

## A SURVEY ON SCHEDULING IN HADOOP FOR BIGDATA PROCESSING

Bhavsar Nikhil, Bhavsar Riddhikesh, Patil Balu, Tad Mukesh

Department of Computer Engineering

JSPM's Imperial College of Engineering and Research,  
Wagholi, Pune, Maharashtra, India

nsbhavsar007@gmail.com, riddh.bhavsar@gmail.com, patilbalu7755@gmail.com

**Abstract:** Scheduler is a way of assigning jobs to various available resources. Job assignment is done in such a way that it should minimize starvation of jobs and maximum utilization of resources. Scheduler efficiency is dependent on realistic workloads and clusters. Here, we are introducing a scheduler technique for real, multi node, complex system as a Hadoop. Based on size, priority is assigned to make scheduling efficient. This makes our algorithm different than conventional scheduling algorithm. Performance of the proposed scheduling technique can be increased by assigning Deadline constraints for local optimality of data. Map Reduce framework is used for execution of jobs.

**Keywords:** Deadline constraints, Hadoop, Local optimality, Map reduce, Scheduler.

### 1. INTRODUCTION

Big Data describes innovative techniques and technologies to capture, store, distribute, manage and analyze petabyte or larger-sized datasets with high-velocity and different structures. Another definition of big data is the collection of clusters or realistic workload that gives the benchmark to the large amount of data. To handle these large dataset such as big data a technology is introduced called Hadoop. Hadoop is open source framework uses Hadoop distributed file system (Hdfs) for storage. In the advent of large scale data the scheduler circums with the task handling with numerous data processes, we proceed with the implementation of new scheduling method which overcomes the problem of handling the large amount of data. Our solution implements a size-based, preemptive scheduling discipline. The scheduler allocates cluster resources such that job size information? Which is not available a-priori? Is inferred while the job makes progress toward its completion. This ensures that neither small nor large jobs suffer from starvation. The outcome of our work

materializes as a full-fledged scheduler implementation that integrates seamlessly in Hadoop. The scheduler works with jobs to whom the priorities must be given first as the large amount of data may contain the files or job with same name, same size, so to decide the job, to calculate the size of the job for processing the jobs we are implementing the scheduler. To make our scheduler more efficient deadline constraints are added with sized based priority.

## 2. BIG DATA

Big data is an unrolling concept which defines the typical large amount of unstructured, semi-structured and structured data. Big data doesn't specify the derived data when speaking about the Zeta bytes, petabytes and Exabyte's size of data.

Nowadays Big data is expanding around us every single minute. Its increasing factor is the use of internet processing and social media generates it. Big data comes from multi sources at an alarming volume, variety and velocity [1].

Big data plays important role in academics as well as industrial sources where the data generation ratio is at the peak. Relational database cannot be use always because it has certain limits so data scientist developed the concept of big data which can be used in efficient manner to accept challenges of bulk data storage.

Mike Guiltier, Forrester Analyst, proposes a definition that attempts to be pragmatic and actionable for IT professionals: "Big Data is the frontier of a firm's ability to store, process, and access (SPA) all the data it needs to operate effectively, make decisions, reduce risks, and serve customers" (Guiltier, December 2012).

### 2.1 3V'S of Big Data

- **Velocity**:-Analysis of the Batch process, Real Time process, data streaming, when one takes chunks of data, it submits the job then there may be certain delay from the server.[1][2]
- **Volume**:-The size of data files is increasing with high ratio. The data size .The storage area is also increased from megabytes to Giga bytes, Zeta Bytes, petabytes and exabytes. Eg. Facebook generates around 500Tb of Data every single day.
- **Variety**:-In the era of internet the variety of data is changing day by day. It follows with structured data like tables and unstructured like videos, xml, audio etc.

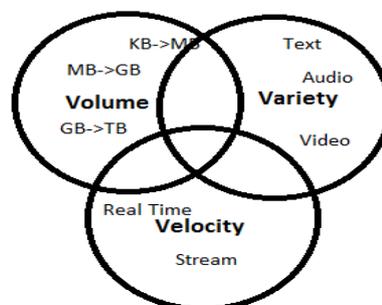


Fig.1: 3 V's of Big Data

## 2.2 Differentiate Between Big Data and DBMS

- Parallel database are actively querying the large data sets.
- DBMS supports only structured data where as Big Data supports structured, semi-structured and unstructured data.
- Parallel database are relational paradigm of rows and columns.

## 2.3 Pillars of Big Data

- Text: The nature of text with all types of structures, unstructured and semi-structured.
- Tables: Table form-rows and columns.
- Graphs: Degree of separation, subject predicates, object prediction, semantic discovery [6].

## 3. HADOOP

### 3.1 Introduction of Hadoop

Hadoop is a framework which is open-source for processing and storing big data in a distributed way on large clusters. It processes two tasks: large storage of data and processes faster with data.

- i. Framework: It is to develop and to run software applications provides– programs, tool sets, connections, etc [7].
- ii. Open-source software: It is downloaded and can be used to develop the commercial software. Its broad [7] and open network can be easily, managed by the developers.
- iii. Distributed Form: Data is stored on multiple computers and divided into equal chunks on the parallel connected nodes.
- iv. Large storage: The Hadoop framework stores large amount of data by dividing into separate blocks and stored on clusters.

Hadoop includes:

1. Hadoop Distributed File System (HDFS): It manages the large amount of data that is stored in distributed fashion.
2. Map Reduce: a framework and a programming model for distributed processing [6] [7] on parallel clusters.

### 3.2 HDFS File Structure

The HDFS is a distributed file system developed to execute on commodity hardware. The feature of HDFS is fault-tolerance and to be deployed on low-cost hardware. HDFS provides great access to data and is suitable for applications that have large data sets. [5].

The distribution of files is divided into chunks of 64MB size.

The HDFS architecture contains a unit Name Node, multiple Data Nodes; it can also be stated as master slave like architecture. Name Node manages the mapping of file system namespace and regulates access to the block of files. Data Nodes are responsible for the process of read and write from the clients. It also performs operations like creation of blocks and deletion of blocks. [Hadoop. Apache]

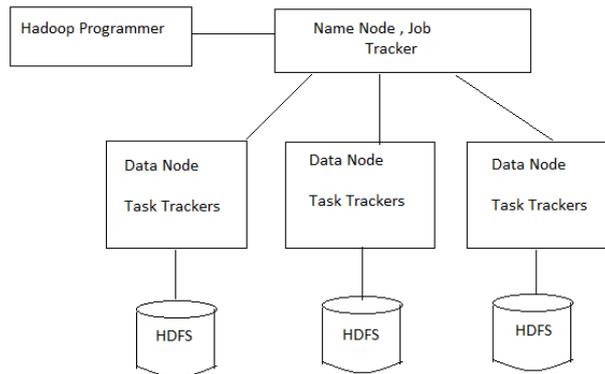


Fig.2: HDFS Architecture

### 3.3 Map reduce

Map Reduce is a programming model and software framework first developed by Google (Google’s Map Reduce paper submitted in 2004). Intended to facilitate and simplify the processing of vast amounts of data in parallel on large clusters of commodity hardware in a reliable, fault-tolerant manner, Petabytes of data, Thousands of nodes. Computational processing occurs on both: Unstructured data file system Structured database. Underlying runtime system automatically parallelizes the computation across large-scale clusters of machines [4]

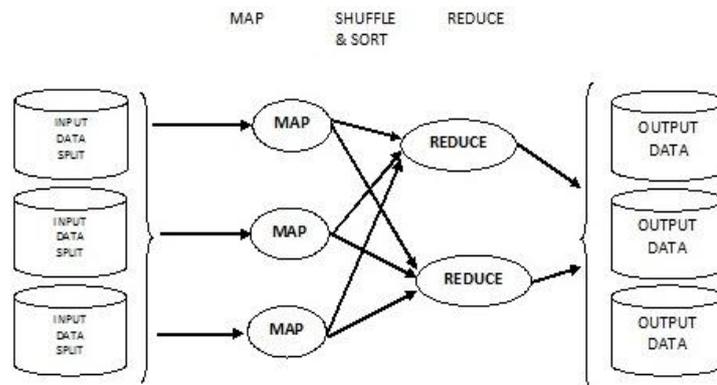


Fig.3: Map Reduce Architecture

## 4. HADOOP SCHEDULING

### 4.1 Scheduling

Hadoop implements the schedulers where resources are assigning to the jobs. Traditional scheduling gives us the scenario that all algorithms are not same and might not be as effective and dependent completely. We have to study various schedulers available in Hadoop and how relevant they are as compared to traditional schedulers in terms of performance, throughputs, time consuming etc. First in First out (FIFO) scheduler is a default scheduler which considers e order for submission of the jobs to get executed.

- i. **FIFO Scheduler:** FIFO is traditional default scheduler which performs with the use of queue. Job is divided into several tasks and then loads to free slots of the queue on the task tracker. In FIFO the job have to wait for execution due to acquisition of

clusters take place this leads to wait for other jobs for their turn. Shared clusters have ability for offering resources to users.

- ii. **Fair Scheduler:** Fair scheduler groups jobs into “pools” and it allots each pool a guaranteed minimum share with segment more capacity equally between pools. Facebook first develop the concept of fair scheduler to manage the access to their Hadoop and get the subsequent relation with the Hadoop environment. Pools have minimum mapping slots and minimum reduced slots. The Fair Scheduler supports. Preemption, so if a pool has not received its fair share for a certain period of time, then the scheduler will kill tasks in pools running over capacity in order to give the slots to the pool running under capacity [8].
- iii. **Capacity Scheduler:** Capacity scheduler organizes jobs into queues. It shares as percent of cluster. FIFO scheduling within each queue and it supports preemption. Yahoo developed the capacity scheduler addresses a usage scenario where the number of users is large, and there is a need to ensure a fair allocation of computation resources amongst users [3][6]. When a Task Tracker slot becomes free, the queue with the lowest load is chosen, from which the oldest remaining job is chosen. A task is then scheduled from that job.

## 5. CONCLUSION

Big Data is in demand now days in the market. Prior to this there was traditional database system but now the large amount of data is been generated and this data may be structured and, unstructured data in the industries. To overcome the issue of Big Data storage and processing the open source framework named Hadoop is developed by Apache can be used. Hadoop gives a source to Big Data processing with its components like Hadoop Distributed File System (HDFS) and Map Reduce. To process with the Big Data the default scheduler called FIFO has been used and research has been made to overcome with the problems related to the FIFO .In this paper we have discussed many techniques for making the efficient scheduler for the map reduce so that we can speed up our system or data retrieval technique like quinsy, Asynchronous Processing, Speculative Execution, Job Awareness, Delay Scheduling, Copy Compute Splitting etc. had made the scheduler effective for the faster processing. So that schedulers must work effectively with rapid processing and the faster execution of job with the big data.

## REFERENCES

- [1]. Big Data definition "<http://searchcloudcomputing.techtarget.com/definition/big-data-Big-Data>"
- [2]. Bigdata 3V's "<http://www.datasciencecentral.com/forum/topics/the-3vs-that-define-big-data>"
- [3]. B.Thirumala Rao, Dr. L. S.S. Reddy,“Survey on Improved Scheduling in Hadoop Map Reduce in Cloud Environments”.
- [4]. Dean, J. and Ghemawat, S. 2008. Map Reduce: simplified data processing on large clusters
- [5]. Hadoop Distributed File System [[http://Hadoop.apache.org/docs/r1.2.1/hdfs\\_design.html](http://Hadoop.apache.org/docs/r1.2.1/hdfs_design.html)].
- [6]. Harshawardhan S. Bhosale<sup>1</sup> , Devendra P. Gadekar<sup>2</sup> Big Data Processing Using Hadoop: Survey on Scheduling
- [7]. [http://www.sas.com/en\\_us/insights/big-data/Hadoop.html](http://www.sas.com/en_us/insights/big-data/Hadoop.html)
- [8]. Scheduling with Fair Scheduler".<http://arxiv.org/ftp/arxiv/papers/1207/1207.0780.pdf>"