

Job schedulers for Big data processing in Hadoop environment: testing real-life schedulers using benchmark programs



Mohd Usama, Mengchen Liu, Min Chen*

Embedded and Pervasive Computing (EPIC) Lab, School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan 430074, China

ARTICLE INFO

Keywords:

Big Data
Hadoop
MapReduce
HDFS
Scheduler
Classification
Locality
Benchmark

ABSTRACT

At present, big data is very popular, because it has proved to be much successful in many fields such as social media, E-commerce transactions, etc. Big data describes the tools and technologies needed to capture, manage, store, distribute, and analyze petabyte or larger-sized datasets having different structures with high speed. Big data can be structured, unstructured, or semi structured. Hadoop is an open source framework that is used to process large amounts of data in an inexpensive and efficient way, and job scheduling is a key factor for achieving high performance in big data processing. This paper gives an overview of big data and highlights the problems and challenges in big data. It then highlights Hadoop Distributed File System (HDFS), Hadoop MapReduce, and various parameters that affect the performance of job scheduling algorithms in big data such as Job Tracker, Task Tracker, Name Node, Data Node, etc. The primary purpose of this paper is to present a comparative study of job scheduling algorithms along with their experimental results in Hadoop environment. In addition, this paper describes the advantages, disadvantages, features, and drawbacks of various Hadoop job schedulers such as FIFO, Fair, capacity, Deadline Constraints, Delay, LATE, Resource Aware, etc, and provides a comparative study among these schedulers.

1. Introduction

1.1. Motivation

A significant amount of research has been done in the field of Hadoop Job scheduling; however, there is still a need for research to overcome some of the challenges regarding scheduling of jobs in Hadoop clusters. Industries estimate that 20% of the data is in structure form while the remaining 80% of data is in semi structure form. This is a big challenge not only for volume and variety but also for data processing, which can lead to problems for IO processing and job scheduling. Fig. 1 shows the architecture of a Hadoop distributed system.

As we know, this is an era of Big Data where humans process a significant amount of data, in the range of terabytes or petabytes, using various applications in fields such as science, business, and commerce. Such applications require a considerable amount of input/output processing and spend most of the time in IO processing, which is a major part of job scheduling. It is reported that at the Facebook and Microsoft Bing data center, IO processing requires 79% of the jobs' duration and 69% of

the resources. Therefore, here we present a comprehensive study of all Hadoop schedulers, in order to provide implementers an idea on which scheduler is the best fit for which job, so that the execution of a job does not take much time for IO processing and job scheduling becomes much easier. This paper will be useful for both beginners and researchers in understanding Hadoop job scheduling in Big data processing. It will also be useful in developing new ideas for innovations related to Hadoop scheduler.

1.2. The definition: big data

Many definitions of big data have been presented by scientists. Among all the definitions, the most popular definition [97] was presented by Doug Laney (2001) in his META Group research note, which describes the characteristics of datasets that cannot be handled by the traditional data management tools described below.

Three V's: volume (size of datasets and storage), velocity (speed of incoming data), and variety (data types). With the development of discussion and increasing research interest in big data, the Three V's have

* Corresponding author.

E-mail addresses: usama6832@gmail.com (M. Usama), mengchenliu.cs@qq.com (M. Liu), minchen2012@hust.edu.cn (M. Chen).

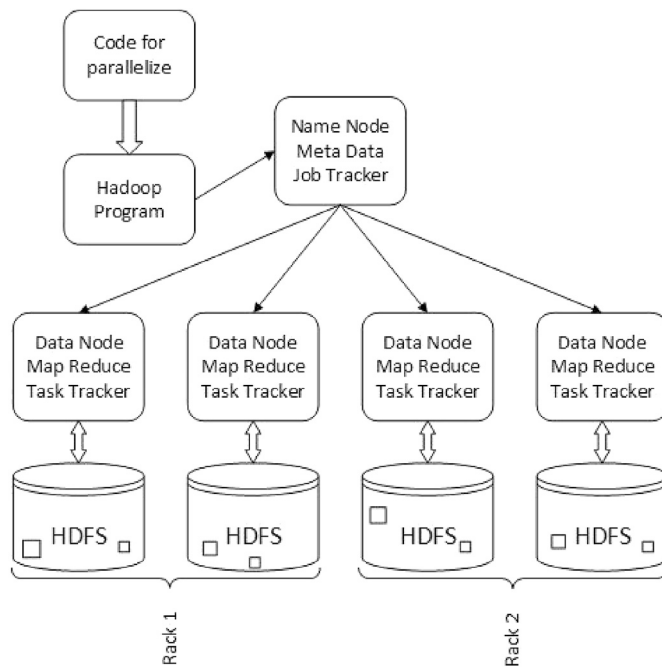


Fig. 1. Hadoop distributed system architecture.

been expanded to the Five V's: volume, velocity, variety, veracity (integrity of data), value (usefulness of data), and complexity (degree of interconnection among data structures).

1.3. Characteristics of big data

Vs of Big data [70]. Velocity of data: velocity means the speed of data processing. The increasing rate at which data flows into an organization has followed a similar pattern to that of volume. Volume of data: Volume refers to the size of data and storage capacity. A large amount of data comes into the storage from many sources such as social media web, business transactions etc. Variety of data: Variety refers to the different types of data. Big data includes a variety of data, both structured and unstructured.

1.4. What is big data?

Big data [24,32] has come into existence since the last few years. Big data is similar to a conventional database system, but the difference is that it exceeds in size and processing capacity [32]. Big data is very big, moves very fast and includes both structured and unstructured formats of data, whereas conventional databases include only structured format of data [10]. To obtain the benefits of big data, an alternative way must be chosen to process the data [65]. Big data is data that exceeds the processing capacity of conventional database systems [77].

Today's big data [82] plays a vital role in various disciplines and has become popular in computer science and technology, business and finance, banking and online purchasing, oceanography, astronomy, health-care and so on [34]. At present, big data has become a necessity in business analytics and many other fields [2,50]. It provides tremendous benefits for business enterprises [59,60]. Big data is made up of a large number of datasets. The size of these datasets continues to increase day by day as data comes in continuously from different sources such as social media sites, business transactions, personal data, digital photos, etc. Big data has massive amount of unwanted data in both structured and unstructured formats. In structured data, the data is stored in a systematic and well defined manner, while in unstructured data, the data is stored in an unsystematic and undefined manner.

The data coming from Wikipedia, Google, and Facebook have an unstructured format, whereas the data coming from E-commerce transactions have a structured format [3]. Due to the presence of both structured and unstructured data, a number of challenges arise in big data such as data capture, sharing, privacy, data transfer, analysis, storage, search, job scheduling, handling of data, visualization, and fault tolerance [2,57]. It is very complicated and difficult to handle these challenges using traditional database management tools [2].

Traditional data management tools [74] are not capable of processing, analyzing, and scheduling jobs in big data [64]. Therefore, we use a different set of tools to handle these challenges. Hadoop is the most suitable tool to handle all the challenges in big data. Hadoop is an open source software framework for processing big data, and was founded by Apache. It can process large amount of data, in the range of petabytes. Hadoop is a highly reliable cloud computing platform [89], and ensures a high availability of data by making copies of data at different nodes. It is scalable and takes care of the detection and handling of bugs.

There are two components of Hadoop, HDFS and MapReduce. The Hadoop distributed file system is used for data storage, while MapReduce is used for data processing. MapReduce has two functions, Map and Reduce. The functions are both written by the user, and the functions take values as input key value pairs and output the result as a set of key value pairs. First, the Map produces intermediate key value pairs, then the MapReduce function combines all the intermediate values having the same intermediate key as a library, and finally it is passed to the Reduce function. The Reduce function receives the intermediate key with a set of values for that key and merges them to make a smaller set of values.

The aim of scheduling of jobs [84] is to enable faster processing of jobs and to reduce the response time as much as possible by using better techniques for scheduling depending on the jobs, along with the best utilization of resources. FIFO scheduling is default scheduling mode in Hadoop; the jobs coming first get higher priority than those coming later. In some situations, this type of scheduling has a disadvantage, that is, when longer jobs are scheduled prior to shorter jobs, it leads to starvation. Fair scheduling shares the resources equally among all jobs. Capacity scheduling was introduced by Yahoo. It maximizes the utilization of resources and throughput in clusters. LATE scheduling policy was developed to optimize the performance of jobs and to minimize the job response time by detecting slow running processes in a cluster and launching equivalence processes as the background. Facebook uses Delay scheduling, to achieve better performance and lower response time for map tasks by applying changes to MapReduce. In deadline scheduler, the deadline constraints are specified by the user before scheduling the jobs in order to increase system utilization. Resource aware scheduling improves resource utilization; it uses node, master node, and worked node to complete job scheduling. In matchmaking scheduling, each node is marked by the locality marker, which ensures that every node gets an equitable chance to seize a local task. Through this scheduling, high data locality and better cluster utilization is achieved.

There have already been a few review papers on job scheduling algorithms for Big data processing. Yoo, D. et al. (2011) presented a comparative study on job scheduling methods and discussed their strengths and weakness [63]. Rao, B. T. et al. (2012) presented a review on scheduling algorithms and provided guidelines for the improvement of scheduling algorithms in Hadoop MapReduce [77]. Sreedhar C et al. (2015) presented a survey on big data management and discussed various scheduling algorithms in Hadoop. They also discussed the latest advancements related to scheduling algorithms [72]. Jyoti V Gautam et al. (2015) presented a paper on the scheduling policies for Hadoop and performed a comparative study on MapReduce optimization techniques [12].

The remaining parts of the paper are organized as follows. Section II describes challenges for job scheduling in Big Data. Section III describes the architecture, working, features, and requirements of job scheduling in Hadoop. Section IV describes the various Hadoop job schedulers along with their advantages and disadvantages, and compares the various

Download English Version:

<https://daneshyari.com/en/article/7111772>

Download Persian Version:

<https://daneshyari.com/article/7111772>

[Daneshyari.com](https://daneshyari.com)