

A Hash Tree Framework for BigData Processing

Dr. A. Suresh Babu, Associate Professor, JNTU-A-Anantapur-India, asureshjntu@gmail.com

S. Sireesha, Student (M.Tech), JNTU-A-Anantapur-India, selapareddisireesha@gmail.com

Abstract: Big Data is colossal quantity of data which is used to scrutinize the data and heap data. Though there are files, Database and Relational Database which these are not used to store the unstructured data so Big Data came into existence. There are frameworks such as Map reduce and Spark; no functionality being exposed so that code cannot be overridden. The HDM (Hierarchically Distributed Data Matrix) came into continuation which is reusable for subsequent development and natively compassable by using dependencies but there is no enhancement in the processing speed. In attendance we improved Hierarchical Distributed Data Matrix i.e.; HDM+ (HASH TREE Framework for BIGDATA processing) which leads to decrease the convolution so that processing speed increases and the execution time reduces. Besides by using a hashing method we can upload and store any type of data i.e. Heterogeneous data Fault-Tolerance can be done by using Self-healing.

Keywords — *Big Data Processing, Parallel programming, Functional Programming, Distributed Systems, Hash tree, System Architecture.*

I. INTRODUCTION

Now a day's Big Data had became a most admired term in this world and without this ; no data is stored if a more amount of data is used to store and access then it is used because to store such an amount of data the DBMS and Relational Data Base are not enough . Previously there is Data Base and relational data it supports only the structured data either only text or images etc. Big Data [5] supports Semi-Structured i.e.; both Unstructured and structured. Face book is the best example for both unstructured and structured data. Big Data is a well-liked term in our daily life and worn to symbolize the growth rate, immense data and used to repossess data. Over the period of 10 years ,the representation of Map Reduce[2] had became the real and standard .It has been used generally become an admired method to organize and utilize the power of large clusters. The primary topic of the Framework moves investigation to the information as opposed to moving the information to a framework to dissect it. It enables software engineers to insightfulness in information driven approach where they can focus on applying change to set of information dealings while the data of scattered achievement and adaptation to internal failure are clearly overseen by the structure. Be that as it may, as of late, with raising applications' apportionments in the information examination area, different regions of the Hadoop [4] structure have been unsurprising and in this way we see an unparalleled enthusiasm to leave upon these difficulties with new arrangements which constituted another signal of for the most part space specific, improved huge information allotment stages .Several systems like Spark , Flink , Pregel and Storm have been exhibited to take part in ever better datasets on utilizing scattered groups of ware gear. These structures significantly lessen the convolution of growing applications and huge information

programs. In authenticity some certifiable situations require amalgamation of various huge information and pipelining. There is more summonses while applying huge information learning. In introduce enormous information arrangement, for example, Map Reduce and Spark[3], there is no fitting method to appropriate and portrayal a sent and all around tuned online segment to facilitate engineers. Accordingly, there is huge and even concealed repetitive change in huge information applications.

II. RELATED WORK

A. *Alma Reduce simplified data processing on large clusters. Jeffrey Dean and Sanjay Ghemawat.*

Map reduce[2] is an influence show and a lin-Oked achievement for regulation and creating colossal datasets that is submissive to a wide assortment of certifiable errands. Clients symbolize the working out regarding a guide and decrease capacities, and the essential runtime framework mechanically parallelizes the calculation crosswise over huge scale groups of utensils, handles contraption disappointments, and calendars between machine report to make capable utilization of plates and systems. In Map Reduce, customers can physically portray Combinators for each mapper to apply depict unite in a Map Reduce work In Spark, customers can use Aggregator API rather than customary group By and lessen exercises to describe the gathering errand before rearranging information over the bunch[1]. Those are on the whole inductions of the nearby collection enhancement, software engineers which are established inside the Google over the previous four years to discover supplementary than ten thousand confined Map Reduce assignments. On a whole

consistently on Google in excess of twenty peta bytes for each day.

B. A Fault-Tolerant Abstraction for In-Memory Cluster Computing

In this paper we present hard-wearing scattered Datasets[5] (Resilient Distributed Datasets-RDD), a disseminated reminiscence generalization that lets programmers achieve in reminiscence calculation on an colossal huddle in a fault-tolerant loom. Iterative algorithms and interactive data mining tools are two types of applications in both cases, keeping data in memory can perk up concert by a categorize of magnitude in RDDs.[9] To accomplish adaptation to internal failure capably, RDDs bear the cost of an obliged type of mutual memory, in view of coarse-grained changes marginally fine-grained updates to shared state. In any case, we demonstrate that RDDs are sufficiently open to restrict a wide class of calculations, including ongoing master programming models for iterative employment, for example, Pregl and creative applications that these models don't assistant. We ordered the RDDs in a coordination which is called Spark, which we familiarize the entire time of variety to client benchmarks and applications. Draw on the reliance chart and the reproducibility of HDM occupations, the current HDM execution motor gives adaptation to non-critical failure to information handling by utilizing pushing heredity, in which the lost or fizzled information allotments would be re-process from its folks or progenitors in the information reliance diagram. Pushing ancestry is an outstanding procedure that is likewise utilized as part of present day information concentrated structures, for example, Spark[3], Tachyon, Nectar and BAD-FS. Later on, we want to include more adaptation to non-critical failure components, for example, snapshotting and replication to help the necessities for various sorts of consumption.

C. AUTONOMIC COMPUTING TOWARDS A SELF-HEALING SYSTEM SHAREE'S LASTER, B.S.AYODEJI O.OLATUNJI

Self-healing in Information Technology describes a system or device that has the ability to perceive that it isn't working legitimately, without human contribution make the basic changes in accordance with restore itself to normal task. Since clients of an item may discover the cost of administration it excessively costly (at times ,much more than the cost of the item itself).Some item designers are irritating to make items that fix themselves. Self-healing [2]systems form a vicinity of research that is impulsively attractive and garnering increased attention, but not very well defined in terms of the compass. A self-healing is that it has the capability to determine, analyze and repair (or at least mitigate) disruptions have the services that it delivers.

D. KNN (K- Nearest Neighbour)

In design acknowledgment, the k-closest neighbours calculation (k-NN) is a non-parametric technique utilized

for grouping and regression.[6] In the two cases, the info comprises of the k nearest preparing cases in the component space. The yield relies upon whether k-NN is utilized for characterization or relapse:

In k-NN characterization, the yield is a class participation. A question is grouped by a lion's share vote of its neighbours, with the protest being doled out to the class most normal among its k closest neighbours (k is a positive whole number, commonly little). On the off chance that k = 1, at that point the question is basically appointed to the class of that solitary closest neighbour.

In k-NN relapse, the yield is the property estimation for the question. This esteem is the normal of the estimations of its k closest neighbours.

k-NN is a kind of occurrence based learning, or lethargic realizing, where the capacity is just approximated locally and all calculation is conceded until order. The k-NN calculation is among the least difficult of all machine learning calculations.

Both for order and relapse, a valuable method can be to allocate weight to the commitments of the neighbours, so that the closer neighbours contribute more to the normal than the more inaccessible ones. For instance, a typical weighting plan comprises in giving each neighbour a weight of 1/d, where d is the separation to the neighbour.[6]

The neighbours are taken from an arrangement of articles for which the class (for k-NN characterization) or the protest property estimation (for k-NN relapse) is known. This can be thought of as the preparation set for the calculation, however no unequivocal preparing step is required.

III. PROPOSED SYSTEM ARCHITECTURE

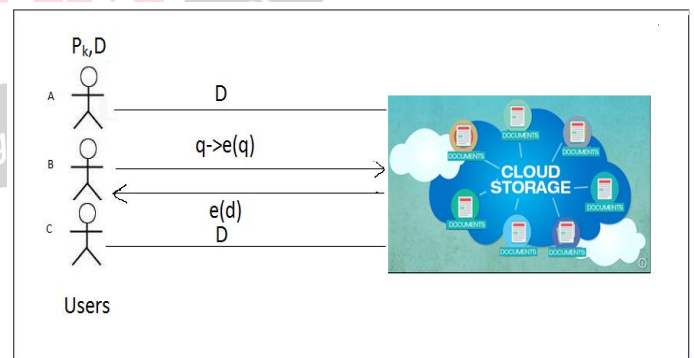


Fig 1: System Overview of HDM+

- D: Data
- Q: Query
- E(q): Encrypted data
- P_k: Patternkey

In general, every system follows a specific process in order to store the information, see fig 1 The HDM+ can be based on the information sources they use. There possible sources of information can be identified as input for the

recommendation process. The available sources are the user data (demographics), the item data (keywords, genres) and the user-item ratings (obtained by transaction data, explicit ratings).

In this paper we propose the HDM+ architecture where the pattern key and data is given to Cloud storage[9] there it is stored and by using the hashing method the keywords are encrypted, can store any type of data and the Encrypted data is $E(q)$. By using the hashing method the process becomes easy while searching the file using the key[6]. So that the processing speed becomes fast and execution becomes easy. AES algorithm is used to provide security for different applications.

A. Secure Retrieval of k -Nearest Neighbours (SRKNN)

In this stage, user initially sends his query q (in encrypted form) called trapdoor to C1. After this, C1 securely retrieve (in encrypted form) the class labels corresponding to the k -nearest neighbours of the input query q . After this, C1 securely retrieve (in encrypted form) the class labels corresponding to the k -nearest neighbours of the input query q . Secure Computation of Majority Class (SCMCK): Following from Stage 1, C1 compute the class label with a majority voting among the k -nearest neighbours of q . At the end of this step, only User knows the class label corresponding to his input query record q .

IV. MODULES DESCRIPTION

A. JOB COORDINATION SERVICE :

This is used to check all the services such as Storage and Resolver .There are two types

HDM+ STORAGE INTERFACE:

This Storage interface is used to store all the types of files except dll and exe files. This is used to store the Index terms called as keywords. To each and every file there will be some keywords based on the data stored in the file. These all files will be stored in the cloud. Let us consider an example.

Example 1: The file belongs to the Apple I phone and there may be two keywords here Apple and I phone .If we give keyword as apple then while searching it shows the apple fruit also. So here the keyword is I phone. Based on file requirement the keywords changes.

2) HDM+ DATA RESOLVER INTERFACE:

This is also called as the Data retrieval/Search item interface [7]; all the keywords are stored w.r.t to the files based on that the user can resolve the keywords .The files are displaced based on the keywords which the user gives.

B.STORAGE SERVICE

There are Data repository and Hash tree in the storage service

i)HDM+ Data repository

The data is stored in the form of the files called Data files and all the keywords are stored in the form Hash code called as Trapdoor. All the data in the file are converted into a Hash code [1]called as Data Signature .There will be the cloud storage all the files will be there .If the data stored is modified in the cloud store then the Self Healing error will occur. If user wants to open the particular modified file; then it displays the self healing error ;it shows the message as “You wants to recover” .If the user clicks on that then it automatically recover that the particular file. Here; the AES (Advanced Encryption Standard)[7] algorithm is used to recover that error which is Symmetric Algorithm.

ii)Hash-Tree :

The Hash-Tree consists of the File index center which has Trapdoor and Doc-No. Trapdoor is the Hash Code for the Keys and the Doc-No is the index number which consists of different one.

C.SEARCHABLE ENCRYPTION

Data Retrieval- Secure of k -Nearest Neighbour(SRKNN):In this summit fig-2, benefactor dominantly entrust enquire q (in encoded frame) called trapdoor.[2] After this, C1(cluster) consistently rescue (in scrambled frame) the set names subsequent to the k -closest neighbours of the info inquiry q . After this, C1 solidly recover (in encoded shape) the class marks resultant to the k -closest abutting of the information question q . Secure Computation of standard Class (SCMCK): Following from Stage 1, C1 process the class name with a typical assurance with the k -closest neighbours of q . At the end phases of this walk, just junkie knows the class name resultant to his info question record q .

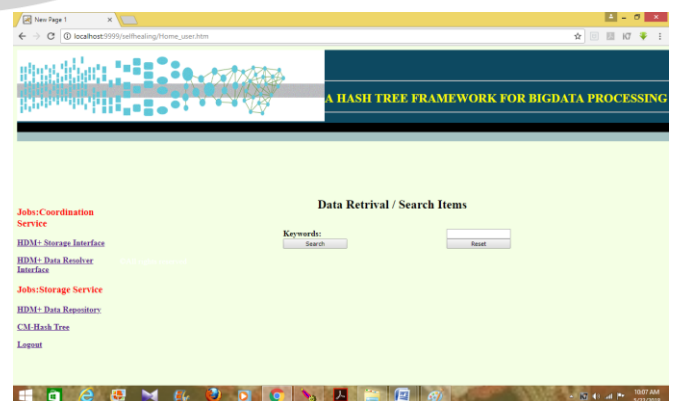


Fig 2: Data Retrieval from cloud storage

D. Multilevel interrelationship-HDM+ Storage interface

The engender encrypted keyword-mapping, multilevel affinity as shown in fig-3. The alliance among elements at an variety of levels keyword is mentioned in some entity metaphors at the component level. Entities at the facet level are allied with a set-level facet via form. A set-level element is contained in a source.

There is an edge flanked by two keywords if two rudiments at the constituent level mentioning these keywords are allied via a path. A ranking scheme is proposed which deals with relevance at many levels.

E. Self healing-AES cryptosystem

AES , the U.S. administration proclaimed that AES [7]might be worn to defend confidential data. It quickly turned into the avoidance encryption calculation for shielding secret information and also the essential widely accessible and open figure acknowledged by the NSA for zenith private information. The NSA picked AES as single of the cryptographic calculations to be worn by its information declaration directorate to safeguard national guard frameworks.[2] Its thriving use by the U.S. organization prompted broad use in the private zone, preminent AES to be changed into the most appreciated calculation worn in symmetric key cryptography. The translucent combination movement helped create an abnormal state of lightness in AES among guard and cryptography specialists. AES is auxiliary protected than its forerunners DES and 3DES as the calculation is more grounded and use longer key lengths additionally empowers more quickly encryption than DES and 3DES, influencing it to ideal for programming applications.

F . Hash Tree

In proposed system all the keywords or the keys are converted into the hash code by utilizing the Hash code method. Once this is done then all the keys are arranged in sorted order after this the searching begins by using the binary search. The searching of keywords becomes easy by using a hash key.

G .Binary Search

It starts by contrasting the objective esteem and the center component. The situation of the cluster is returned if the esteem coordinates the center component. In the event that the esteem is more noteworthy or not as much as the center esteem, the looking procedure proceeds in the upper or lower half of the cluster, disposing of the other half from exhibit. The process continues till the keyword or the trapdoor given by the user is found. Once the key is found then it displays the related file.



Fig 3: Data items are stored using trapdoors

Kw: Keywords

Td: Trapdoors

Trapdoors are the hash code for the particular keywords

V.COMPARATIVE ANALYSIS

EXISTING SYSTEM	PROPOSED SYSTEM
1.Supports Homogeneous data type	Supports Heterogeneous data type also
2. No self-healing	Self-healing is provided
3.Doesn't supports Fault tolerance	Supports Fault tolerance

V. CONCLUSION AND FURURE SCOPE

In this the Heterogeneous Data sets are stored and retrieved using Hash keys .The Self healing (using AES)is used to avoid the errors and produce security form the unauthorized users. The processing speed decreases so that time saves .To fulfill encrypted search, IDCrypt builds search indexes at proxies with the identifiers of encrypted data. We also design the token-adjustment search scheme to search across different indexes. To share encrypted data between different users, we propose the two-layer encryption (trapdoor---homorphic and symmetric (only one key---AES)---document encryption) scheme TLES to broadcast clandestine keys between dissimilar proxies.

REFERENCES

[1]CraigChambers,AshishRaniwala,FrancesPerry,StephenAdams, Robert R. Henry, Robert Bradshaw, and Nathan Weizenbaum. Flume Java: easy, efficient data-parallel pipelines. In PLDI, 2010.

[2] Jeffrey Dean and Sanjay Ghemawat. Map Reduce: simplified data processing on large clusters. Commun.ACM, 51(1), 2008.

[3] Yin Haul, Ashutosh Chauhan, Alan Gates, Günther Hagleitner, Eric N. Hanson, Owen O'Malley, Jitendra Pandey, Yuan Yuan, Rubio Lee, and Xuedong Zhang. Major technical advancements in Apache Hive. In *SIGMOD*, pages 1235–1246, 2014.

[4] Mohammad Islam, Angelo K. Huang, Mohamed Battista, Michelle Chiang, Santhosh Srinivasan, Craig Peters, Andreas Neumann, and Alejandro Abdelnur. Oozie: towards a scalable workflow management system for hadoop. In *SIGMOD Workshops*, 2012

[5] B. Saha et al. Apache Tez: A Unifying Framework for Modelling and Building Data Processing Applications. In *SIGMOD*, 2015.

[6] Altman, N.S.(1992). "An introduction to kernel and nearest-neighbour nonparametric regression" ,The American Statistician 46(3).

[7] Announcing the ADVANCED ENCRYPTION STANDARD (AES) November 26, 2001. from the original on March 12, 2017. Retrieved October 2, 2012.

[8] M. Zaharias et al. Spark: Cluster Computing with Working Sets. In *Hot Cloud*, 2010.

[9] M. Zaharias et al. Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. In *NSDI*, 2012.

