



INTERNATIONAL JOURNAL OF ADVANCE RESEARCH, IDEAS AND INNOVATIONS IN TECHNOLOGY

ISSN: 2454-132X

Impact factor: 4.295

(Volume 3, Issue 6)

Available online at www.ijariit.com

Proficient Analysis of Mining Big Data Using Map Reduce Framework

Ajinkya Molke

D. Y Patil Engineering College, Pune, Maharashtra

molke.ajinkya@gmail.com

Abstract: Information now stream from everyday life from telephones and charge cards and TVs and PCs; from the framework of urban areas from sensor-prepared structures, trains, transports, planes, extensions, and production lines. The information stream so quick that the aggregate gathering of the previous two years is currently a zettabyte. This colossal volume of information is known as large information. Huge Data alludes to advancements and activities that include information that is excessively various, quick changing or gigantic for traditional advances, aptitudes and framework to address productively. Said in an unexpected way, the volume, speed or assortment of information is excessively extraordinary. The volume of information with the speed it is created makes it troublesome for the present processing foundation to deal with enormous information. To conquer this downside, enormous information handling can be performed through a programming worldview known as MapReduce. Commonplace, execution of the mapReduce worldview requires arranged connected stockpiling and parallel preparing. Hadoop and HDFS by apache are generally utilized for putting away and overseeing huge information. In this exploration paper the creators recommend different strategies for taking into account the issues close by through MapReduce structure over HDFS. MapReduce strategy has been learned at in this paper which is required for actualizing Big Information examination utilizing HDFS and minimizes time.

Keywords: Big Data, MapReduce, Hadoop, HDFS, Zettabyte.

I. INTRODUCTION

Big data is a popular term used to describe the exponential growth and availability of data, both structured and unstructured. Big data may be important to business and society as the Internet has become. Big Data is so large that it's difficult to process using traditional database and software techniques. More data may lead to more accurate analyses. More accurate analyses may lead to more confident decision making. And better decisions can mean greater operational efficiencies, cost reductions and reduced risk

Today enormous measure of computerized information is being gathered in numerous vital ranges, including web based business, interpersonal organization, fund, social insurance, training, and environment. It has turned out to be progressively famous to mine such enormous information keeping in mind the end goal to pick up bits of knowledge to help business choices or to give better customized, higher quality administrations. Lately, an extensive number of figuring structures [1], [2], [3], [4], [5], [6] have been produced for enormous information investigation. Among these structures, MapReduce(with its open-source executions, for example, Hadoop) is the most broadly utilized as a part of generation in view of its straightforwardness, consensus, and development. This paper will concentrate on enhancing MapReduce. Enormous information is continually developing. As new information and upgrades are being gathered, the info information of a major information mining calculation will step by step change, and the registered outcomes will get to be distinctly stale and old after some time. Much of the time, it is alluring to intermittently invigorate the mining calculation with a specific end goal to

stay up with the latest. For instance, the PageRank calculation processes positioning scores of site pages in light of the web chart structure for supporting web look. In any case, the web diagram structure is continually developing; Web pages and hyper-connections are made, erased, and upgraded. As the hidden web chart develops, the PageRank positioning outcomes bit by bit get to be distinctly stale, conceivably bringing down the nature of web pursuit. In this manner, it is attractive to invigorate the PageRank calculation routinely. Incremental handling is a promising way to deal with reviving mining comes about. Given the extent of the information huge information, it is regularly extremely costly to rerun the whole calculation sans preparation. Incremental handling abuses the way that the info information of two ensuing calculations an and B are comparative. Just a little division of the information has changed. The thought is to spare states in calculation A, re-utilize A's states in calculation B, and perform re-calculation just for states that are influenced by the changed information. The acknowledgment of this standard with regards to the MapReduce processing system is examined. Various past reviews (counting Percolator [7], CBP [8], and Naiad [9]) have taken after this standard and planned new programming models to bolster incremental handling.

Then again, Incoop [10] stretches out MapReduce to bolster incremental preparing. In any case, it has two principle constraints. To begin with, Incoop underpins just undertaking level incremental preparing. That is, it spares and reuses states at the granularity of individual Map and Reduce errands. Every assignment regularly forms a substantial number of key-value sets (kv-sets). In the event that Incoop identifies any information changes in the contribution of an assignment, it will rerun the whole errand. While this approach effectively influences existing MapReduce highlights for state investment funds, it might bring about a lot of repetitive calculation if just a little division of kv-sets has changed in an assignment. Second, Incoop underpins just a single stride calculation, while critical mining calculations, for example, PageRank, require iterative calculation. Incoop would regard every cycle as a different MapReduce work. Nonetheless, a little number of info information changes may bit by bit spread to influence a substantial bit of middle of the road states after various cycles, bringing about costly worldwide re-calculation a while later. This paper proposes i2MapReduce, an augmentation to MapReduce that backings fine-grain incremental preparing for both one stage and iterative calculation. Contrasted with past arrangements, i2MapReduce joins the accompanying three novel components:

Fine-grain incremental preparing utilizing MRBG-Store. Dissimilar to Incoop, i2MapReduce support kv-pair level fine-grain incremental handling so as to minimize the measure of re-calculation however much as could be expected. This paper models the kv-pair level information stream and information reliance in a MapReduce calculation as a bipartite diagram, called MRBGraph. A MRBG-Store is intended to protect the fine-grain states in the MRBGraph and bolster productive inquiries to recover fine-grain states for incremental handling.

Universally useful iterative calculation with unobtrusive

augmentation to MapReduce API. Past work proposed iMapReduce [6] to effectively bolster iterative calculation on the MapReduce stage. Be that as it may, it targets sorts of iterative calculation where there is a balanced/all-to-one correspondence from Reduce yield to Map include. In examination, this paper gives broadly useful support, including coordinated, as well as one-to-numerous, many-to-one, and many-to-numerous correspondence. The framework upgrades the Map API to permit clients to effectively express circle invariant structure information, and proposes a Project API capacity to express the correspondence from Reduce to Map. While clients need to somewhat alter their calculations with a specific end goal to take full preferred standpoint of i2MapReduce, such adjustment is unobtrusive contrasted with the push to re-execute calculations on a totally unique programming worldview, for example, in Percolator [7], CBP [8], and Naiad [9].

Incremental preparing for iterative calculation. Incremental iterative preparing is significantly more difficult than incremental one-stage handling in light of the fact that even a little number of overhauls may engender to influence an expansive part of middle of the road states after various cycles [1]. To address this issue, this paper proposes to reuse the merged state from the past calculation and utilize a change spread control (CPC) system. This paper likewise upgrades the MRBG-Store to better bolster the get to designs in incremental iterative handling. As far as anyone is concerned, i2MapReduce is the principal MapReduce-based arrangement that proficiently underpins incremental iterative calculation.

Analysts executed i2MapReduce by altering Hadoop-1.0.3. Analysts assess i2MapReduce utilizing a one-stage calculation (A-Priori) and four iterative calculations (PageRank, SSSP, Kmeans, GIM-V) with various calculations attribute. Exploratory outcomes on Amazon EC2 demonstrate critical execution upgrades of i2MapReduce contrasted with both plain and iterative MapReduce performing re-calculation. For instance, for the iterative PageRank calculation with 10 percent information changed, i2MapReduce enhances the run time of re-calculation on plain MapReduce by an eight overlay speedup [1]. This paper utilizes an adjusted adaptation of the A-priori calculation, named as Top K rules, which finds and suggests just the best K guidelines of the framework, not considering the repetitive principles, and giving just the tenets which are better to describe the framework conduct.

II. SURVEY RELATED DETAILS

Incoop: MapReduce for Incremental Computations

In this concept, the execution and evaluation of a general MapReduce framework, Incoop, for incremental computations is described. Incoop identifies changes in the inputs and performs the automatic upgrade to the outputs by applying fine-grained re-use mechanism. Incoop, allows existing MapReduce programs which are not designed for incremental processing, to execute clearly in an incremental way. Computations in Incoop can respond automatically and efficiently to modifications in their input data by reusing intermediary results from past computations, and incrementally update the output based on the changes in the

input.

Execution of Incoop by extending the Hadoop framework, and evaluated it through a diversity of applications, including two case studies of higher-level services: incremental query and log processing systems. Our results show important performance improvements without varying a single line of application code.

Limitations:

- Supports only task level incremental processing.
- Supports only one step computation.

IMapReduce

Iterative computation is generally known for many applications such as data mining, online social network analysis, graph analysis, etc. These iterative applications generally include large data sets consisting huge data records. This indicates the requirement of distributed computing frameworks for analysing large data sets on a cluster of machines. As a solution for such problems, iMapReduce permits users to define the iterative computation with the distinguished map and reduce functions, and lends the support for automatic iterative processing within a single job. More precisely, iMapReduce remarkably enhances the performance of iterative executions by reducing the expense of resulting new MapReduce tasks again and again, thus avoiding the shuffling of static data, and hence permitting asynchronous execution of map tasks.

Limitations:

- The operations to each iteration are the same.
- Unnecessary communication overhead.
- Unnecessary synchronization overhead.

Past work Incoop, [10] supports just assignment level incremental handling. That is, it spares and reuses states at the granularity of individual Map and Reduce undertakings. Every errand normally forms an extensive number of key-value sets (kv-sets). On the off chance that Incoop distinguishes any information changes in the contribution of an assignment, it will rerun the whole undertaking. While this approach effortlessly influences existing MapReduce highlights for state investment funds, it might cause a lot of excess calculation if just a little part of kv-sets have changed in an undertaking.

Past work proposed iMapReduce, [6] to proficiently bolster iterative calculation on the MapReduce stage. In any case, it targets sorts of iterative calculation where there is a balanced/all-to-one correspondence from Reduce yield to Map enter.

Past work Incoop, [10] supports incremental one-stage handling. Incoop would regard every cycle as a different MapReduce work. Be that as it may, a little number of information changes may progressively engender to influence an extensive part of halfway states after various cycles, bringing about costly worldwide re-calculation a short time later.

In the past work [1] the specialists have grown quick methods

for developing information, and its mapping. Be that as it may, the diminishment part still needs change. In the work, they have portrayed different mapping and decreasing strategies, however in the event that lessening is not improved then the general framework productivity is low and may prompt to a moderate reaction for a constant framework. Past works have taking after issues:

Does not underpins key-value pair level incremental preparing and underpins just errand level incremental preparing. Does not support General-reason iterative calculation and just backings coordinated/all-to-one correspondence from Reduce yield to Map input.

Does not support incremental handling for iterative calculations and just backings incremental one-stage preparing. Speed of the mining procedure is low.

III. PROPOSED WORK

Current paper proposes, a framework which conquers the downside of slower diminishment times, and uses a technique which decreases the info information speedier when contrasted with any proposed calculation, consequently enhancing the general productivity of the framework. The proposition utilizes a changed form of the A-priori calculation, named as Top K rules, which finds and suggests just the best K standards of the framework, not considering the repetitive principles, and giving just the guidelines which are better to describe the framework conduct and giving the most ideal proposals for the framework. This technique will enhance the general speed and precision of the administer mining procedure and make the whole Map Reduce structure perform progressively with most elevated amount of exactness. Proposed approach works in the accompanying way,

Step 1: Gathering of Advancing Datasets

The advancing datasets will be gathered for mapping what's more, lessening.

Step 2: Improvement of Mapping Procedure

MapReduce takes into consideration dispersed handling of the guide and lessening operations. Given that every mapping operation is autonomous of the others, all maps can be performed in parallel however practically speaking this is constrained by the quantity of free information sources and additionally the quantity of CPUs close to every source.

Delineate Maps input key/esteem sets to an arrangement of halfway key/esteem sets. Maps are the individual undertakings which change input records into halfway records. The changed middle of the road records require not be of an indistinguishable sort from the information records. A given info combine may guide to zero or many yield sets. Guide () is run precisely once for each K1 key esteem, creating yield sorted out by key qualities K2.

Step 3: Advancement of Reduction Technique

An arrangement of “Reducers” can play out the decrease stage, gave that all yields of the guide operation that have a similar key are exhibited to a similar reducer in the meantime, or that the diminishment capacity is affiliated. Reduce () is run precisely once for every K2 key esteem created by the Map step.

Intelligent View of MapReduce Process:

The Map and Reduce elements of MapReduce are both characterized concerning information organized in (key, esteem) sets. Delineate one sets of information with a sort in one information area, and returns a rundown of sets in an alternate space:

Map(k1,v1) list(k2,v2)

The Map capacity is connected in parallel to each pair in the info dataset. This creates a rundown of sets for every call. After that, the MapReduce system gathers all sets with a similar key from all rundowns and gatherings them together, making one gathering for every key.

The Reduce capacity is then connected in parallel to every gathering, which thusly creates an accumulation of qualities in a similar space: Reduce (k2, list (v2)) list (v3)

Each Reduce call commonly creates it is possible that one esteem v3 or an unfilled return, however one call is permitted to return more than one esteem. The profits of all calls are gathered as the fancied outcome list.

Therefore the MapReduce structure changes rundown of (key, esteem) sets into a rundown of qualities.

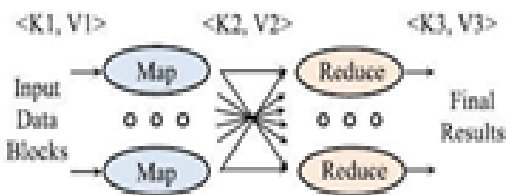


Fig. 1. MapReduce Computation

Step 4: Change in Diminishment Procedure Utilizing Top K

Rules Change is done in diminishment procedures utilizing Top K Rules calculation which is the alters variant of from the earlier calculation. The Top K Rules calculation acts as takes after:

The TopKRules calculation takes as info an exchange database, a number k of tenets that the client needs to find, and the minconf threshold.

The calculation fundamental thought is the accompanying. TopKRules first sets an interior minsup variable to 0. At that point, the calculation begins hunting down standards. When a

manage is discovered, it is added to a rundown of standards L requested by the support. The rundown is utilized to keep up the top-k rules found as of not long ago. When k substantial standards are found, the inner minsup variable is raised to the support of the run with the most reduced support in L. Raising the minsup esteem is utilized to prune the look space when scanning for more principles. From there on, every time a substantial lead is found, the manage is embedded in L, the standards in L not regarding minsup any longer are expelled from L, and minsup is raised to the estimation of the minimum intriguing guideline in L. The calculation keeps scanning for more principles until no govern are discovered, which implies that it has found the top-k rules.

To look for guidelines, TopKRules does not depend on the traditional two stages way to deal with produce rules since it would not be effective as a top-k calculation (as clarified in the presentation). The technique utilized by TopKRules rather comprises of producing tenets containing a solitary thing in the predecessor and a solitary thing in the subsequent. At that point, every control is recursively developed by adding things to the forerunner or subsequent. To choose the things that are added to a govern to develop it, TopKRules examines the exchanges containing the administrator to discover single things that could extend its left or right part. Two procedures for extending rules in TopKRules are left extension and right development. These procedures are connected recursively to investigate the inquiry space of affiliation guidelines.

Another thought joined in TopKRules is to attempt to create the most encouraging principles first. This is on account of if tenets with high support are discovered before, TopKRules can raise its inside minsup variable speedier to prune the inquiry space. To play out this, TopKRules utilizes an inner variable R to store every one of the guidelines that can be extended to have a possibility of discovering more substantial tenets. TopKRules utilizes this set to decide the guidelines that are the destined to deliver substantial standards with a high support to raise minsup all the more rapidly and prune a bigger part of the pursuit space.

Step 5: Result Analysis and Comparison

The Processing of time will be broke down and will be compared with existing outcomes.

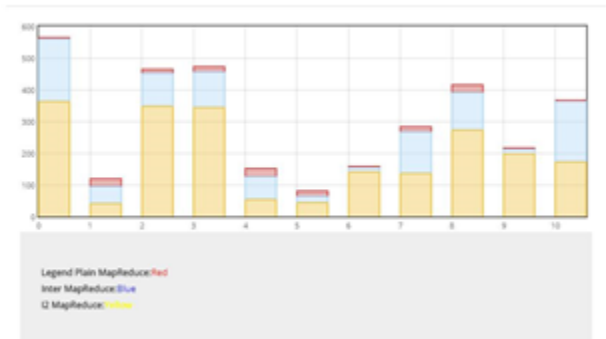


Fig. 2: Time Reduced Maps

IV. OBJECTIVES

The fundamental point of the Paper is to study, model and execute execution upgrades for Big Data systems. The destinations are recorded beneath.

- To cut down the Response Time required for Data Processing.
- To scale down the amount of Computation needed to be performed.

The standpoint of this venture work is entirely limited to take in the current Hadoop MapReduce framework and fortifying Hadoop by upholding incremental MapReduce for Big Data preparing. At long last, comes about created will be contrasted and the consequences of existing frameworks.

V. HADOOP SYSTEM

Hadoop is a structure which grants stockpiling and huge information preparing in circulated way over bunches of figuring machines with the assistance of programming models. Hadoop comprises of four sorts to be specific Datanodes, NameNode, Task Tracker and Job Tracker. In disseminated document framework gave by these hubs, the Job Tracker deals with the employments and Task Tracker runs the undertakings.

VI. ADVANTAGES OF HADOOP

Ability to store and process extensive pieces of information Since extensive volumes and assortments of data is consistently expanding, especially from online business and web-based social networking, fast stockpiling and preparing turns into a critical element.

High Figuring Capacity: Hadoop's dispersed registering model aides in preparing enormous information quicker.

Adaptability: Unlike customary databases, there is no compelling reason to preprocess information before putting away it. Any measure of information can be put away and later chose how to use it. Such information incorporates content, pictures and recordings.

Low Cost: The open-source system is free and uses reasonable PC equipment to store enormous measures of information.

Scalability: More information can be effortlessly taken care of by including hubs and thus improving the framework. Little fixation is required over organization.

VII. HDFS ARCHITECTURE

HDFS [4], the Hadoop Distributed File System, is a distributed file system designed to hold very large amounts of data (petabytes or even zettabytes), and provide high through put access to this information. Files are stored in a redundant fashion across multiple machines to ensure their durability to failure and high availability to very parallel applications. In particular, it ensures Big Datas durability to failure and high availability to parallel applications. Figure 1 shows HDFS has

a master/slave architecture.

An HDFS [4] cluster consists of a single NameNode, a master server that manages the file system namespace and regulates access to files by clients. In addition, there are a number of DataNodes, usually one per node in the cluster, which manage storage attached to the nodes that they run on. HDFS exposes a file system namespace and allows user data to be stored in files. Internally, a file is split into one or more blocks and these blocks are stored in a set of DataNodes.

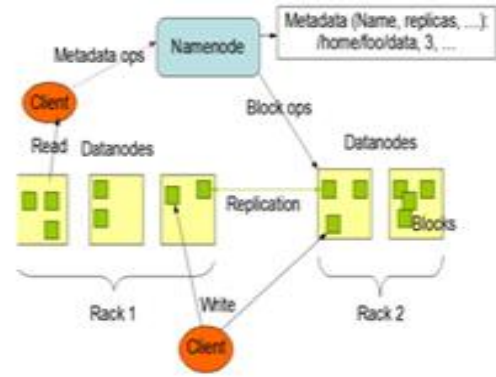


Fig. 3: HDFS Architecture

it into little settled size squares and stores numerous duplicates of every piece on bunch hub plate HDFS takes after the ace slave design The ace hub is known as NameNode which deals with the document framework and gives data to the customer. There are a few of slave hubs called DataNode which is utilized to store real information in units of pieces. The NameNode oversee mapping table which maps information pieces of DataNode keeping in mind the end goal to prepare compose and read demand of HDFS customers.

The DataNode likewise perform operations, for example, shutting, renaming and opening document and catalogs. If there should arise an occurrence of NameNode disappointment, optional NameNode work occasionally and spare a duplicate of the metadata. The DataNode stores the information squares and executes direction like creation, erasure, information substitution and replication from the NameNode.

A DataNode occasionally send a flag and get guideline from NameNode. A flag is additionally helpful for NameNode. A flag is additionally valuable for NameNode to recognize network with its DataNode. On the off chance that a flag is not got by NameNode in the designed timeframe, it denote the hub down.

VIII. ADVANTAGES OF MAPREDUCE

The benefits of MapReduce writing computer programs are Scalability, Cost-powerful arrangement, Flexibility in environment, Fast Execution, Security and Authentication, Parallel Preparing, Availability and strong nature. These focal points of MapReduce approach are regularly exploited to

execute different calculations.

IX. INCREMENTAL PROCESSING

The In incremental processing, the framework has an enormous vault of information that is redesigned as the approaching information is handled. The measure of the approaching information is little when contrasted with the span of the current information. Incremental preparing is across the board utilized as a part of regions, for example, seeking, ordering. At the point when another substance is included as opposed to figuring the entire information once more, the framework engenders the progressions sparing the execution time. Incremental handling of information is done in lumps so that the information adjustment is done just for the important ranges preferences of MapReduce writing computer programs are Scalability, Cost-viable arrangement, Flexibility in environment, Fast Execution, Security and Authentication, Parallel preparing, Availability and versatile nature. These focal points of MapReduce approach are regularly exploited to actualize different calculations.

X. INCREMENTAL MAPREDUCE

Incremental processing is a urging way to deal with conquer the weaknesses of conventional MapReduce. Incremental MapReduce is an augmentation to MapReduce which ties fine-grain incremental handling to both one-stage calculation and iterative calculation. Incremental MapReduce is the interesting MapReduce that effectively support incremental iterative calculation.

XI. SYSTEM ARCHITECTURE

In framework execution we need to characterize essential squares of framework with the goal that we can make modularized perspective of framework as appeared beneath figure. The four primary commitments made by our proposed work.



Fig. 4: System Architecture

Are Clarified as Takes After
PageRank
K-means
MapReduce

Incremental MapReduce

The venture characterizes and tackles the issue of preparing enormous information by utilizing the incremental MapReduce. For the calculation of incremental MapReduce a few calculations must be executed which ought to be run iteratively.

Above all else, client submits base document and its iterative records. At that point PageRank calculation, which is utilized to rank the pages, takes input records and produces the outcome. Aftereffects of the PageRank are given as contribution to the k-implies calculation. The K-implies calculation, which is utilized for the grouping, takes the info and gives the Output as Cluster.

MapReduce takes the input (base record and its iterative documents) forms on info documents and gives us the yield, which is put away in the HDFS (Hadoop Distributed File System). Incremental MapReduce takes contribution as yield consequences of the k-implies bunching calculation and goes ahead with it. In this way, the aftereffects of incremental MapReduce are additionally put away in HDFS. Correlation of the consequences of MapReduce Incremental MapReduce is done from the HDFS where the outcomes are put away.

XII. MATHEMATICAL MODEL

Where S=system, C=cache manager, Mi= no. of modules, Fi=no. of functions, Pi=no. of processors The system can be

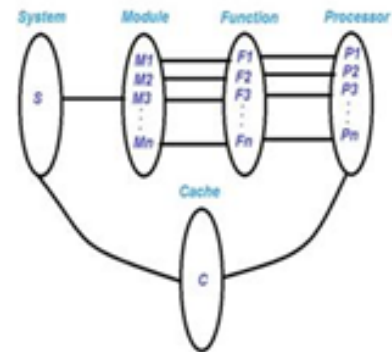


Fig. 5: Mathematical Model

Represented as
 $S = fM_1, M_2, M_3, M_4g$
 Where $M_1 = \text{split}, M_2 = \text{Map}, M_3 = \text{Reduce}, M_4 = \text{Cache Store}$ as
 $S = fM_i > 01 M_i < 04g$
 $S = ffc/01g$
 Where S is system
 Now consider another set
 $Q: F(s)$
 InputQuery it set of InputQuery: Which can be represented?
 as
 $Q = ffQ_1, Q_2, Q_3 \dots Q_n g$
 Where Q=Query, F=Functionality
 Then Q_n is Initilize to P_n Where P_n is Number of Processes

Qn!Pn

Pn= ff01 Pn ng

XIII. IMPLEMENTATION

Many different implementations of the MapReduce inter-face are possible. The right choice depends on the environment. For example, one implementation may be suitable for a small shared-memory machine, another for a large NUMA multi-processor, and yet another for an even larger collection of networked machines.

Large clusters of commodity PCs connected together with switched Ethernet. In our environment:

- (1)Machines are typically dual-processor x86processors running Linux, with2-4GB of memory per machine.
- (2)Commodity networking hardware is used typically either 100 megabits/second or 1 gigabit/second at the machine level, but averaging considerably less in overall bisection bandwidth.
- (3)A cluster consists of hundreds or thousands of ma-chines, and therefore machine failures are common.
- (4)Storage is provided by inexpensive IDE disks attached directly to individual machines. A distributed le system developed in-house is used to manage the data stored on these disks. The le system uses replication to provide availability and reliability on top of unreliable hardware.
- (5)Users submit jobs to a scheduling system. Each job consists of a set of tasks, and is mapped by the scheduler to a set of available machines within a cluster.

XIV. SOFTWARE REQUIREMENT SPECIFICATION

Software	OS	Windows and Linux
Hadoop		Hadoop 1.2.1 and Thrity Party
Java		Hadoop Sun JDK
Visual studio IDE		Microsoft .NET technology

Fig. 6: Experiment Setup

XV. PERFORMANCE MEASURES USED

In this method we measure the performance of MapReduce on two computations. One computation searches through approximately one terabyte of data looking for a particular data. The other computation sorts approximately one tera byte of data. These two programs are representative of a large subset of the real programs written by users of MapReduce one class of programs shuffles data from one representation to another class extracts a small amount of interesting data from a large dataset.

XVI. RESULT TABLES

The consequence of calculation will be broke down and will be contrasted with existing outcomes.

XVII. CONCLUSION

We are in a period of Big Data. For the capacity and preparing of huge information Hadoop innovation is utilized. HDFS is utilized for the capacity and MapReduce for the calculation. Incremental MapReduce is a MapReduce based system planned for incremental iterative calculations. Incremental MapReduce support more perplexing state-to-structure connections. Likewise it underpins huge information sets of terabytes and petabytes.

The principal show utilizes i2MapReduce, which joins a fine-grain incremental motor, a universally useful iterative model, and an arrangement of successful systems for incremental iterative calculation. The new model uses a changed rendition of the A-priori calculation, named as Top K rules, which finds and suggests just the best K principles of the framework. Contrasted and the principal show, the new model is significantly more proficient and accomplished the palatable execution also. The fundamental target of this paper was to toss some light on the proposed work. It gives a promising system to enhance the general speed and exactness of the lead mining procedure and make the whole Map Reduce structure perform progressively with largest amount of precision by utilizing Top K Rules.

REFERENCES

1. Yanfeng Zhang, Shimin Chen, Qiang Wang, and Ge Yu, "i2MapReduce: Incremental MapReduce for Mining Evolving Big Data", IEEE Transactions On Knowledge And Data Engineering, Vol. 27, No. 7, July 2015.
2. M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica, "Resilient distributed datasets: A fault-tolerant abstraction for, in-memory cluster computing," in Proc. 9th USENIX Conf. Netw.Syst. Des. Implementation, 2012, p. 2.
3. G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski, "Pregel: A system for large-scale graph processing," in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2010, pp. 135146.
4. Y. Bu, B. Howe, M. Balazinska, and M. D. Ernst, "Haloop: Efficient iterative data processing on large clusters," in Proc. VLDB Endowment, 2010, vol. 3, no. 12, pp. 285296.
5. J. Ekanayake, H. Li, B. Zhang, T. Gunarathne, S.-H. Bae, J.Qiu, and G. Fox, "Twister: A runtime for iterative mapreduce," in Proc. 19th ACM Symp. High Performance Distributed Comput., 2010, pp. 810818.
6. Y. Zhang, Q. Gao, L. Gao, and C. Wang, "imapreduce: A distributed computing framework for iterative computation," J. Grid Comput., vol. 10, no. 1, pp. 4768, 2012.
7. D. Peng and F. Dabek, "Large-scale incremental processing using distributed transactions and notifications," in Proc. 9th USENIX Conf. Oper. Syst. Des. Implementation, 2010, pp. 115.
8. D. Logothetis, C. Olston, B. Reed, K. C. Webb, and K. Yocum, "Stateful bulk processing for incremental analytics," in Proc. 1st ACM Symp. Cloud Comput., 2010, pp. 5162.
9. D. G. Murray, F. McSherry, R. Isaacs, M. Isard, P. Barham, and M. Abadi, "Naiad: A timely dataflow system," in Proc. 24th ACM Symp. Oper. Syst. Principles, 2013, pp. 439455.

10. P. Bhatotia, A. Wieder, R. Rodrigues, U. A. Acar, and R. Pasquin, "Incoop: Mapreduce for incremental computations," in Proc. 2nd ACM Symp. Cloud Comput., 2011, pp. 7:17:14.