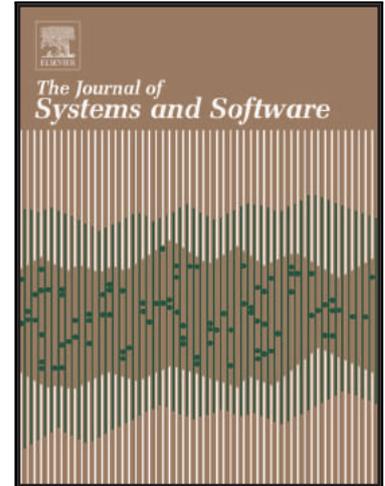


Accepted Manuscript

Performance evaluation of cloud-based log file analysis with Apache Hadoop and Apache Spark

Ilias Mavridis, Eleni Karatza

PII: S0164-1212(16)30237-0
DOI: [10.1016/j.jss.2016.11.037](https://doi.org/10.1016/j.jss.2016.11.037)
Reference: JSS 9889



To appear in: *The Journal of Systems & Software*

Received date: 22 April 2016
Revised date: 19 October 2016
Accepted date: 23 November 2016

Please cite this article as: Ilias Mavridis, Eleni Karatza, Performance evaluation of cloud-based log file analysis with Apache Hadoop and Apache Spark, *The Journal of Systems & Software* (2016), doi: [10.1016/j.jss.2016.11.037](https://doi.org/10.1016/j.jss.2016.11.037)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Highlights

- We focus on three performance indicators, the execution time, resource utilization and scalability.
- We conducted realistic log file analysis experiments in both frameworks.
- We proposed a power consumption model and an utilization-based cost estimation.
- We experimentally confirmed Spark's best performance.

ACCEPTED MANUSCRIPT

Performance evaluation of cloud-based log file analysis with Apache Hadoop and Apache Spark

Ilias Mavridis^a, Eleni Karatza^b

^a I. Mavridis Department of Informatics, Aristotle University of Thessaloniki
54124 Thessaloniki, Greece

^b E. Karatza Department of Informatics, Aristotle University of Thessaloniki
54124 Thessaloniki, Greece
Tel.: +30-2310-997974
Fax: +30-2310-997974

Abstract

Log files are generated in many different formats by a plethora of devices and software. The proper analysis of these files can lead to useful information about various aspects of each system. Cloud computing appears to be suitable for this type of analysis, as it is capable to manage the high production rate, the large size and the diversity of log files. In this paper we investigated log file analysis with the cloud computational frameworks ApacheTM Hadoop[®] and Apache SparkTM. We developed realistic log file analysis applications in both frameworks and we performed SQL-type queries in real Apache Web Server log files. Various experiments were performed with different parameters in order to study and compare the performance of the two frameworks.

Keywords: Log Analysis, Cloud, Apache Hadoop, Apache Spark, Performance Evaluation

1. Introduction

Log files are a very important source of information, useful in many cases. However as the scale and complexity of the systems are being increased, the analysis of log files is becoming more and more demanding. The effort of collecting, storing and indexing a large number of logs is aggravated even more when the logs are heterogeneous.

The log production rate can reach to several TeraBytes (TB) or PetaBytes (PB) per day. For example Facebook dealt with 130 TB of logs every day in 2010 [1] and in 2014 they have stored 300 PB of logs [2]. Due to the size of these datasets conventional database solutions cannot be used for the analysis, instead

Email addresses: imavridis@csd.auth.gr (Ilias Mavridis^a, Eleni Karatza^b),
karatza@csd.auth.gr (Ilias Mavridis^a, Eleni Karatza^b)

virtual databases combined with distributed and parallel processing systems are considered more appropriate [3].

Although there are still open challenges, cloud or even interconnected cloud systems [4] meet the needs of log analysis. Cloud is a field proven solution that is used for years by many big companies like Facebook, Amazon, ebay, etc. to analyze logs. Also in academia there are many studies that investigated cloud computing (mainly Hadoop) to analyze logs [5] - [15]. By the aforementioned, we can assume that log analysis is a big data use case and as a result many of the important challenges of cloud-based log file analysis are actually big data challenges. These challenges like data variety, data storage, data processing, and resource management are studied in many works [16].

The data variety research issues are related to the handling of the increasing volume of data, to the extraction of meaningful content and to the aggregation and correlation of streaming data from multiple sources. Data storage related issues are about efficient methods to recognize and store important information, techniques to store large volumes of information in a way it can be easily retrieved and migrated between data centers. The data processing and resource management are about programming models optimized for streaming and multidimensional data, engines able to combine multiple programming models on a single solution and resource usage and energy consumption optimization techniques.

Also the analysis of log files comes with some extra challenges. Since many systems are distributed and heterogeneous, logs from a number of components must be correlated first [17]. Moreover, maybe there are missing, duplicate or misleading data that makes log analysis more complex or even impossible. Furthermore many analytical and statistical modeling techniques do not always provide actionable insights [17]. For example, statistical techniques could reveal an anomaly in network utilization, but they do not explain how to address this problem.

Hadoop is the framework that has mainly been used for log analysis in cloud by many big companies and in academia too [5]-[15]. Hadoop was originally designed for batch processing providing scalability and fault tolerance but not fast performance [18]. It enables applications to run even in thousands of nodes with Petabytes of data. Hadoop manages the large amount of log data by breaking up the files into blocks and distributes them to the nodes of the Hadoop cluster. It follows a similar strategy for computing by breaking jobs into a number of smaller tasks. Hadoop supports quick retrieval and searching of log data, scales well, supports faster data insertion rates than traditional database systems and is fault tolerant [19].

Hadoop follows a processing model that frequently writes and reads data from the disk; this affects its performance and makes it unsuitable for real-time applications [20]. On the other hand, Spark uses more effectively the main memory and minimizes these data transfers from and to disk, by performing in-memory computations. Also Spark provides a new set of high-level tools for SQL queries, stream processing, machine learning and graph processing [21].

The main contribution of this work is the performance evaluation of log file

analysis with Hadoop and Spark. Although log analysis in cloud has been studied in many works like [5]-[15], most of them are about Hadoop based tools and they ignore the powerful Spark. Our work complements the existing studies by investigating and comparing the performance of realistic log analysis applications in both Hadoop and Spark. To the best of our knowledge this is the first work that investigates and compares the performance of real log analysis applications in Hadoop and Spark with real world log data. We conducted several experiments with different parameters in order to study and compare the performance of the two frameworks. This paper extends our previous work [22] by introducing two additional performance indicators, scalability and resource utilization. Also we proposed a power consumption model and we make an utilization-based cost estimation. In order to further investigate the performance of log analysis we developed additional applications and we conducted a new series of experiments in both frameworks.

The rest of the paper is organized as follows. Section II provides an overview of related research. Section III describes briefly what is a log file and log file analysis in cloud. Section IV outlines the two open source computing frameworks Hadoop and Spark. Section V describes the setup of our experiments. Section VI presents the experimental results and the evaluation of the frameworks. Finally Section VII concludes this paper.

2. Related work

Cloud computing for log file analysis is presented in several studies. In [5], the authors highlighted the differences between traditional relational database and big data. They claimed that log files are produced in higher rate than traditional systems can serve and they proposed and presented log file analysis with Hadoop cluster.

A weblog analysis system based on Hadoop HDFS, Hadoop MapReduce and Pig Latin Language was implemented in [6]. The proposed system aims to overcome the bottleneck of massive data processing of traditional relational databases. The system was designed to assist administrators to quickly analyze log data and take business decisions. It provides an administrator monitoring system, problem identification and systems future trend prediction.

Narkhede et al. presented a Hadoop based system with Pig for web log applications [7]. A web application was created to store log files on Hadoop cluster, run MapReduce jobs and display results in graphical formats. The authors concluded that Hadoop can successfully and efficiently process large datasets by moving computation to the data rather than moving data to computation.

In agreement with [6] and [7], [8] proposed mass log data processing and data mining methods based on Hadoop to achieve scalability and high performance. To achieve scalability and reliability, log data are stored in HDFS and are processed with Hadoop MapReduce. In this case, the experimental results showed that the Hadoop based processing and mining system improved the performance of statistics query.


```

1 157.55.33.xx0 - - [16/Jan/2013:17:57:48 -0800] "GET /webboard/ HTTP/1.0" 503 441
2 27.130.254.xx9 - - [16/Jan/2013:17:57:48 -0800] "GET / HTTP/1.1" 200 6210
3 66.249.75.xx3 - - [16/Jan/2013:17:57:51 -0800] "GET /webboard/index.php?area=statistics;u=24404 HTTP/1.1" 503 441
4 157.56.93.xx3 - - [16/Jan/2013:17:58:45 -0800] "GET /webboard/index.php?topic=2635.0 HTTP/1.0" 503 441
5 61.19.227.xx0 - - [16/Jan/2013:17:58:53 -0800] "GET /uploads/menu_1_324.gif HTTP/1.1" 200 6018
6 223.205.40.xx9 - - [16/Jan/2013:17:58:54 -0800] "GET /index1.php?id=13 HTTP/1.1" 200 15095
7 69.171.228.xx7 - - [16/Jan/2013:17:58:54 -0800] "GET /index1.php?id=13 HTTP/1.1" 200 15095
8 223.205.40.xx9 - - [16/Jan/2013:17:58:55 -0800] "GET /poll/chartpoll2.php HTTP/1.1" 200 1220
9 223.205.40.xx9 - - [16/Jan/2013:17:58:55 -0800] "GET /poll/chartpoll.php HTTP/1.1" 200 1209
10 223.205.40.xx9 - - [16/Jan/2013:17:58:55 -0800] "GET /FCKeditor/UserFiles/Image/dr_sda.jpg HTTP/1.1" 200 17158

```

Figure 1: Apache web access log.

knowledge this is the first work that studies the performance of real log analysis applications in Hadoop and Spark with real world log data.

3. Log file analysis

Log file analysis is the analysis of log data in order to extract some useful information [21]. A proper analysis requires a good knowledge of the device or software that produces the log data. It must be clear how the system that produces the logs works and what is good, suspicious or bad for it. Worth noting that a same value in two different systems maybe is bad for one but completely normal for the other.

3.1. Log files

Each working computer system collects information about various operations. This information is stored in specific files which called log files [23]. Log files are consisting of log messages or simply log(s). A log message is what a computer system, software, etc. generates in response to some sort of stimulation [19]. The information that pulled out of a log message and declares the reason that the log message was generated is called log data [19].

A common log message contains the timestamp, the source, and the data. The timestamp indicates the time at which the log message was created. The source is the system that created the log message and the data is the core of the log message. Unfortunately this format is not a standard; a log message can be significantly different from system to system.

One of the most common used types of log files in the web is the log file that is generated by the Apache HTTP Server [24]. Figure 1 shows the first ten lines of a real Apache web access log file.

As shown in Figure 1, the first element of each row is the ip address of the client (e.g., 157.55.33.xx0) or may be the name of the node that made the request to the server. Then there are two dashes (- -) that means that there is no value for this two fields [24]. The first dash stands for the identity of the client specified in RFC 1413 [25], and the second dash represents the user id of the person requesting the server. The fourth element in each of these log messages indicates the date and time that the clients request had been served by the server and in the same brackets there is the servers time zone (e.g., -0800). Next in double quotes is the request line from the client. The request line first contains the method that has been used by the client (e.g., GET), then

is the requested source (e.g., /poll/chartpoll2.php) and finally the used protocol (e.g., HTTP/1.1). At the end of each row there are two numbers that follow the request line. The first number is the status code that the server returns to the client (e.g., 200) and the last number indicates the size of the object returned to the client and is usually expressed in bytes (e.g., 6210).

3.2. Log analysis in the Cloud

The growing need for processing and storing resources for log analysis can be fulfilled by the virtually unlimited resources of cloud. By the combination of cloud and log analysis, a new term emerged the Logging as a Service (LaaS). With LaaS anybody can collect, store, index and analyze log data from different systems, without the need of purchasing, installing and maintaining his own hardware and software solutions.

There are some cloud service providers that maintain globally distributed data centers who undertake to analyze log files for one of their clients [26]. Usually these providers offer real-time, birds-eye view of the logs, a customized dashboard with easy to read line, bar, or pie charts. Users of such services can collect logs from various devices and software and submit them to the cloud using pre-built connectors for different languages and environments, as shown in Figure 2.

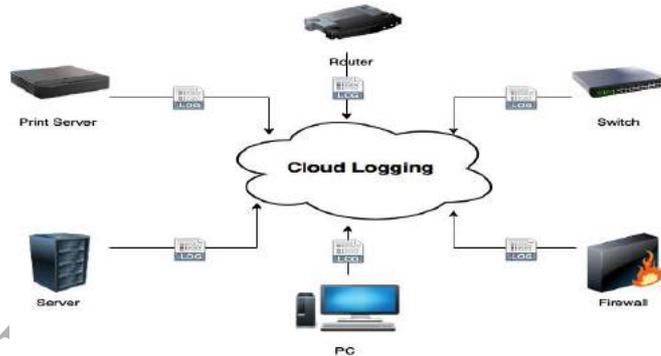


Figure 2: Logging as a Service (LaaS).

To uploading logs to the cloud, most LaaS providers support Syslog (the most widespread logging standard), TCP/UDP, SSL/TLS and a proprietary API typically RESTful over HTTP/HTTPS [19]. Also many providers have custom APIs in order to support not Syslog logs and a customizable alerting system.

The LaaS is a quite new cloud service but there are already providers like [27] and [28] which offer different service products. There are some features and capabilities common to all and some others that may vary from provider to provider. Every LaaS must support file uploading from the user to the cloud,

indexing of data (for fast search, etc.), long-term storage and a user interface for searching and reviewing the data [19]. Most providers also support various types of log formats and have their own API [27] [28]. Moreover, the majority of providers charge their services with the model pay as you go, where the user is charged depending on the use of services he has made.

On the other hand, not every LaaS provider is the same. For example there are differences in the way that each provider has developed its system [19], some providers have built their services to another cloud infrastructure, while others in their own cloud infrastructure. Also, although all LaaS providers offer the possibility of long-term storage of data, the charges are not the same and the highest possible storage time differs also. Furthermore as it was expected the charges vary from provider to provider.

4. Cloud computational frameworks

We studied two of the most commonly used cloud computational frameworks, Hadoop and Spark. Hadoop is a well-established framework that is used for storing, managing and processing large data volumes for many years by web behemoths like Facebook, Yahoo, Adobe, Twitter, ebay, IBM, LinkedIn and Spotify [29]. Hadoop is based on MapReduce programming model for batch processing. However the need for faster real-time data analysis led to a new general engine for large-scale data processing, Spark. Spark was developed by AMPLab [30] of UC Berkeley and it can achieve up to one hundred times higher performance compared to Hadoop MapReduce [31] by using more effectively the main memory.

4.1. Apache Hadoop

Hadoop origins from Apache Nutch [32], which is an open source search engine. Basic elements to the creation of Hadoop were two Google studies, the first one was published in 2003 and describes the Google Distributed Filesystem (GFS) [33] and the second one was published in 2004, and describes MapReduce [34]. In February 2006 a part of Nutch became independent and created Hadoop. In 2010 a team from Yahoo began to design the next generation of Hadoop, the Hadoop YARN (Yet Another Resource Negotiator) or MapReduce2 [35].

YARN changed the resource management of Hadoop and supported a wide range of new applications with new features. YARN is more general than MapReduce (Figure 3), in fact MapReduce is a YARN application. There are other YARN applications like Spark [36], which can run parallel to the MapReduce, under the same resource manager.

4.1.1. Hadoop Distributed File System

Hadoop Distributed File System (HDFS) stores files into blocks of a fixed size in different nodes of Hadoop cluster [37] (Figure 4). By dividing files into smaller blocks, HDFS can store much bigger files than the disk capacity of each

To run the experiments we saved the data in HDFS, in blocks of 128 MB. Because the number of active nodes will be different during the experiments every data block was copied to all Datanodes. By this way even if there is only one active Datanode still it will have all the necessary blocks to recreate the original file.

Also in order to execute SQL queries the data must be structured in tables. As the logs are already in HDFS there is no need to re-upload the data. Hive can easily use this data, define the relationships and create a virtual relational database. Using Hive [57] and the appropriate regular expressions we created the database tables in a few seconds.

5.3. Metrics and monitoring tools

In our experiments we took measurements related to the application execution time and resource utilization. These metrics help us to understand the performance impact of different parameters such as input dataset, number of active nodes and type of application. We recorded utilization information about CPU, main memory, disk and network for every node for every moment that the experiments run.

Sysstat [58] and Ganglia [59] are the main tools that were used to analyze the resource utilization of the cluster nodes. Sysstat is a monitoring package for Linux systems with low overhead that can be easily used to gather system performance and usage activity information. There are Cloud specific monitoring tools such as Amazon CloudWatch, AzureWatch, Nimsoft and Cloud-Kick [60] but the main disadvantage is that many of them are commercial and provider-dependent. For instance, Amazon Cloud Watch monitors Amazon Web Services(AWS) such as Amazon EC2, Amazon RDS DB instances, and applications running on AWS.

On the other hand there are General purpose monitoring tools that were developed before the advent of Cloud Computing for monitoring of IT infrastructure resources. Many of these tools are evolved and adopted in Cloud. Such tools are Ganglia, Nagios, Collectd, Cacti and Tivoli. General purpose infrastructure monitoring tools usually follow a clientserver model by installing an agent in every system to be monitored.

We selected the open source Ganglia rather than other powerful tool like JCatascopia [61] because Ganglia is lightweight and in many cases adds the lowest overhead at the monitoring hosts compared to other tools. It has features such as scalability, portability and extensibility [60] and it's well documented. Also Ganglia fulfills the monitoring requirements of our experiments as it supports Hadoop and Spark with more than one hundred special built-in metrics like HDFS and YARN metrics.

5.4. Experimental scenarios

In order to investigate the log file analysis with Hadoop and Spark we developed realistic applications in both frameworks. All applications are developed in java and are developed with the Eclipse IDE and Apache Maven. The log

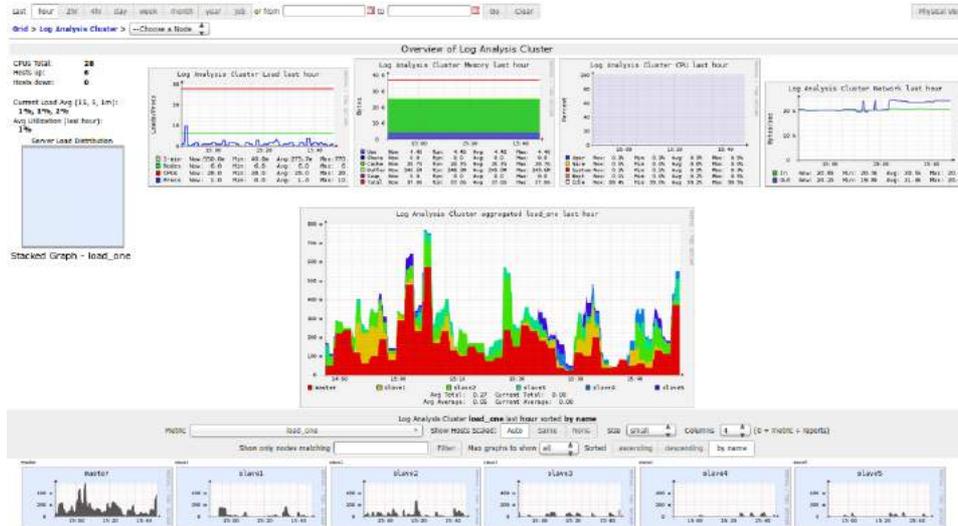


Figure 6: Ganglia monitoring tool.

files that we used for our experiments are stored in HDFS and are accessible by both Hadoop and Spark applications. All applications read the input data from HDFS and produce output results in simple format that can be represented easily in charts and graphs. To parse the log data we used regular expressions.

As we are dealing with http log files, we considered that it would be quite possible for an administrator of a webpage to ask for information about traffic distribution, detection of possible threats and error identification.

Concerning the traffic distribution we developed two applications. The first application measures the requests to the web server and sorts them by day. By this way we can find out how the overall traffic is distributed into the seven days of the week and with minor changes in the source code it could be possible to find the traffic per seconds till years. The second application finds the ten most popular requests and displays them in descending order.

Taking into consideration that one of the most common attacks is denial of service attack (DoS attacks) [62], we developed two applications for detecting and help preventing this kind of malicious attack. So the third application that we developed for both frameworks searches for moments with abnormal big number of identical requests to the server and displays the time and the number of requests. The fourth application examines the suspicious for DoS attack moments, counts and displays the requests that every client made. By this way it is possible to detect the id of suspicious clients and take prevention measures like updating the systems firewall.

Finally we developed another two applications for error detection and correction. In the http log file when a response code is bigger than 300 that means

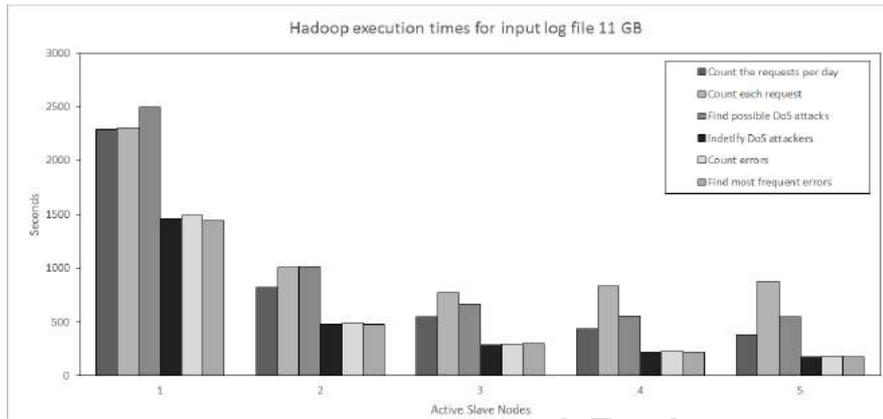


Figure 9: Execution time of Hadoop applications with 11 GB input file.

Map and Reduce execution times

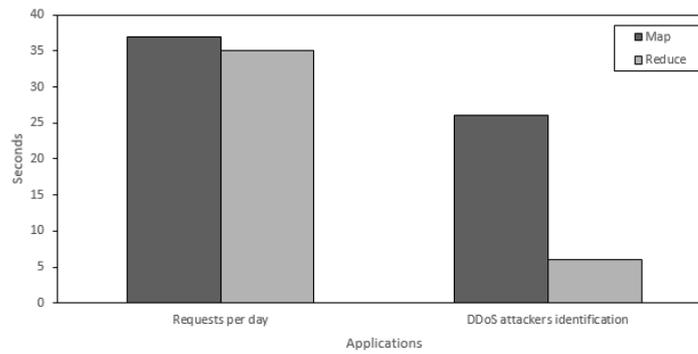


Figure 10: Execution time of Map and Reduce phases for two different applications.

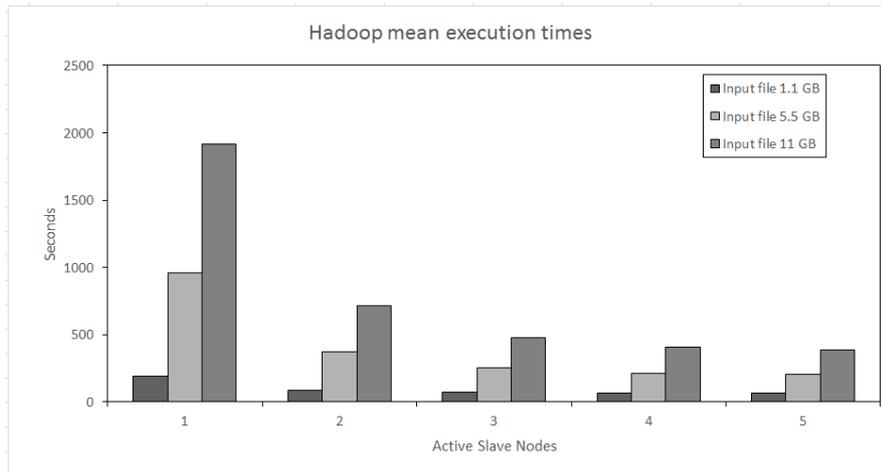


Figure 11: Mean execution time of all Hadoop applications.

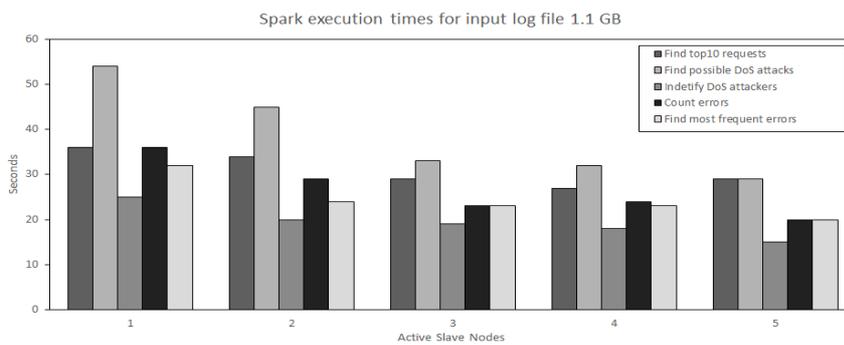


Figure 12: Execution time of Spark applications with 1.1 GB input file.

CPU is considered the most power demanding component. There are detailed analytical power models like [85] which are real-time and highly accurate, but are not so general and portable because they are based on microarchitectural knowledge of a particular processor. Also these models rely on the assumption that CPU is the main power consumer, which is invalid for workloads that involve a large amount of I/O activity, as in the case of MapReduce [86].

Many studies are based on the correlation between power consumption and resource utilization. Like Lee and Zomaya [94] that in their model assume that the relation between CPU utilization and energy consumption is linear. Chen et al [89] propose a linear power model too. Their model take into consideration the power consumption from individual components to a single work node. Joulemeter [90] monitors the resource usage of VMs and then converts it to energy consumed based on the power model of each individual hardware component.

Bohra et al. [86] also proposed a power consumption model that is based on the correlation between the power consumption and components utilization. The authors proposed a four-dimensional linear weighted power model (eq. 1) for the total power consumption.

$$P_{\text{total}} = c_0 + c_1 P_{\text{CPU}} + c_2 P_{\text{cache}} + c_3 P_{\text{DRAM}} + c_4 P_{\text{disk}} \quad (1)$$

In eq. 1 the P_{CPU} , P_{cache} , P_{DRAM} and P_{disk} are specific performance parameters for CPU, cache, DRAM and disk, and c_0, c_1, c_2, c_3, c_4 are weights.

In [87] Chen et al. modified the proposed model of Bohra et al. to the following model:

$$P_{\text{total}} = P_{\text{idle}} + c_1 P_{\text{CPU}} + c_2 P_{\text{HDD}} \quad (2)$$

We also further modify the Chen et al. model by adding the contribution of network hardware and the model transformed to eq. 3.

$$P_{\text{total}} = P_{\text{idle}} + c_1 P_{\text{CPU}} + c_2 P_{\text{RAM}} + c_3 P_{\text{HDD}} + c_4 P_{\text{Network}} \quad (3)$$

To compare Spark and Hadoop-MapReduce power consumption we applied our proposed model based on the CPU, memory, disk and network utilizations. We assumed that the hardware weights are the same for Hadoop and Spark applications. Based on the utilization measurements for the count errors application with input file of 5.5 GB and five slave nodes cluster, we found the relationships of power consumption of Hadoop and Spark for the specific configuration (eq. 4).

$$P_S = P_{\text{CPUHadoop}} + 1.9P_{\text{RAMHadoop}} + 3.8P_{\text{HDDHadoop}} + 1.9P_{\text{NetworkHadoop}} \quad (4)$$

Spark. In order to evaluate the performance of the frameworks we developed, executed and evaluated the performance of realistic log analysis applications in both of them.

Hadoop is one of the first frameworks for cloud computing, is widely known and is used for years by many big companies. Over the years Hadoop evolved and improved in order to meet the new era needs. These new needs led also to the creation of Spark. Spark is faster and more flexible than Hadoop. Spark achieved to dramatically reduce the execution time by saving the intermediate results in memory instead of the disk. Also, Spark except of MapReduce functions, supports a wide range of new capabilities that can be combined to generate new hitherto been impossible applications.

We conducted several experiments with different number of slave nodes, size of input file and type of application in order to confirm Sparks best performance. Initially we studied the total execution time of Hadoop and Spark log analysis applications and we analyzed the factors that affect it. We found that Spark due to the effective exploitation of main memory and the use of efficiency optimizing techniques is faster than Hadoop in every case.

Also as both frameworks were originally designed for scalability we conducted experiments in order to test the scalability of the two frameworks. In the experimental results we observed that the doubling of the active slave nodes leads almost to the halving of total execution time and the doubling of the input file leads almost to the doubling of total execution time for both frameworks. With these observations we made some cost estimations for the two frameworks and we investigated how we can achieve a significant cost reduction by accepting a burden in execution time.

Finally we took measurements related to the mean resource utilization of the cluster. From the various experiments we concluded that Spark has a higher mean utilization of the available resources compared to Hadoop. As it was expected Spark has higher memory utilization than Hadoop because saves intermediate results to memory. Also we found that Spark has a higher mean disk and network utilization which was quite unexpected. These unexpected experimental results can be explained by the fact that Spark processes the same amount of data but in significantly less time than Hadoop does and this leads to larger mean values of utilization related parameters. We also proposed a power consumption model and by using the utilization measurements we compared the power consumption of the two frameworks under specific experimental configuration.

The overall experiments showed Sparks best performance. However, the applications were implemented in such a way to make possible the comparison between the two frameworks. As future work could be implemented applications that make full use of Spark capabilities in order to evaluate the performance of the framework for more complex log analysis applications.

Acknowledgements

The authors would like to thank Okeanos the GRNETs cloud service for the valuable resources.

References

- [1] <https://www.facebook.com/notes/facebook-engineering/scaling-facebook-to-500-million-users-and-beyond/409881258919>
- [2] <https://code.facebook.com/posts/229861827208629/scaling-the-facebook-data-warehouse-to-300-pb/>
- [3] P. Kumar, A. Bhardwaj, A. Doegar, Big Data Analysis Techniques and Challenges in cloud Computing Environment, International Journal of Advance Foundation And Research In Science & Engineering (IJAFRSE) Volume 1, Special Issue , ICCICT 2015.
- [4] I.A. Moschakis, H.D. Karatza, A meta-heuristic optimization approach to the scheduling of Bag-of-Tasks applications on heterogeneous Clouds with multi-level arrivals and critical jobs, Simulation Modelling Practice and Theory, Elsevier, Vol. 57, Pages: 1 - 25, 2015.
- [5] B. Kotiyal, A. Kumar, B. Pant, R. Goudar, 'Big Data: Mining of Log File through Hadoop, IEEE International Conference on Human Computer Interactions (ICHCI'13), Pages: 1-7, Chennai, India, 23 Aug - 24 Aug 2013.
- [6] C. Wang, C. Tsai, C. Fan, S. Yuan, A Hadoop based Weblog Analysis System, 7th International Conference on Ubi-Media Computing and Workshops (U-MEDIA 2014), Pages: 72-77, Ulaanbaatar, Mongolia, 12-4 July 2014.
- [7] S. Narkhede, T. Baraskar, HMR log analyzer: Analyze web application logs over Hadoop MapReduce, International Journal of UbiComp (IJU), Vol.4, No.3, Pages: 41-51, July 2013.
- [8] H. Yu, D. Wang, Mass Log Data Processing and Mining Based on Hadoop and Cloud Computing, 7th International Conference on Computer Science and Education (ICCSE 2012), Pages: 197, Melbourne, Australia, 14-17 July 2012
- [9] J. Wei, Y. Zhao, K. Jiang, R. Xie, Y. Jin, Analysis farm: A cloud-based scalable aggregation and query platform for network log analysis, International Conference on Cloud and Service Computing (CSC),Pages: 354 - 359, Hong Kong, China, 12 - 14 Dec 2011.
- [10] J. Yang, Y. Zhang, S. Zhang, D. He, Mass flow logs analysis system based on Hadoop, 5th IEEE International Conference on Broadband Network and Multimedia Technology (IC-BNMT), Pages: 115 - 118, Guilin, China, 17-19 Nov. 2013.

- [11] X. LIN, P. WANG, B. WU, Log analysis in cloud computing environment with Hadoop and Spark, 5th IEEE International Conference on Broadband Network and Multimedia Technology (IC-BNMT 2013), Pages: 273 - 276, Guilin, China, 17-19 November 2013.
- [12] Y. Liu, W. Pan, N. Cao, G.i Qiao, System Anomaly Detection in Distributed Systems through MapReduce-Based Log Analysis, 3rd International Conference on Advanced Computer Theory and Engineering (ICACTE), Pages: V6-410 - V6-413, Chengdu, China, 20-22 Aug. 2010.
- [13] S. Vernekar, A. Buchade, MapReduce based Log File Analysis for System Threats and Problem Identification, Advance Computing Conference (IACC), 2013 IEEE 3rd International, Pages: 831 - 835, Patiala, India, 22-23 February 2013.
- [14] M. Kumar, M. Hanumanthappa, Scalable Intrusion Detection Systems Log Analysis using Cloud Computing Infrastructure, 2013 IEEE International Conference on Computational Intelligence and Computing Research (ICCIC), Pages: 1-4, Tamilnadu, India, 26-8 December 2013.
- [15] H. Kathleen, R. Abdelmounaam, SAFAL: A MapReduce Spatio-temporal Analyzer for UNAVCO FTP Logs, IEEE 16th International Conference on Computational Science and Engineering (CSE), Pages: 1083 - 1090, Sydney, Australia, 3-5 Dec. 2013.
- [16] M. Assunoa, R. Calheirosb, S. Bianchic, M. Nettoc, R. Buyya, Big Data Computing and Clouds: Trends and Future Directions, Journal of Parallel and Distributed Computing Volumes 7980, May 2015, Pages 315, Special Issue on Scalable Systems for Big Data Management and Analytics
- [17] A. Oliner, A. Ganapathi, W. Xu, Advances and challenges in log analysis, Communications of the ACM, Volume 55 Issue 2, February 2012, Pages 55-61, ACM New York, NY, USA
- [18] <http://wiki.apache.org/hadoop/>
- [19] Dr. Anton A. Chuvakin, Kevin J. Schmidt, Christopher Phillips, Patricia Moulder, Logging and Log Management: the authoritative guide to understanding the concepts surrounding logging and log management, Elsevier Inc. Waltham, 2013.
- [20] G.L. Stavrinides, H.D. Karatza, A cost-effective and QoS-aware approach to scheduling real-time workflow applications in PaaS and SaaS clouds, In Proceedings of the 3rd International Conference on Future Internet of Things and Cloud (FiCloud'15), Rome, Italy, Aug. 2015.
- [21] <https://databricks.com/spark>

- [22] I. Mavridis and H. Karatza, Log File Analysis in Cloud with Apache Hadoop and Apache Spark, 2nd International Workshop on Sustainable Ultrascale Computing Systems (NESUS 2015), Pages: 51 - 62, Krakow, Poland, 10-11 Sept. 2015.
- [23] Jorge Pinto Leite, Analysis of Log Files as a Security Aid, 6th Iberian Conference on Information Systems and Technologies (CISTI), Pages: 1-6, Lousada, Portugal, 15-18 June 2011.
- [24] <http://httpd.apache.org/docs/1.3/logs.html>
- [25] <http://www.rfc-base.org/rfc-1413.html>
- [26] D. Jayathilake, Towards Structured Log Analysis, 9th International Joint Conference on Computer Science and Software Engineering (JCSSE 2012), Pages: 259 - 264, Bangkok, Thailand, 30 May - 1 June 2012.
- [27] <http://www.loggly.com>
- [28] <http://www.splunkstorm.com/>
- [29] <http://wiki.apache.org/hadoop/PoweredBy>
- [30] <https://amplab.cs.berkeley.edu/software/>
- [31] R. Xin, J. Rosen, M. Zaharia, M. J. Franklin, S. Shenker, I. Stoica, Shark: SQL and Rich Analytics at Scale, SIGMOD 2013, Pages 13-24, New York, USA, 22-27 June 2013.
- [32] <http://nutch.apache.org/>
- [33] G. Sanjay, G. Howard, L. Shun-Tak, "The Google File System", 19th ACM Symposium on Operating Systems Principles, Lake George, NY, October, 2003.
- [34] J. Dean, S. Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters", OSDI'04: Sixth Symposium on Operating System Design and Implementation, San Francisco, CA, December, 2004.
- [35] <http://hortonworks.com/hadoop/YARN/>
- [36] <http://wiki.apache.org/hadoop/PoweredByYARN>
- [37] <http://hadoop.apache.org/docs/r2.4.1/hadoop-project-dist/hadoop-hdfs/HdfsUserGuide.html>
- [38] <http://hive.apache.org/>
- [39] <http://mahout.apache.org/>
- [40] <http://zookeeper.apache.org/>

- [56] https://wiki.apache.org/hadoop/FAQ#Does_Hadoop_require_SSH.3F
- [57] <https://hive.apache.org/javadocs/r0.10.0/api/org/apache/hadoop/hive/serde2/SerDe.html>
- [58] <http://sebastien.godard.pagesperso-orange.fr/>
- [59] M. Massiea, B. Chunb, D. Cullera, The ganglia distributed monitoring system: design, implementation, and experience, *Parallel Computing*, vol. 30, no 7, pp.817-840, 2004
- [60] K. Fatemaa, V. C. Emeakaroha, Ph. D. Healya, J. P. Morrisona, Th. Lynn, A survey of Cloud monitoring tools: Taxonomy, capabilities and objectives, *Journal of Parallel and Distributed Computing*, vol. 74, no 10, pp. 2918-2933, 2014
- [61] D. Trihinas, G. Pallis, M. Dikaiakos, JCatascope: Monitoring Elastically Adaptive Applications in the Cloud, *Cloud and Grid Computing (CCGrid)*, 2014 14th IEEE/ACM International Symposium on Cluster, Chicago, Illinois, USA, 26-29 May 2014.
- [62] F. Lau, S. H. Rubin, M. H. Smith, and L. Trajkovic., Distributed Denial of Service Attacks, *IEEE International Conference on Systems, Man, and Cybernetics*, Pages: 2275-2280, Nashville, TN, USA, October 2000.
- [63] <http://www.w3.org/Protocols/rfc2616/rfc2616-sec10.html>
- [64] J.Conejero, B. Caminero, C. Carron, Analysing Hadoop Performance in a Multi-user IaaS Cloud, *High Performance Computing and Simulation (HPCS)*, Pages 399 - 406, Bologna, Italy, 21-25 July 2014
- [65] G. Velkoski, M. Simjanoska, S. Ristov, M. Gusev, CPU Utilization in a Multitenant Cloud, *IEEE EUROCON 2013*, Pages 242-249, Zagreb, Croatia, 1-4 July 2013
- [66] L. Gu, H. Li, Memory or Time: Performance Evaluation for Iterative Operation on Hadoop and Spark, *2013 IEEE 10th International Conference on High Performance Computing and Communications and 2013 IEEE International Conference on Embedded and Ubiquitous Computing (HPCC EUC)*, Pages 721-727, Zhangjiajie, China, 13-15 Nov. 2013
- [67] P.R. Magalhaes Vasconcelos, G. Azevedo de Araujo Freitas, Performance analysis of Hadoop MapReduce on an OpenNebula cloud with KVM and OpenVZ virtualizations, *2014 9th International Conference for Internet Technology and Secured Transactions (ICITST)*, Pages 471-476, London, 8-10 Dec. 2014
- [68] Eug. Feller, Lav. Ramakrishnan, Chr. Morin, Performance and energy efficiency of big data applications in cloud environments: A Hadoop case study,

