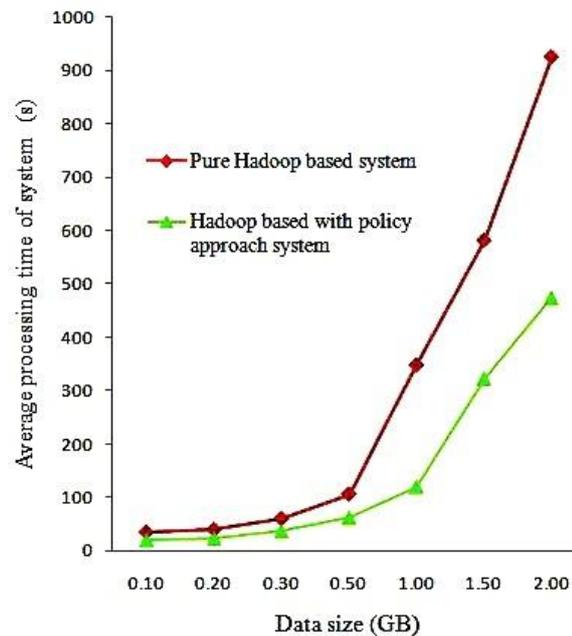FIGURE 3. Comparison processing time of 3 systems in different data size

FIGURE 4. Comparison processing time of 2 systems in different data size

The range of 95% confidence interval for comparison pure Hadoop based system and Hadoop based with policy approach, is about (22.83, 273.39) that presents significant differences between processing time of two systems.

## 5. CONCLUSION

In recent network management systems, the volume of data is so high and traditional, centralized system cannot support these amounts of data. Therefore a solution is needed for storing and processing network management data. This paper is presented a policy based system work on Hadoop which can manage widespread network and support hierarchical management for managers.

Experimental have proved that Hadoop based system can improve system computational speed up to 8 times in some situations. As regards, presented system is policy based in other situations hence computational speed is more than 3 times faster rather than pure Hadoop based system. In addition, there is no concern about storing large volumes of data. At last flexibility and reliability with high fault tolerance is considered which can recover lost data easily.

## REFERENCES

[1]  J. Ding, *Advances in Network Management*, USA: Auerbach Publications, 2009.

[2]  S. Abeck, et al., *Network Management: Know It All*, 1st. ed. USA: Elsevier by Morgan Kaufmann Publisher, 2008.

[3]  T. White, *Hadoop: The Definitive Guide*, 3d. ed. USA: O'Reilly Media / Yahoo Press, 2012.

[4]  S. Ghemawat, H. Gobioff, S. Leung, "The google file system," in *Nineteenth ACM Symposium on Operating Systems Principles*, New York, NY, USA, ACM, 2003, pp. 29-43.

[5]  J. Dean, S. Ghemawat, "MapReduce: simplified data processing on large cluster," in *Sixth Symposium on Operating System Design and Implementation,* Vol. 6, *Berkeley, CA, USA, USENIX Association*, 2004, pp. 1-13.

[6]  MapReduce Tutorial, *Hadoop*, [online] 2013,https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html#Purpose, (Accessed: May, 2014).

[7]  D. Wu, G. Zhang, J. Yang, "Cloud to cloud: a frame work model for next generation network management," in*10^{th}IEEE International Conference Telecommunication*, Vol. 1, 2003, pp. 240-245.

[8]  H. Wang, W. Ding and Z. Xia, "A cloud-pattern based network traffic analysis platform for passive measurement", in *IEEEInternational Conference on Cloud Computing and Service Computing (CSC)*, 2012, pp. 1-7.

[9]  Y.Bao et al. "Massive sensor data management framework in cloud manufacturing based on Hadoop,"in *10^{th}IEEE Industrial Informatics (INDIN) International Conference*, 2012, pp. 397 – 401.

[10] H. Ding, et al., "Distributed storage of network measurement data on Hbase," in 2^{nd}*IEEE Cloud Computing and Intelligent Systems (CCIS) International Conference*, 2012, pp. 716 – 720.

[11] Y.Lee, W.Kang and H.Son"An internet traffic analysis method with mapReduce", in*IEEE/IFIP Network Operations and Management Symposium Workshops (NOMS Wksps)*, 2010, pp. 357 – 361.

[12] Z. Quing, et al., "A distributed network measurement system based on hadoop," in8^{th}*Wireless Communications, Networking and Mobile Computing (WiCOM) International Conference on*, 2012, pp. 1 – 4.

[13] Yeonhee Lee and Youngseok Lee, "Toward scalable internet traffic measurement and analysis with hadoop," *ACM SIGCOMM Computer Communication Review*, Vol. 43, 2013, pp. 5-13.

[14] T. Issariyakul and E. Hossain, *Introduction to Network Simulator NS2*, 1st. ed., USA:Springer, 2008.

[15] J.Gosling et al., *The Java Language Specification Java SE 7 Edition*, 1st. ed., English:Addison-Wesley Professional, 2013.