# Chapter 5 Hadoop Distributed File System (HDFS)

## ABSTRACT

Hadoop Distributed File System, which is popularly known as HDFS, is a Java-based distributed file system running on commodity machines. HDFS is basically meant for storing Big Data over distributed commodity machines and getting the work done at a faster rate due to the processing of data in a distributed manner. Basically, HDFS has one name node (master node) and cluster of data nodes (slave nodes). The HDFS files are divided into blocks. The block is the minimum amount of data (64 MB) that can be read or written. The functions of the name node are to master the slave nodes, to maintain the file system, to control client access, and to have control of the replications. To ensure the availability of the name node, a standby name node is deployed by failover control and fencing is done to avoid the activation of the primary name node during failover. The functions of the data nodes are to store the data, serve the read and write requests, replicate the blocks, maintain the liveness of the node, ensure the storage policy, and maintain the block cache size. Also, it ensures the availability of data.

DOI: 10.4018/978-1-5225-3790-8.ch005

Copyright © 2019, IGI Global. Copying or distributing in print or electronic forms without written permission of IGI Global is prohibited.

## INTRODUCTION

HDFS is a Java based distributed file system running on commodity machines. It holds very large amount of data. In order to store large data, the files are split into smaller blocks and stored across multiple machines. This allows parallel processing. For example, if India wants to store aadhar card details of all people in its country, the names starting with 'A' can be stored in one server, the names starting with 'B' can be stored in server2, etc. HDFS demonstrated 200 PB of storage and a single cluster of 4500 servers. This chapter explains the architecture of HDFS. Also the salient features of HDFS are explained so that any reader can easily understand the architecture and use it.

# BACKGROUND

Hortonworks (Hortonworks, 2017) stated, "HDFS is a scalable, fault-tolerant, distributed storage system that works closely with a wide variety of concurrent data access applications, coordinated by YARN".

Cloudera (Cloudera Inc, 2017) supported Hadoop by stating:

HDFS is a fault-tolerant and self-healing distributed filesystem designed to turn a cluster of industry-standard servers into a massively scalable pool of storage. Developed specifically for large-scale data processing workloads where scalability, flexibility, and throughput are critical, HDFS accepts data in any format regardless of schema, optimizes for high-bandwidth streaming, and scales to proven deployments of 100PB and beyond.

Vangie Beal (Webopedia 2017) stated, "The primary objective of HDFS is to store data reliably even in the presence of failures including NameNode failures, DataNode failures and network partitions. ".

TechTarget (TechTarget, 2013) stated "Hadoop Distributed file system is designed to be highly fault-tolerant, facilitating the rapid transfer of data between compute nodes".

## **ARCHITECTURE OF HDFS**

HDFS cluster consists of:

25 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: <u>www.igi-</u> <u>global.com/chapter/hadoop-distributed-file-system-</u> hdfs/216600

### **Related Content**

### Machine Intelligence of Pi From Geometrical Figures With Variable Parameters Using SCILab

Ankita Mandal, Soumi Duttaand Sabyasachi Pramanik (2021). *Methodologies and Applications of Computational Statistics for Machine Intelligence (pp. 38-63).* www.irma-international.org/chapter/machine-intelligence-of-pi-from-geometrical-figures-with-variable-parameters-using-scilab/281161

#### PNC in 4D Object and Multi-Dimensional Data Modeling

(2017). Probabilistic Nodes Combination (PNC) for Object Modeling and Contour Reconstruction (pp. 209-234).

www.irma-international.org/chapter/pnc-in-4d-object-and-multi-dimensional-datamodeling/180359

#### Big Data Applications in Healthcare Administration

Joseph E. Kasten (2020). International Journal of Big Data and Analytics in Healthcare (pp. 12-37).

www.irma-international.org/article/big-data-applications-in-healthcare-administration/259986

#### Ontology-Based IoT Healthcare Systems (IHS) for Senior Citizens

Sakshi Guptaand Umang Singh (2021). *International Journal of Big Data and Analytics in Healthcare (pp. 1-17).* 

www.irma-international.org/article/ontology-based-iot-healthcare-systems-ihs-for-seniorcitizens/287604

# Big Data for Satellite Image Processing: Analytics, Tools, Modeling, and Challenges

Remya S., Ramasubbareddy Somula, Sravani Nalluri, Vaishali R.and Sasikala R. (2022). *Research Anthology on Big Data Analytics, Architectures, and Applications (pp. 1110-1124).* 

www.irma-international.org/chapter/big-data-for-satellite-image-processing/291028