

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 5, May 2018

I²mapreduce: Incremental Mapreduce for Evolving Big Data

Sayali Lonkar

M.E Student (CE), JSPMs Imperial college of Engg and Research, wagholi, Pune, Maharashtra, India

ABSTRACT: Many online data sets grow gradually as new entries are added and existing entries are deleted or changed. With the constant arrival of new data and updates, the results of data mining applications become obsolete over time. To keep away from the scenario it is necessary to refreshing mining results so that it can avoid the cost of re-computation from scratch. MapReduce for efficient iterative computations, it is too expensive to perform an entirely new large-scale MapReduce iterative job to timely accommodate new changes to the underlying data sets. The i2MapReduce is an extension to MapReduce that supports key-value pair level incremental processing and also more sophisticated iterative computation, which is widely used in data mining applications. Incremental interactive processing is a small number of updates which propagate to affect a large portion of intermediate states after a number of iterations. The technique helps in improving the job running time and reduces the running time of refreshing the results of big data.

KEYWORDS: Big Data, MapReduce, Incremental Processing, Iterative Computation

I. INTRODUCTION

Modern Internet applications have created a need to manage immense amounts of data quickly. For example devices and communication means like social networking sites, the amount of data produced by mankind is growing rapidly every year. It has become increasingly popular to mine such big data, which helps in taking business decisions or to provide better personalized good quality services. A large number of frameworks have been developed for big data analysis. MapReduce is one of the simple, generalized, framework used in production. Implementations of map-reduce enable many of the most common calculations on large-scale data to be performed on large collections of computers, efficiently and in a way that is tolerant of hardware failures during the computation. Here the main focus is on improving MapReduce technique.

The data lying in the servers of the company was just data until yesterday – sorted and filed. Suddenly, the slang Big Data got popular and now the data in the company is Big Data. The term covers each and every piece of the data that organization has stored till now. Big data involves the data produced by different devices and applications. Big data technologies are important in providing more accurate analysis, which may lead to more concrete decision-making resulting in greater operational efficiencies, cost reductions, and reduced risks for the business. In many situations, it is desirable to periodically refresh the mining computation in order to keep the mining results up-to-date.

Incremental processing is an advanced approach to refreshing mining results. Given the size of the input big data, it is very heavy weighted to return the entire computation from scratch. Incrementally processing the new data of a large data set, takes state as implicit input and combines it with new data. MapReduce [1] programming model is widely used for large scale and one-time data-intensive distributed computing, but it lacks for built-in support for the iterative process.

II. REVIEW OF LITERATURE

It is a framework for machine learning and data mining in the cloud. Graph Lab abstraction which naturally expresses the asynchronous, dynamic, graph-parallel computation while ensuring data consistency and achieving a high degree of

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 5, May 2018

parallel performance in the shared-memory setting. With the exponential growth in the scale of Machine Learning and Data Mining (MLDM) problems and increasing sophistication of MLDM techniques, there is an increasing need for systems that can execute MLDM algorithms efficiently in parallel on large clusters. Simultaneously, the availability of Cloud computing services like Amazon EC2 provide the promise of on-demand access to affordable large-scale computing and storage resources without substantial upfront investments. Unfortunately, designing, implementing, and debugging the distributed MLDM algorithms needed to fully utilize the Cloud can be prohibitively challenging requiring MLDM experts. To address race conditions, deadlocks, distributed state, and communication protocols while simultaneously developing mathematically complex models and algorithms. MLDM have focused on modeling the dependencies between data. By modeling data dependencies, then it is able to extract more signals from noisy data. Unfortunately, data parallel abstractions like MapReduce are not generally well suited for the dependent computation typically required by more advanced MLDM algorithms. Dynamic asynchronous graph parallel computation will be a key component in large-scale machine learning and data-mining systems, and thus further research into the theory and application of these techniques will help in defining the emerging field of Big Datalearning [1].

Parallel dataflow systems are an increasingly popular solution for analyzing large data volumes. They offer a simple programming abstraction based on directed acyclic graphs, and relieve the programmer from dealing with the complicated tasks of scheduling computation, transferring intermediate results, and dealing with failures. Most importantly, they allow dataflow programs to be distributed across large numbers of machines which is imperative when dealing with today's data volumes. While dataflow systems were originally built for tasks like indexing, filtering, transforming, or aggregating data, their simple interface and powerful abstraction have made them popular for other kinds of applications, like machine learning or graph analysis. Many of these algorithms are of iterative or recursive nature, repeating some computation until a condition is fulfilled. Naturally, these tasks pose a challenge to dataflow systems, as the flow of data is no longer acyclic. It focuses on two different kind of iteration. Bulk iteration which computes a completely new partial solution from the previous iteration's result, optionally using additional data sets that remain constant in the course of the iteration. Prominent examples are machine learning algorithms like Batch Gradient Descent and Distributed Stochastic Gradient Descent, many clustering algorithms (such as K-Means), and the well-known PageRank algorithm and another one is Incremental Iteration. The major drawbacks of the system are dataflow systems are practically inefficient for many iterative algorithms. The systems are, however, still required for other typical analysis and transformation tasks [2].

Hadoop technique for efficient data processing on large clusters. The need for highly scalable parallel data processing platforms is rising due to an explosion in the number of massive-scale data intensive applications both in industry (e.g., web-data analysis, click-stream analysis, network-monitoring log analysis) and in the sciences (e.g., analysis of data produced by massive-scale simulations, sensor deployments, high-throughput lab equipment). MapReduce is a well-known framework for programming commodity computer clusters to perform large-scale data processing in a single pass. A MapReduce cluster can scale to thousands of nodes in a fault-tolerant manner. Although parallel database systems may also serve these data analysis applications, they can be expensive, difficult to administer, and lack fault-tolerance for long-running queries. Hadoop, an open-source MapReduce implementation, has been adopted by Yahoo!, Facebook, and other companies for large-scale data analysis. The main drawback of the experiment is MapReduce framework does not directly support these iterative data analysis applications. Instead, programmers must implement iterative programs by manually issuing multiple MapReduce jobs and orchestrating their execution using a driver program. Even though much of the data may be unchanged from iteration to iteration, the data must be re-loaded and re-processed at each iteration, wasting I/O, network bandwidth, and CPU resources[3].

There is a growing demand for large-scale processing of unstructured data, such as text, audio, and image files. It is estimated that unstructured data is now accumulating in data centers at three times the rate of traditional transaction-based data. These environments often require data processing systems to absorb terabytes of new information every day while running a variety of complex data analytics. The data "deluge" presents major data management challenges, and non-relational information analysis is quickly emerging as a bedrock technology for these large-scale data processing

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 5, May 2018

efforts. The work carried out addresses the need for stateful dataflow programs that can rapidly sift through huge, evolving data sets. These data-intensive applications perform complex multi-step computations over successive generations of data inflows, such as weekly web crawls, daily image/video uploads, log files, and growing social networks. While programmers may simply re-run the entire dataflow when new data arrives, which is grossly inefficient, so increasing the result latency and squandering hardware resources and energy. Alternatively, programmers may use prior results to incrementally incorporate the changes. However, current large-scale data processing tools, such as Map-Reduce or Dryad, limits how programmers incorporate and use state in data-parallel programs. Straightforward approaches to incorporating state can result in custom, fragile code and disappointing performance [4].

The advances in sensing, storage, and networking technology have created huge collections of high-volume, high dimensional data. Making sense of these data is critical for companies and organizations to make better business decisions, and even bring convenience to the daily life. Recent advances in scientific computing, such as computational biology, have led to a flurry of data analytic techniques that typically require an iterative refinement process. Common to these proposed distributed frameworks, iterative updates are performed iteration by iteration. Specifically, an iterative update in iteration k is performed after all updates in iteration $k - 1$ have completed. Within the same iteration, the results from an earlier update cannot be used for a later update. While it is the traditional model for iterative computations, such a model does slow down the progress of the computation. First, the concurrent update model does not utilize the already retrieved partial results from the same iteration, which prolongs the iteration process and lead to slow convergence. Second, a synchronization step is required for every iteration. And it is a time consuming in a heterogeneous distributed environment [5].

In this paper, implementation of MapReduce runs on a large cluster of commodity machines and is highly scalable: a typical MapReduce computation processes many terabytes of data on thousands of machines. Programmers find the system easy to use: hundreds of MapReduce programs have been implemented and upwards of one thousand MapReduce jobs are executed on Google's clusters every day [6].

This part will get you started in thinking about designing and analyzing algorithms. It is intended to be a gentle introduction to how we specify algorithms, some of the design strategies we will use throughout this book, and many of the fundamental ideas used in algorithm analysis [7].

In this paper, we present Google, a prototype of a large-scale search engine which makes heavy use of the structure present in hypertext. Google is designed to crawl and index the Web efficiently and produce much more satisfying search results than existing systems. To engineer a search engine is a challenging task. Search engines index tens to hundreds of millions of Web pages involving a comparable number of distinct terms. They answer tens of millions of queries every day. Despite the importance of large-scale search engines on the Web, very little academic research has been done on them. Furthermore, due to rapid advance in technology and Web proliferation, creating a Web search engine today is very different from three years ago. This paper provides an in-depth description of our large-scale Web search engine the first such detailed public description we know of to date [8].

III. SYSTEM OVERVIEW

In proposed work, the map task in the iteration starts, before finishing all reduce task in the previous iteration. Map task should be started as soon as their input data are available. The main loop in the I2MapReduce implementation not requires the completion of previous iteration job before starting the next iteration job .The results in the avoidable synchronization overhead. I2MapReduce implementation starts a new job for each iteration and implement incremental processing.

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 5, May 2018

SYSTEM ARCHITECTURE:

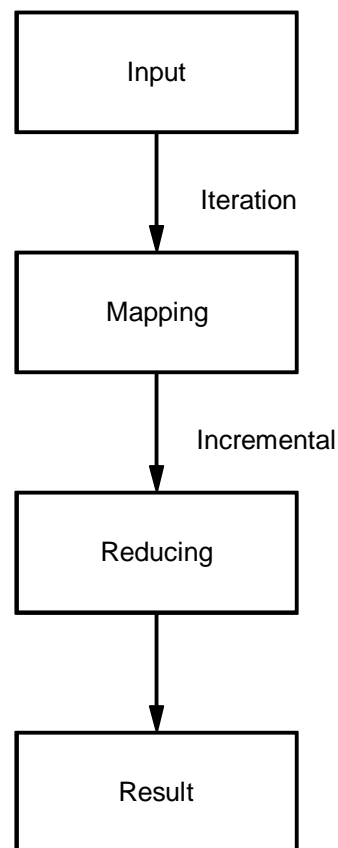


Fig1. Proposed System Architecture

IV.ALGORITHM

Algorithm -I:Query Algorithm in MRBG-Store

The MRBG-Store supports the preservation and retrieval of fine-grain MRBGraph states for incremental processing. We see two main requirements on the MRBG-Store. First, the MRBG-Store must incrementally store the evolving MRBGraph. Consider a sequence of jobs that incrementally refresh the results of a big data mining algorithm.

Algorithm-II: PageRank

PageRank is a well-known iterative graph algorithm for ranking web pages. It computes a ranking score for each vertex in a graph. After initializing all ranking scores, the computation performs a MapReduce job per iteration.

Algorithm-III: K-means Clustering

K-means is a commonly used clustering algorithm that partitions points into k clusters

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 5, May 2018

V. COMPARATIVE RESULTS

Stages	PlainMR	IterMR	I2MR
Map	1200s	850s	450s
Reduce	550s	400s	250s

Table 1.comparative result

VI.RESULT

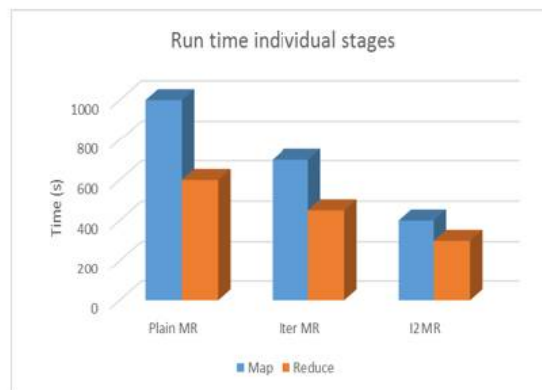


Fig.2 Graph

VII. CONCLUSION

The conclusion of the survey is the following, i2MapReduce that explicitly supports iterative processing of large relational data and addresses all problems in the implementation of MapReduce of iterative processing. Provides a framework for programmers to explicitly model iterative algorithms and proposes the concept of persistent tasks to perform iterative calculations to repeatedly avoid creating, destroying and scheduling tasks. Facilitates the asynchronous execution of activities within the same iteration, to speed up processing speed.

REFERENCES

- [1] Y. Low, D. Bickson, J. Gonzalez, C. Guestrin, A. Kyrola, and J. M. Hellerstein, "Distributed graphlab: a framework for machine learning and data mining in the cloud," PVLDB, vol. 5, no. 8, pp. 716–727, 2012.
- [2] S. Ewen, K. Tzoumas, M. Kaufmann, and V. Markl, "Spinning fast iterative data flows," PVLDB, vol. 5, no. 11, pp. 1268–1279, 2012.
- [3] Y. Bu, B. Howe, M. Balazinska, and M. D. Ernst, "Haloop: efficient iterative data processing on large clusters," PVLDB, vol. 3, no. 1- 2, pp. 285–296, 2010.
- [4] D. Logothetis, C. Olston, B. Reed, K. C. Webb, and K. Yocum, "Stateful bulk processing for incremental analytics," in Proc. of SOCC '10, 2010.
- [5] Y. Zhang, Q. Gao, L. Gao, and C. Wang, "experiment to Accelerate large-scale iterative computation" vol. 3, no. 1- 2, pp. 281–294, 2013.
- [6] J. Dean and S. Ghemawat, "Mapreduce: simplified data processing on large clusters," in Proc. of OSDI '04, 2004.
- [7] Cormen, T.H., Stein, C., Rivest, R.L., Leiserson, C.E.: Introduction to Algorithms, 2nd edn. McGraw-Hill Higher Education (2001).
- [8] Brin, S., Page, L.: The anatomy of a large-scale hypertextual web search engine. Comput Networks ISDN 30, 107–117 (1998).