

A Hadoop-driven Data Analysis System

Ashish Jindal
Rutgers University
Piscataway, NJ, USA
Email: ashish.jindal@rutgers.edu

Yikun Xian
Rutgers University
Piscataway, NJ, USA
Email: siriusxyk@gmail.com

Sanjivi Muttana
Rutgers University
Piscataway, NJ, USA
Email: sm1727@scarletmail.rutgers.edu

Abstract—The complexity of modern analytics needs is outstripping the available computing power of legacy systems. Hadoop complements this need for storing and analyzing huge sets of information by providing a platform for parallel processing of large data sets stored over multiple machines. This project aims to setup a hadoop infrastructure and demonstrate its power for big data analysis.

Index Terms—Hadoop-driven, Data Analysis, Infrastructure

I. PROJECT DESCRIPTION

Apache Hadoop is an open-source software framework written in Java for distributed storage and distributed processing of very large data sets on computer clusters built from commodity hardware. Hadoop stores data as it comes in - structured or unstructured - saving time on configuring data for relational databases. In our project we will store a large dataset in hadoop filesystem and use it for analysis of events from various sources like news, twitter etc. Our project falls in the category *Scalable Algorithms Infrastructure*. The main stumbling points for this project are identifying a large dataset for hadoop, configuring hadoop for multiple clusters and building a data analytics system over hadoop.

The project has four stages: requirement gathering, design, infrastructure implementation, and user interface.

A. Stage1 - The Requirement Gathering Stage

Following are the deliverables for this stage:

a) *General Description*: This project is about setting up an infrastructure for big data analytics using hadoop and demonstrate a data analytics application using event data from various sources like new, twitter etc.

b) *User Type 1*: People interested in analysing data that can't be stored on a single machine.

- User Interaction Modes: Hive query language.
- Real World Scenarios:
 - Scenario 1 Description: Processing of large server logs.
 - System Data Input for Scenario 1: Log files.
 - Input Data Types for Scenario 1: Structured log files.
 - System Data Output for Scenario 1: Information based on query.
 - Scenario 2 Description: Text mining.
 - System Data Input for Scenario 2: Big sources of textual information.
 - Input Data Types for Scenario 2: Text.
 - System Data Output for Scenario 2: Information based on query.

c) *Project Timeline and Division of Labor*: We will start with studying about hadoop and map reduce systems and then set up a hadoop system on local host. When we are successful in doing that we will find a small dataset, large enough to fit on a single system and start test analysis over it using hadoop. Then we will repeat the same process by setting up hadoop in pseudo distributed mode. Parallely two team members will start building the demo analytics application over hadoop. In the end we will setup hadoop on multiple clusters and deploy our application over it.

- Week 1 and Week 2: Work on requirement gathering and design of system using hadoop. Setup hadoop on local host as a single cluster and perform dummy analysis to test its functionality.
- Week 3 and week 4: Setup hadoop infrastructure over multiple clusters and create an analytics application demonstrating big data analysis using hadoop.
- Week 5: Test the application and prepare the necessary documentation project report.