Chapter 3 Hadoop History and Architecture

ABSTRACT

As the name indicates, this chapter explains the evolution of Hadoop. Doug Cutting started a text search library called Lucene. After joining Apache Software Foundation, he modified it into a web crawler called Apache Nutch. Then Google File System was taken as reference and modified as Nutch Distributed File System. Then Google's MapReduce features were also integrated and Hadoop was framed. The whole path from Lucene to Apache Hadoop is illustrated in this chapter. Also, the different versions of Hadoop are explained. The procedure to download the software is explained. The mechanism to verify the downloaded software is shown. Then the architecture of Hadoop is detailed. The Hadoop cluster is a set of commodity machines grouped together. The arrangement of Hadoop machines in different racks is shown. After reading this chapter, the reader will understand how Hadoop has evolved and its entire architecture.

INTRODUCTION

Hadoop is an open source framework used for storing and processing big data. It is developed by Apache Software Foundation. Hadoop environment can be setup with commodity hardware alone. It is a framework that supports distributed environment with cluster of commodity machines. It can work with single server or can scale up including thousands of commodity machines.

DOI: 10.4018/978-1-5225-3790-8.ch003

Hadoop has undergone number of revisions also. This chapter gives the novice users an idea about how Hadoop was initiated and what are the major revisions of it. Also this chapter describes in detail the architecture of Hadoop.

BACKGROUND

The most acute information management challenges stem from organizations (e.g., enterprises, government agencies, libraries, "smart" homes) relying on a large number of diverse, interrelated data sources, but having no way to manage their *dataspaces* (Franklin, Halevy, & Maier, 2005) in a convenient, integrated, or principled fashion. Michael Franklin et.al, (2005) highlighted the need for storage systems to accept all data formats and to provide APIs for data access that evolve based on the storage system's understanding of the data.

In the past years, (Dean & Ghemawat, 2004) at Google have implemented hundreds of special- purpose computations that process large amounts of raw data, such as crawled documents, web request logs, etc., to compute various kinds of derived data, such as inverted indices, various representations of the graph structure of web documents, summaries of the number of pages crawled per host, the set of most frequent queries in a given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time.

Robert Kallman et.al, (2008) developed H-Store, a next-generation OLTP system that operates on a distributed cluster of shared-nothing machines where the data resides entirely in main memory. But it needs a separate database design for the attributes- Table replication and Data partitioning. To solve the above problems, (Chang, Dean, Ghemawat, Hsieh, Wallach, Burrows... Gruber, 2008) developed BigTable which is distributed storage system for maintaining structured data of petabytes size across thousands of commodity servers. Later an open source equivalent to BigTable was created and it was called "Hadoop". Hadoop is an open source platform that brings the ability to cheaply process large amounts of data and it is more suitable for storing voluminous unstructured data.

11 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-history-and-architecture/216598

Related Content

Development of a New Means to Improve the Performance of Self-Organizing Maps

Vijaya Prabhagar Murugesanand Punniyamoorthy M. (2022). *International Journal of Data Analytics (pp. 1-16).*

www.irma-international.org/article/development-of-a-new-means-to-improve-the-performance-of-self-organizing-maps/307065

Evaluation of Optimum Parameters for Casting of Birla Lance Pipes

Dinesh S. Shinde, Ashnut Dutt, Ranjan Kumar Ghadai, Kanak Kalitaand Amer Nasr A. Elghaffar (2021). *Data-Driven Optimization of Manufacturing Processes (pp. 13-23).*

www.irma-international.org/chapter/evaluation-of-optimum-parameters-for-casting-of-birla-lance-pipes/269303

Comprehensive Analysis of State-of-the-Art CAD Tools and Techniques for Chronic Kidney Disease (CKD)

Mynapati Lakshmi Prasudha, Rakesh Kasumollaand Deepak Sukheja (2021). International Journal of Big Data and Analytics in Healthcare (pp. 1-12). www.irma-international.org/article/comprehensive-analysis-of-state-of-the-art-cad-tools-and-techniques-for-chronic-kidney-disease-ckd/287605

The Effect of Monetary Policy on Price Stability and Gross Domestic Product in Ghana: A Predictive Analytic Approach

Yaw Bediako, Patrick Ohemeng Gyaaseand Frank Gyimah Sackey (2022). *International Journal of Data Analytics (pp. 1-17).*

www.irma-international.org/article/the-effect-of-monetary-policy-on-price-stability-and-grossdomestic-product-in-ghana/307066

On the Use of Digital Platforms to Support SME Internationalization in the Context of Industrial Business Associations

Eric Costa, António Lucas Soaresand Jorge Pinho de Sousa (2019). *Handbook of Research on Expanding Business Opportunities With Information Systems and Analytics (pp. 66-94).*

www.irma-international.org/chapter/on-the-use-of-digital-platforms-to-support-smeinternationalization-in-the-context-of-industrial-business-associations/208559