

## ENHANCED ESTIMATION MODEL FOR HADOOP WORDCOUNT USING MAP REDUCING METHOD

<sup>1</sup>Neha Chouhan, <sup>2</sup>Dr. Pankaj Dashore

<sup>\*1</sup>Research Scholar, Sanghvi Innovative Academy, Indore, M.P, India

<sup>2</sup>Professor, Sanghvi Innovative Academy, Indore, M.P, India

<sup>\*</sup>Department of Computer Science & Engg.

### ABSTRACT:

*As a result of the rapid development in cloud computing, it's fundamental to investigate the performance of extraordinary Hadoop MapReduce purposes and to realize the performance bottleneck in a cloud cluster that contributes to higher or diminish performance. It is usually primary to research the underlying hardware in cloud cluster servers to permit the optimization of program and hardware to achieve the highest performance feasible. Hadoop is founded on MapReduce, which is among the most popular programming items for huge knowledge analysis in a parallel computing environment. In this paper, we reward a particular efficiency analysis, characterization, and evaluation of Hadoop MapReduce WordCount utility.*

**Keywords:** Performance analysis, cloud computing, Hadoop WordCount.

### I. INTRODUCTION

We are residing within the era of large data. In these days a tremendous amount of knowledge is generating everywhere as a result of advances within the web and verbal exchange applied sciences and the pursuits of men and women using smartphones, social media, internet of things, sensor contraptions, online offerings and lots of more. In a similar way, in improvements in knowledge applications and broad distribution of application, a couple of govt and commercial organizations such as monetary institutions, healthcare institution, schooling and research division, power sectors, retail sectors, lifestyles sciences and environmental departments are all producing a enormous amount of information everyday. For examples, international data enterprise (IDC) said that 2.8 ZB (zettabytes) knowledge of universe had been saved in the year of 2012 and this may reach up to forty ZB through 2020 [1]. In a similar fashion Facebook processes round 500 TB (terabytes) knowledge per day [2] and Twitter generates eight TB data daily [3]. The huge datasets no longer handiest comprise structured form of knowledge but greater than seventy five% of the dataset includes uncooked, semi-structured and unstructured type of data [4]. This large quantity of information with one of a kind codecs can be viewed as giant information. The derivation of big knowledge is indistinct and there are a lot of definitions on huge data. For examples, Matt Aslett outlined massive knowledge as “tremendous data is now virtually universally understood to refer to the recognition of larger business intelligence through storing, processing, and examining data that was previously ignored because of problem of normal data management applied sciences” [5]. Recently, the term of giant data has got a brilliant momentum from governments, industry and research communities. In [6], significant information is outlined as a term that encompasses using tactics to capture, approach, analyze and visualize potentially significant datasets in a cheap timeframe now not obtainable to usual IT applied sciences.

## II. LITERATURE SURVEY

[1] Samneet Singh and Yan Liu, "A Cloud Service Architecture for Analyzing Big Monitoring Data", ISSNII1007-0214II05/10Ipp55-70 Volume 21, Number 1, February 2016

In this paper, author proposed a structure that integrates search-headquartered clusters and semantic media wiki by using relaxation APIs to help the exploration of cloud monitoring data. This structure advantages from an internet-based Media-Wiki interface and enables a person to outline the entry to monitoring knowledge and prepare the processing results. The quest-based cluster developed on SolrCloud permits indexing of significant size of knowledge, and thus makes the entire architecture compatible to explore and display the ever-gathering data such as the traces constructed from knowledge centers. The structure additionally involves an extension, which runs spark on Yarn cluster for deploying effective evaluation ways for gigantic knowledge set. It utilizes the spark's MapReduce paradigm to establish the cluster in the dataset utilizing k-way clustering approach.

[2] JOSEPH A. ISSA, "Performance Evaluation and Estimation Model Using Regression Method for Hadoop WordCount", Received November 19, 2015, accepted December 12, 2015, date of publication December 18, 2015, date of current version December 29, 2015.

In this paper, the writer offered a distinct performance analysis and analysis for Hadoop WordCount workload utilizing different processors similar to Intel's ATOM D525, Xeon X5690, and AMD's Bobcat E350. Our analysis suggests that Hadoop WordCount is compute-sure workload in both map segment and scale down segment. The outcome exhibit that enabling HT and growing the number of sockets have a high impact on the Hadoop WordCount performance even as reminiscence velocity and capacity does now not have an impact on efficiency vastly.

[3] Yaxiong Zhao, Jie Wu, and Cong Liu, "Dache: A Data Aware Caching for Big-Data Applications Using the MapReduce Framework", ISSNII10070214II05/10Ipp39-50 Volume 19, Number 1, February 2014

In this paper, author recommends Dache, a knowledge-conscious cache framework for big-data functions. In Dache, tasks publish their intermediate outcome to the cache manager. A project queries the cache supervisor before executing the specific computing work. A novel cache description scheme and a cache request and reply protocol are designed. We enforce Dache by means of extending Hadoop. Test bed experiment results show that Dache tremendously improves the completion time of MapReduce jobs.

[4] Zhuoyao Zhang Ludmila Cherkasova, "Benchmarking Approach for Designing a MapReduce Performance Model", *ICPE'13*, April 21-24, 2013

In this work, author presents a novel efficiency analysis framework for answering this question. We observe that the execution of every map (lessen) duties consists of distinctive, good-defined knowledge processing phases. Handiest map and scale back services are customized and their executions are consumer-outlined for exclusive MapReduce jobs. The executions of the remaining phases are time-honored and rely on the amount of information processed by means of the phase and the performance of underlying Hadoop cluster. First, we design a suite of parameterizable micro benchmarks to measure normal phases and to derive a platform performance model of a given Hadoop cluster.

## III. EXISTING SYSTEM

### ESTIMATION MODEL USING AMDAHL'S LAW METHOD

Based on Amdahl's law definition we discussed before, the performance of a given processor can be divided into two parts, the part which increases with the performance improvement and is said to scale is defined as variable  $a$ , and the other part which does not improve due to the performance enhancement and is said to not scale is defined as variable  $b$ . The  $a$  and  $b$  variables can be derived using the basic definition of Amdahl's law which can be written in the form of:

$$T = T_o + (T_1 - T_o) \times \frac{I_1}{I} \quad (1)$$

where  $T_1$  is the measured execution time at a given input size  $I_1$ , and  $T_o$  be the non-scale execution time. We can write  $T_o$  in terms of a second measurement  $T_2$  at  $I_2$ :

$$T_o = \frac{T \times I_2 - T_1 \times I_1}{(I_2 - I_1)} \quad (2)$$

When we substitute Eq(2) for  $T_o$  in Eq(1) we obtain Amdahl's law in terms of two specific measurements without reference to:

$$T = a + b \times \frac{1}{I} \quad (3)$$

$$a = \frac{I_2 \times T_2 - I_1 \times T_1}{(I_2 - I_1)} \quad (4)$$

$$b = \frac{(I_1 \times I_2)(T_2 - T_1)}{(I_2 - I_1)} \quad (5)$$

The variables  $a$  and  $b$  can be transformed to the performance instead of the time domain by using  $P = 1/T$ . This will give us  $a$  and  $b$  variables in terms of performance and input size as shown in Eq(6) and Eq(7).

$$a = \frac{P_1 \times I_2 - P_2 \times I_1}{(P_1 \times P_2)(I_2 - I_1)} \quad (6)$$

$$b = \frac{(I_1 \times I_2)(P_2 - P_1)}{(P_1 \times P_2)(I_2 - I_1)} \quad (7)$$

For two data points, we will have  $(I_1; P_1)$  and  $(I_2; P_2)$ , and for  $n$  data points, we will have  $(I_1; P_1); \dots; (I_n; P_n)$ . We expect these points to satisfy an equation of the form (except for noise):

$$P_i = \frac{I_i}{a \times I_i + b} \quad (8)$$

Because of noise, we cannot expect to find values for  $a$  and  $b$  that produce equality for each point  $i$ . In this case, we resort to the theory of linear least-squares estimation to obtain best estimates for  $a$  and  $b$ . In particular, given  $a$  and  $b$ , we take the error in our estimate for  $P_i$  in terms of  $I_i$  to be the difference between the measured and estimated value for  $P_i$ :

$$e_i = \left( \frac{1}{P_i} - \left( a + b \cdot \frac{1}{I_i} \right) \right) \quad (9)$$

The best estimates for  $a$  and  $b$  are those that minimize the sum of the squares of these errors:

$$E = \sum_{i=1}^n 2e_i^2 \quad (10)$$

The estimates for  $a$  and  $b$  are those at which the values of the partial derivatives  $\frac{\partial E}{\partial a}$  and  $\frac{\partial E}{\partial b}$  are simultaneously zero. By computing these derivatives explicitly, we obtain equations satisfied by the best choices for  $a$  and  $b$ , which is the best functional fit to the measured data.

$$\frac{\partial E}{\partial a} = \sum_{i=1}^n 2e_i \frac{\partial}{\partial a} \left( \frac{1}{P_i} - \left( a + b \cdot \frac{1}{I_i} \right) \right) \quad (11)$$

Thus,  $\frac{\partial E}{\partial a} = 0$  implies

$$\left( \sum_{i=1}^n \frac{1}{P_i} \right) = n \cdot a + \left( \sum_{i=1}^n \frac{1}{I_i} \right) \cdot b, \quad (12)$$

$$c_{11} \equiv n, \quad (18)$$

$$c_{12} \equiv c_{21} \equiv \sum_{i=1}^n \frac{1}{I_i}, \quad (19)$$

$$c_{22} \equiv \sum_{i=1}^n \frac{1}{I_i^2}, \quad (20)$$

$$d_1 \equiv \sum_{i=1}^n \frac{1}{P_i}, \quad (21)$$

$$d_2 \equiv \sum_{i=1}^n \frac{1}{P_i} \frac{1}{I_i}. \quad (22)$$

Given these best estimates for  $a$  and  $b$  in terms of  $(I_1, P_1); \dots; (I_n, P_n)$ , we have the following best estimate for  $P$  in terms of  $I$ :

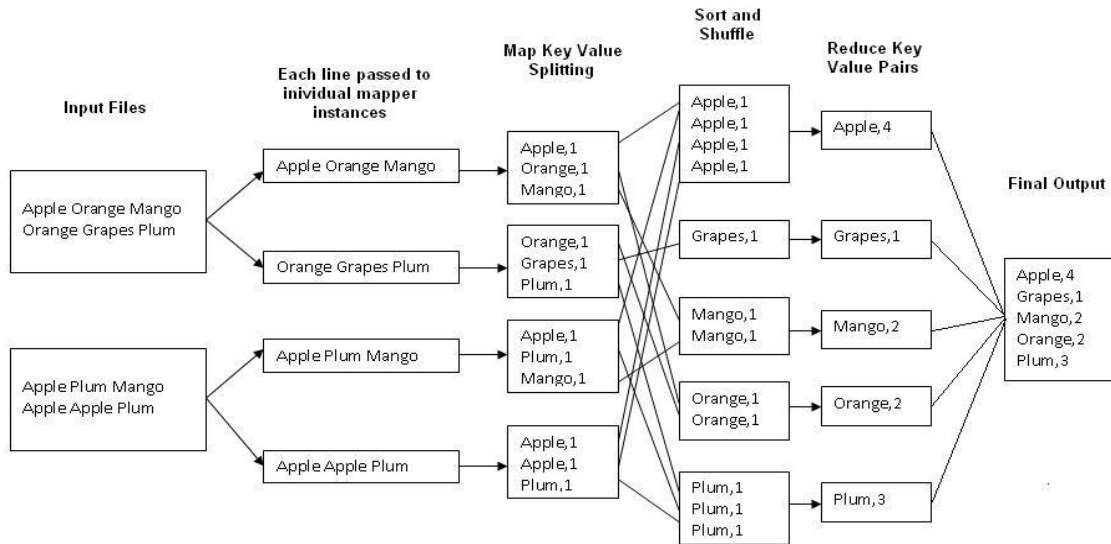
## V. PROPOSED WORK

Word count is a typical example where Hadoop map reduce developers start their hands on with. This sample map reduce is intended to count the no of occurrences of each word in the provided input files.

The word count operation takes place in two stages a mapper phase and a reducer phase. In mapper phase first the text is tokenized into words then we form a key value pair with these words where the key being the word itself and value '1'. For example consider the sentence

“tring tring the phone rings”

In map phase the sentence would be split as words and form the initial key value pair as <tring,1>



- <tring,1>
- <the,1>
- <phone,1>
- <rings,1>

In the reduce phase the keys are grouped together and the values for similar keys are added. So here there are only one pair of similar keys ‘tring’ the values for these keys would be added so the out put key value pairs would be

- <tring,2>
- <the,1>
- <phone,1>
- <rings,1>

This would give the number of occurrence of each word in the input. Thus reduce forms an aggregation phase for keys.

## VI. CONCLUSION

Map-Reduce have become an important platform for a variety of data processing applications. Word Count Mechanisms in Map-Reduce frameworks such as Hadoop, suffer from performance degradations in the presence of faults. Word Count Map-Reduce, proposed in this paper provides an online, on-demand and closed-loop solution to managing these faults. The control loop in word count mitigates performance penalties through early detection of anomalous conditions on slave nodes. Anomaly detection is performed through a

novel sparse-coding based method that achieves high true positive and true negative rates and can be trained using only normal class (or anomaly-free) data. The local, decentralized nature of the sparse-coding models ensures minimal computational overhead and enables usage in both homogeneous and heterogeneous Map-Reduce environments.

## VII. REFERENCES

- [1] Samneet Singh and Yan Liu, "A Cloud Service Architecture for Analyzing Big Monitoring Data", ISSN 11007-0214/105/101pp55-70 Volume 21, Number 1, February 2016
- [2] JOSEPH A. ISSA, "Performance Evaluation and Estimation Model Using Regression Method for Hadoop WordCount", Received November 19, 2015, accepted December 12, 2015, date of publication December 18, 2015, date of current version December 29, 2015.
- [3] Yaxiong Zhao, Jie Wu, and Cong Liu, "Dache: A Data Aware Caching for Big-Data Applications Using the MapReduce Framework", ISSN 110070214/105/101pp39-50 Volume 19, Number 1, February 2014
- [4] Zhuoyao Zhang Ludmila Cherkasova, "Benchmarking Approach for Designing a MapReduce Performance Model", *ICPE'13*, April 21-24, 2013
- [5] Nikzad Babaii Rizvandi, Albert Y. Zomaya, Ali Javadzadeh Bolori, Javid Taheri, "On Modeling Dependency between MapReduce Configuration Parameters and Total Execution Time", 2012
- [6] Nikzad Babaii Rizvandi, Javid Taheri, Reza Moraveji, Albert Y. Zomaya, "On Modelling and Prediction of Total CPU Usage for Applications in MapReduce Environments", 2011
- [7] A. Baratloo, M. Karaul, Z. Kedem, and P. Wyckoff, "Charlotte: Meta computing on the Web," in *Proc. 9th Int. Conf. Parallel Distrib. Comput. Syst.*, 1996, pp. 1\_13.
- [8] J. Bent, D. Thain, A. C. Arpaci-Dusseau, R. H. Arpaci-Dusseau, and M. Livny, "Explicit control in the batch-aware distributed file system," in *Proc. 1st USENIX Symp. Netw. Syst. Design Implement. (NSDI)*, Mar. 2004, pp. 365\_378.
- [9] A. Fox, S. D. Gribble, Y. Chawathe, E. A. Brewer, and P. Gauthier, "Cluster-based scalable network services," in *Proc. 16th ACM Symp. Oper. Syst. Principles*, Saint-Malo, France, 1997, pp. 78\_91.
- [10] S. Ghemawat, H. Gobioff, and S.-T. Leung, "The Google file system," in *Proc. 19th Symp. Oper. Syst. Principles*, New York, NY, USA, 2003, pp. 29\_43.
- [11] S. Ibrahim, H. Jin, L. Lu, L. Qi, S. Wu, and X. Shi, "Evaluating MapReduce on virtual machines: The Hadoop case," in *Proc. Int. Conf. Cloud Comput.*, vol. 5931. 2009, pp. 519\_528.
- [12] J. Issa and S. Figueira, "Graphics performance analysis using Amdahl's law: IEEE/SCS SPECTS," in *Proc. Int. Symp. Perform. Eval. Comput. Telecommun. Syst.*, Ottawa, ON, Canada, 2010, pp. 127\_232.