

# Implementation of Map-Reduce For Huge Data System

Manju Sharma<sup>1</sup>

<sup>1</sup>College of Computer Science and Information Technology, Jazan University, Jazan, Saudi Arabia

---

**Abstract:** - In the fields of engineering, technology, commerce, industry economics there is need to handle massive amounts of data. The capacity to analyze massive amounts of data derived from a variety of data sets continues to be a significant obstacle. Coping with a significant volume of data presents challenges for a wide variety of companies. They will be incapable to manage, alter, process, distribute, or retrieve vast volumes of data using standard software applications. As a result of the conventional software tools' high cost and extensive length of time required for data acquisition. The concept known as "big processing" refers to the management of programs that deal with very big datasets. These types of programs spend the majority of their runtime moving data from their data store to the computational node in some kind of a computer system. The primary obstacles faced by such apps are the limits placed on the amount of data that can be stored and the processing power available. Developers need 100s of processing units and a huge quantity of storage media to process complicated functions with huge files. These sorts of application areas process datasets that range from multiple terabytes to in size, but also traditional data management methods including such producing images and centralized information processing are ineffective at solving the problems caused by these types of application.

**Keywords:-** Map reduction, big data processing, Data efficiency techniques etc.

---

Date of Submission: 15-01-2023

Date of Acceptance: 31-01-2023

---

Google has made the Map-reduce - based framework available to the public in order to enable the analysis of huge datasets. It is used for the purpose of doing distributed analytics on huge datasets all over a collection of computing nodes. In needed to finalize each component of the calculation in an acceptable length of time, this computation has to be divided among thousands of computers that are contained inside a cluster. This is necessary because the quantity of inputs is too vast. This distributed idea entails readily parallelizing computations including using re-execution also as primary strategy for high availability in order to make it more robust. Google has never made its MapReduce implementation available to the public. At long last, an implementation mechanism of MapReduce known as Hadoop was made accessible to the public by the Apache Business. Developers are given the ability to carry out complicated calculations in a straightforward manner thanks to MapReduce, which hides the complexities of data distribution, multithreading, and high availability. Because it makes use of a networked share-nothing design, MapReduce is able to comprehend and analyze both regular and large datasets across a large number of nodes, which is one of the program's distinguishing characteristics. A distributed computer architecture known as share everything architecture is made up of several nodes. Every network is self-sufficient that has its personal disks, ram, and input/output (I/O) devices. Because each cluster in this form of design is soul and does not share anything with the other nodes across the networks, the architecture does not have any sticking points throughout the system.

## A MapReduce programming model drives from three fundamental phases:

1. The mapping step involves partitioning into different map functions, known as mappers, and running each mapping in parallel. Interim data item pairs are what are produced after the Map step is finished.
2. The Shake and Sort step involves hashing the input key of every Mapper in order to split the outputs of the Mappers. In this step, the quantity of partitioning is equivalent to the amount of reducers; inside the surmise, all data item pairs carry a same name, which is associated with the same split. Just after Mapping output has been partitioned, each split is then saved by a number in order to combine all data associated with that key.
3. Reduce phase: split into R Standard output (Reducer); moreover, each Reducer worked in conjunction and examines a unique set of intermediate keys.

Many other programming languages, such as LISP, Matlab, Languages such as java, Twitter, Ruby, and Rust, have each had their own version of the MapReduce library. A description of the Dremel workflow involves the following steps: a database is separated into numerous pieces of measurement; next, the Representation of an image processes each byte of output individually; and finally, the Reducing phase combines the results of all of the handled pieces of measurement. In the last step, known as the Feature extraction step, all of these

components are brought together to generate the final product. The Mapping function which takes in couples of input arguments, generates a series of intermediary name - value pairs, but then sends those pairs onto another Standard output so that almost all of its values connected to almost the same core may be combined into a single value. Its Reduce function takes an initial key and a collection of numbers for that keyword as input. It then combines these numbers in order to generate an output image that has a more manageable but still sufficient subset of data.

Processing as well as data analysis of this kind using typical relational database management system (DBMS) technology may be challenging for businesses dealing with significant volumes of cross data. Companies that operate on the Internet, such as Amazon, Yahoo, Facebook, among LinkedIn, are particularly susceptible to this kind of issue since they are tasked with the responsibility of processing vast volumes of data has quickly and cheaply as possible. In order to circumvent this problem, a significant number of enterprises of this kind have created their own nosql systems. For instance, Google is responsible for the creation of MapReduce and also the Google Shared Folder. Additionally, it developed a database management system referred to as BigTable. Only with assist of a machine learning that drive some of these powerful search amenities, it is now necessary to look reams and reams and return values in msec maybe less. The above algorithms began with Amazon's Mapreduce model, and they are responsible for both of these capabilities. The analysis of massive quantities of data with in modern world seems to be a very difficult challenge to solve. Working with large data is a significant undertaking, and as a result, doing analytics using big datasets is a challenging task. There seems to be a substantial bit of investment in sophisticated analytics methodologies such as Mapper, Hadoop, and Prism, as well as MapReduce enhancements to established relational database management systems. Technologies enabling analyzing large amounts of data are advancing at a fast pace.

To successfully manage such large amounts of data, hadoopMapReduce architecture has recently come to the forefront of widespread attention. Over the course of recent decades several years, the Dremel computing architecture has emerged as a preferred option for doing parallel, batch-based processing on massive amounts of data. Mapper became widely utilized once Google realized its potential and began effectively implementing it. In practice, that is a information processing tool that is both scalable nor fault-tolerant. It gives users the capacity to process massive amounts of data in concurrently using a large number of computational processors with a lower level of performance. MapReduce is quickly have become an industry standard because of its ease of use, scalability, and accident, and it really is getting substantial traction both from the business sector and the academic community. By dividing the computation into a number of smaller tasks that could be executed concurrently on different nodes within the network, we will be able to attain a high - level performance. This distributed information system, also known as a Dds, is used in the Jacquard architecture to first split data over numerous computers. These data is then represented using (key, value) pairings. The critical role is conducted individually master machine using the Software tool. On this device, we are able to preprocess that will be sent to the mapping functions and recording studios the data that will be returned by the reduction functions. The features of a program will determine the number of times that a maps and reduction method pair is run, which might be just once or any amount of instances. For the processing of huge datasets, Hdfs is an available as an expansive version of such Mapreduce. When it comes to managing physical storage all across group, it makes use of a migration flows filesystem. However, the system generates an undesirable speedup with far less important data sources, but it generates a decent pace with a greater data gathering that reinforces the count of virtual machines and reduces this same execution whilst also 30 percent in comparison to standard data mineral extraction or other postprocessing.

### Map Reduce Programming Model

As can be seen in Figure 1, the calculation is broken up into maps, scramble, and reduce stages when using the Mapreduce paradigm.

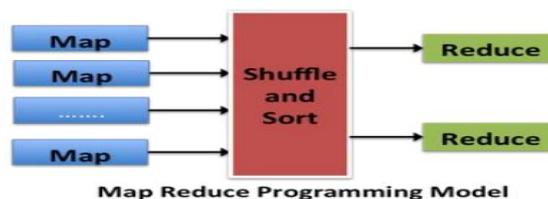


Figure 1: The MapReduce Programming Model

During the mapper, the output data is divided up into a wide range of input divides and then stored automatically among a computer that are part of this event. After the inputs have been dispersed over the clusters, the runtime will construct a substantial proportion of map jobs that will run concurrently in order to handle the inputs. The vital pairs that are provided as input are read by the map jobs, which then create one or

maybe more public keys that are intermediate. With in map job, the fundamental unit of information is referred to this as a tuple. As an instance, we will utilize the Word Frequency tool, which analyzes a collection of text files and determines the frequency with which each individual word appears in those files. A secret key may take the form of a single line in such a word doc when it comes to the processing for text documents. Every definition of such a public key is open for interpretation and modification by the operator. The users are responsible for deciding what should be done with every important combination, and this is all most crucial aspect of the process.

That MapReduce concept was first presented by Google. Mapper runtime is just a software implementation of both the MapReduce paradigm. Because it automatically manages job allocation, uptime, and other features of cloud services, it makes it considerably simpler for coders to build programs that use data parallelism. It also makes it possible for Facebook to take use of a great amount of inexpensive commodity machines to achieve great performance at such a quarter of the expense of a system that is created from a smaller number of much more costly high-end machines. Performance may be scaled using Dremel by executing parallel jobs on the units that are responsible for storing the task data. Each node is responsible for carrying out the duties independently and only has limited connection with another nodes.

Hadoop is indeed a MapReduce system that runs on free software. In order to make advantage of the HdfsMapreduce, one user must first compose a Mapper application by following the programming methodology that was outlined in the prior section. After that, when user sends the MapReduce task to a task queue, which really is a Program code that operates inside its own own specialized instance of the Java Virtual Machine (JVM). It is the job tracker's duty to ensure that the work run is properly coordinated. It then organizes the completion of those map-reduce tasks over a large cluster and divides the original job up into its component map-reduce jobs.

There seems to be a jobtracker process present on each computer, and this process is accountable for the schedule and execution of map/reduce activities on the computer system. In essence, the jobtracker is responsible for delegating map/reduce jobs to the ambiguous role, and the life's work are responsible for delegating responsibilities to the various cores. Additionally, the jobtracker process operates with its own someone else's standalone JVM. Fig 2 shows the device in its entirety.

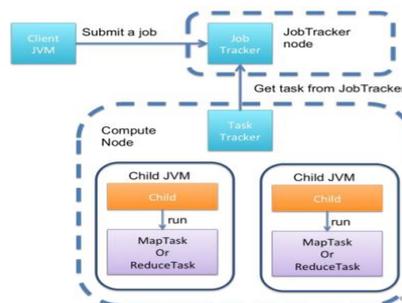


Figure 2: How Hadoop runs a MapReduce job

### Hadoop Distributed File System

“HDFS is the file system component of Hadoop (Refer Figure 5). While the interface to HDFS is patterned after the UNIX file system, faithfulness to standards was sacrificed in favor of improved performance for the applications at hand.”[22]

Hadoop is composed of Hadoop Map Reduce, an implementation of Map Reduce designed for large clusters, and the Hadoop Distributed File System(HDFS), a file system optimized for batch oriented workloads such as Map Reduce. Hadoop installation consists of a single master node and many Worker Nodes. The master is called as the Job Tracker, is responsible for accepting jobs from clients, dividing those jobs into tasks, and assigning those tasks to be executed by worker nodes. Each worker node runs a Task Tracker process that manages the execution of the tasks currently assigned to that node. [23]

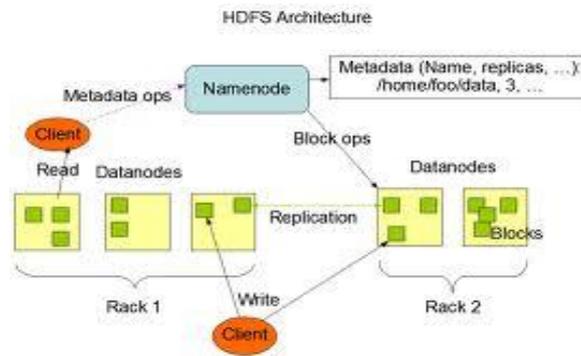


Fig 3 : HDFS Architecture Source: Apache Hadoop Website[22]

### What is Hadoop and why Hadoop?

Hadoop (a cloud computing framework) is a Java-based distributed system, is a new framework in the market. Since Hadoop is new and still being developed to add more features, there are many security issues which need to be addressed. Researchers have identified some of the issues and started working on this. Some of the notable outcomes[25]

Hadoop is a product framework that can be introduced on a ware Linux cluster to allow substantial scale distributed data investigation. Hadoop gives the strong Hadoop Distributed File System (HDFS) and in addition a Java-based API that permits parallel processing over the nodes of the cluster. Programs utilize a Map/Reduce execution engine which works as a fault-tolerant distributed figuring system over huge data sets - a strategy advanced by use at Google. There are separate Map and Reduce steps, each progression done in parallel, each working on sets of key-value sets. Processing can be parallelized more than a huge number of nodes chipping away at terabyte or bigger measured data sets. The Hadoop framework naturally plans map undertakings near the data on which they will work, with "close" which means a similar node or, at any rate, a similar rack. Node disappointments are likewise taken care of consequently.

Hadoop is an open source tool from the ASF – Apache Software Foundation. Its capacity is to storing data and running applications on clusters of product equipment.

For example, if a file is deleted in HDFS, the NameNode will immediately record this in the EditLog. It regularly receives a Heartbeat and a block report from all the DataNodes in the cluster to ensure that the DataNodes are live. It keeps a record of all the blocks in HDFS and in which nodes these blocks are stored. [24]

#### 1.1.5 Hadoop MapReduce for Multi-core

The amount of available slots for maps solutions to reduce activity slots in a tasktracker is always the same. Each slot is occupied by a JVM that is working on a certain job. One computing thread is all that is used by each and every JVM. So user has to select the number pf data slots in accordance with the total multiprocessor and the quantity of RAM that is available upon every node before they are allowed to have more least 1 core. This mapred-site.xml file is where the configuration settings may be changed.

Figure 5 provides several examples that demonstrate how to configure each computation node with two reduction slots as well as four mapped slots. It is possible to utilize the option to represent the variety of the computers that are part of the program. Each computer node may have a unique configuration for this parameter. The reason for this is because each individual computer in the group may have a varied quantity of storage as well as a unique number of processor cores.

```
<property>
  <name>mapred.tasktracker.map.tasks.maximum</name>
  <value>4</value>
</property>

<property>
  <name>mapred.tasktracker.reduce.tasks.maximum</name>
  <value>2</value>
</property>
```

Figure 4: Setting the map slot and reduce slot property in mapred-site.xml

When it initially starts up, your Map-reduce system uses the parameters that were supplied by the user. After the network has been booted up, the tasktracker will begin to collect the tasks that have been allotted to it. To do this, it will make regular calls towards the job finder heartbeat technique to evaluate whether or not there are any jobs that are ready to be carried out. When a job has been given to the jobtracker, it immediately begins working on completing that task. Its first thing that has to be done is to move the user's program code that has been compressed into such a File extension to the storage of the particular node and then unpack that into a locally current directory. This tasktracker will then construct an instance of a task racer in order to carry out the work. A new Jdk is started up by the jobtracker whenever the diagram task has to be carried out. This is done to avoid the possibility of the subscriber map tasks functions causing the tasktracker processes to break. As a direct consequence of this, each map-reduce job is carried out in its own individual JVM instance. Because of the cumulative impacts they have on Speed and ram consumption, the quantity of JVMs that are produced in a particular node (machine) may have a substantial influence on the system's overall performance. In order to maximize use of available Cpus in inter systems, there's also a trend toward spawning a high number of activities and, as either a result, JVMs. On a system with eight cores, for instance, that is not unheard of to have 24 map jobs running simultaneously. However, the effectiveness of this strategy is limited to applications that do not need a substantial quantity of storage. That's because it's difficult to generate a ton of jobs upon every virtual machine if each work requires a big amount of storage.

Because each Jdk is contained inside its own internal memory, it really is impossible for one Map job to share resources with another Map task executing on the very same node while that Mapper is operating within its own Jdk. As a consequence of this, the data structures that are utilized in map jobs are replicated across JVMs. This includes the in-memory learn data structures that are necessary for map/reduce workloads. Duplicating identical data model can significantly cut down on the cumulative amount of storage that is available per each map-reduce job in recollection applications. This is due to the fact that duplicated through learn data structures that are required by the proposal can start taking up a massive proportion of system memory, as displayed in Fig 6. With in case of a normal hash join program, for instance, it is necessary for every map job to keep a memory replica of the data structure. Duplicating your hash table will result in a reduction in the 4 gb Ram that may be allocated per each map operation

Conclusion:- Performance assessment analysis is an important component for enhancing the quality of any system. As for as education is concerned, quality is more essential than quantity. Very few assessments are based on quantity measures, which only predicts cognitive domains, whereas the affective domain is being neglected due to shortage of assessment mechanics. The feedback data collected from various academic institutions are heterogeneous, amorphous and semi-structured consist of alphabets, numeric, alpha numeric along with symbols. Processing or scrutinizing the information from this type of colossal data is complicated by using the traditional conventional relational data base techniques.

The automatic removal of unwanted terms from academic feedback data is a crucial step since unwanted terms may affect the knowledge gained from the data. The existing techniques applied for pre-processing the feedback data have a lot of drawbacks. Text categorization is the process of automatically selecting more predetermined category topics to a given text written in a natural language, according to its similarity with respect to a previously labelled corpus used as a reference set.

Complex information retrieval and categorization systems are needed to process queries, filter, and store and organize data which are mainly composed of texts. In this type of academic feedback data analysis, unstructured texts should be converted into organized and structured information. Classification of text is the key for processing, retrieval of queries and efficient management of information.

Many big data technologies and analytics tools are available to do individual task of processing, storing, extraction, visualize etc., However, the challenge lies in identifying and integrating real time problems and stack it for improved accuracy. This motivated the researcher to focus on developing a framework for analyzing the academic feedback by using Big Data Analytics.

### References/ Bibliography:-

- [1]. A.Moturi, C., & K. Maiyo, S. (2012). Use of Mapreduce for Data Mining and Data Optimization on a Web Portal. International Journal of Computer Applications, 56(7), 39–43. <https://doi.org/10.5120/8906-2945>
- [2]. Abdullah M Alghamdi, & Fahad Alghamdi. (2019). Enhancing Performance of Educational Data Using Big Data and Hadoop. International Journal of Applied Engineering Research, 14(19), 3814–3819.
- [3]. Abhishek, K., Ahad, P., Nikita, K., Rushikesh, H., & Student, B. E. (2018). Survey on MapReduce and Hadoop. 3(11), 422–424.
- [4]. Abouzeid, A., Bajda-Pawlikowski, K., Abadi, D., Silberschatz, A., & Rasin, A. (2009). HadoopDB: An architectural hybrid of mapreduce and DBMS technologies for analytical workloads. Proceedings of the VLDB Endowment, 2(1), 922–933. <https://doi.org/10.14778/1687627.1687731>
- [5]. Algorithm, T. (2014). MapReduce Inputs and Outputs ( Java Perspective ).
- [6]. Babu, S. (2010). Towards automatic optimization of MapReduce programs. Proceedings of the 1st ACM Symposium on Cloud Computing, SoCC '10, 137–142. <https://doi.org/10.1145/1807128.1807150>
- [7]. Babu, S., & Herodotou, H. (2012). Massively parallel databases and MapReduce systems. Foundations and Trends in Databases, 5(1), 1–104. <https://doi.org/10.1561/19000000036>

- [8]. Bhaskar, A., & Ranjan, R. (2019). Enhanced and efficient memory model for Hadoop MapReduce. *International Journal of Innovative Technology and Exploring Engineering*, 9(1), 161–168. <https://doi.org/10.35940/ijitee.A3958.119119>
- [9]. del Río, S., López, V., Benítez, J. M., & Herrera, F. (2015). A MapReduce Approach to Address Big Data Classification Problems Based on the Fusion of Linguistic Fuzzy Rules. *International Journal of Computational Intelligence Systems*, 8(3), 422–437. <https://doi.org/10.1080/18756891.2015.1017377>
- [10]. Dhamecha, M. V. (2020). Improved MapReduce algorithm with efficiency in structure and relations in the large body of the data (bigdata) (Issue 159997107017).
- [11]. Dhankhar, A., & Solanki, K. (2019). A comprehensive review of tools & techniques for big data analytics. *International Journal of Emerging Trends in Engineering Research*, 7(11), 556–562. <https://doi.org/10.30534/ijeter/2019/257112019>
- [12]. Dr.Siddaraju, Sowmya, C. L., Rashmi, K., & Rahul, M. (2014). Efficient Analysis of Big Data Using Map Reduce Framework. *International Journal of Recent Development in Engineering and Technology*, 2(6), 64–68.
- [13]. Foundation, T. A. S. (2008). MapReduce Tutorial. Source, 1–43.
- [14]. Ganguly, P. (2020). Big Data Analytics: Using Hadoop Inspired MapReduce and Apache Spark. 7(2), 72–82.
- [15]. Gao, T., Yang, M., Jiang, R., Li, Y., & Yao, Y. (2019). Research on Computing Efficiency of MapReduce in Big Data Environment. *ITM Web of Conferences*, 26(2019), 03002. <https://doi.org/10.1051/itmconf/20192603002>
- [16]. Gong, N. A. N. (2011). Using Map-Reduce for Large Scale Analysis of Graph-Based Data.
- [17]. González-Vélez, H., & Kontagora, M. (2011). Performance evaluation of MapReduce using full virtualisation on a departmental cloud. *International Journal of Applied Mathematics and Computer Science*, 21(2), 275–284. <https://doi.org/10.2478/v10006-011-0020-3>
- [18]. McCreddie, R., MacDonald, C., & Ounis, I. (2012). MapReduce indexing strategies: Studying scalability and efficiency. *Information Processing and Management*, 48(5), 873–888. <https://doi.org/10.1016/j.ipm.2010.12.003>
- [19]. Zamrodah, Y. (2016). No Title No Title No Title (Vol. 15, Issue 2).
- [20]. Zhang, J., Zhang, X., & Zhang, W. (2013). Microseismic search engine. Society of Exploration Geophysicists International Exposition and 83rd Annual Meeting, SEG 2013: Expanding Geophysical Frontiers, 2140–2144. <https://doi.org/10.1190/segam2013-1277.1>
- [21]. Zhou, F., Pham, H., Yue, J., Zou, H., & Yu, W. (2015). SFMapReduce: An optimized MapReduce framework for Small Files. *Proceedings of the 2015 IEEE International Conference on Networking, Architecture and Storage, NAS 2015*, 23–32. <https://doi.org/10.1109/NAS.2015.7255218>
- [22]. The MapReduce Framework : Analysis and The Research Perspectives Madhavi Vaidya Dr.Shrinivas Deshpande Dr. Vilas Thakare Mumbai(Author) Amravati [vamadhavi@gmail.com](mailto:vamadhavi@gmail.com)
- [23]. Tyson Condie, Neil Conway, Peter Alvaro, Joseph M.Hellerstein, Khaled Elmeleegy, Russel Sears, Map Reduce Online, pages 1-14 , Oct 9,2009
- [24]. REVIEW PAPER ON HADOOP AND MAP REDUCE Nishant Rajput1 , Nikhil Ganage2 , JeetBhavesh Thakur3, IJRET: International Journal of Research in Engineering and Technology eISSN: 2319-1163 | pISSN: 2321-7308
- [25]. Big Data Analytics Challenges and Solutions in Cloud, Ms. Manju Sharma, American Journal of Engineering Research (AJER) e-ISSN: 2320-0847 p-ISSN : 2320-0936 Volume-6, Issue-4, pp-46-51

Manju Sharma. "Implementation of Map-Reduce For Huge Data System." *IOSR Journal of Computer Engineering (IOSR-JCE)*, 25(1), 2023, pp. 50-55.