

Chapter 7

Hadoop MapReduce Programming

ABSTRACT

The second major component of Hadoop is MapReduce. It is the software framework for Hadoop environment. It consists of a single resource manager, one node manager per node, and one application manager per application. These managers are responsible for allocating necessary resources and executing the jobs submitted by clients. The entire process of executing a job is narrated in this chapter. The architecture of MapReduce framework is explained. The execution is implemented through two major operations: map and reduce. The map and reduce operations are demonstrated with an example. The syntax of different user interfaces available is shown. The coding to be done for MapReduce programming is shown using Java. The entire cycle of job execution is shown. After reading this chapter, the reader will be able to write MapReduce programs and execute them. At the end of the chapter, some research issues in the MapReduce programming is outlined.

INTRODUCTION

The major components of Hadoop systems are Hadoop Distributed File System (HDFS) and MapReduce. The MapReduce is a software framework for easily writing applications. It is mainly used for processing larger amounts of data in-parallel on large clusters (thousands of nodes) in a reliable and fault-tolerant manner. This chapter gives an overview of MapReduce programming.

DOI: 10.4018/978-1-5225-3790-8.ch007

Also it explains clearly the different APIs available for programming. The programming can be done in different languages like C, C++, C#, Java, Perl, PHP, Python and Ruby. But this chapter focuses on programming with Java only.

BACKGROUND

The growth of data in the recent applications in the Internet is highly alarming. Analyzing such kind of large data (data analytics) is the demand in the business process requirements. Though many algorithms and techniques have been developed to mine such kind of large data and have been invested in the analytics, the turnaround time is not satisfactory. The huge storage requirements and computing requirements have dictated the distributed computing environment. The Hadoop based Distributed File System has enabled this. The principle involves dividing the jobs into small independent pieces (in many cases split manually) and mapping to various computing system and combining the solution back in a synchronized manner.

WORKING OF MAPREDUCE

The MapReduce framework consists of a single master ResourceManager, one slave NodeManager per cluster-node, and AppMaster per application. The user/application can submit the work to be executed as a job. The input and output of the job are stored in file system. The framework takes care of splitting the job into number of smaller tasks, scheduling the tasks across different nodes and monitoring them. If the task fails, the framework re-executes the job automatically without user intervention. The tasks are normally scheduled in the nodes where data is already present and hence the network bandwidth is properly utilized.

The applications can specify the input/output locations and other job parameters in “job configuration”. Then the client submits the jar/executable file of the job along with its configuration to the ResourceManager. The ResourceManager then:

- Distributes software/configuration to the slaves
- Schedules the tasks

20 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/hadoop-mapreduce-programming/216602

Related Content

Efficient String Matching Algorithm for Searching Large DNA and Binary Texts

Abdulraakeeb M. Al-Ssulami, Hassan I. Mathkourand Mohammed Amer Arafah (2020). *Data Analytics in Medicine: Concepts, Methodologies, Tools, and Applications* (pp. 298-324).

www.irma-international.org/chapter/efficient-string-matching-algorithm-for-searching-large-dna-and-binary-texts/243117

Using Call Detail Records of Mobile Network Operators for Transportation Studies

Erki Saluveerand Rein Ahas (2014). *Mobile Technologies for Activity-Travel Data Collection and Analysis* (pp. 224-238).

www.irma-international.org/chapter/using-call-detail-records-of-mobile-network-operators-for-transportation-studies/113213

Big Data and Cloud Interoperability

Ahmad Yusairi Bani Hashim (2016). *Managing Big Data Integration in the Public Sector* (pp. 59-69).

www.irma-international.org/chapter/big-data-and-cloud-interoperability/141103

Focused Error Analysis: Examples from the Use of the SHEEP Model

Deborah J. Rosenorn-Lanngand Vaughan A. Michell (2016). *International Journal of Big Data and Analytics in Healthcare* (pp. 30-48).

www.irma-international.org/article/focused-error-analysis/171403

Predictive Optimized Model on Money Markets Instruments With Capital Market and Bank Rates Ratio

Bilal Hungundand Shilpa Rastogi (2023). *International Journal of Data Analytics* (pp. 1-20).

www.irma-international.org/article/predictive-optimized-model-on-money-markets-instruments-with-capital-market-and-bank-rates-ratio/319024