# Apache Tajo

## tutorialspoint
### SIMPLY EASY LEARNING

## About the Tutorial

Apache Tajo is an open-source distributed data warehouse framework for Hadoop. Tajo was initially started by Gruter, a Hadoop-based infrastructure company in south Korea. Later, experts from Intel, Etsy, NASA, Cloudera, Hortonworks also contributed to the project.

Tajo refers to an ostrich in Korean language. In the year March 2014, Tajo was granted a top-level open source Apache project. This tutorial will explore the basics of Tajo and moving on, it will explain cluster setup, Tajo shell, SQL queries, integration with other big data technologies and finally conclude with some examples.

## Audience

Before proceeding with this tutorial, you must have a sound knowledge on core Java, any of the Linux OS, and DBMS.

## Prerequisites

This tutorial has been prepared for professionals aspiring to make a career in big data analytics. This tutorial will give you enough understanding on Apache Tajo.

## Disclaimer & Copyright

# Table of Contents

# 1. APACHE TAJO — INTRODUCTION

## Distributed Data Warehouse System

Data warehouse is a relational database that is designed for query and analysis rather than for transaction processing. It is a subject-oriented, integrated, time-variant, and non-volatile collection of data. This data helps analysts to take informed decisions in an organization but relational data volumes are increased day by day.

To overcome the challenges, distributed data warehouse system shares data across multiple data repositories for the purpose of Online Analytical Processing(OLAP). Each data warehouse may belong to one or more organizations. It performs load balancing and scalability. Metadata is replicated and centrally distributed.

Apache Tajo is a distributed data warehouse system which uses Hadoop Distributed File System (HDFS) as the storage layer and has its own query execution engine instead of MapReduce framework.

## Overview of SQL on Hadoop

Hadoop is an open-source framework that allows to store and process big data in a distributed environment. It is extremely fast and powerful. However, Hadoop has limited querying capabilities so its performance can be made even better with the help of SQL on Hadoop. This allows users to interact with Hadoop through easy SQL commands.

Some of the examples of SQL on Hadoop applications are Hive, Impala, Drill, Presto, Spark, HAWQ and Apache Tajo.

## What is Apache Tajo

Apache Tajo is a relational and distributed data processing framework. It is designed for low latency and scalable ad-hoc query analysis.

- Tajo supports standard SQL and various data formats. Most of the Tajo queries can be executed without any modification.

- Tajo has **fault-tolerance** through a restart mechanism for failed tasks and extensible query rewrite engine.

- Tajo performs the necessary **ETL (Extract Transform and Load process)** operations to summarize large datasets stored on HDFS. It is an alternative choice to Hive/Pig.

The latest version of Tajo has greater connectivity to Java programs and third-party databases such as Oracle and PostGreSQL.

# Features of Apache Tajo

Apache Tajo has the following features:

- Superior scalability and optimized performance
- Low latency
- User-defined functions
- Row/columnar storage processing framework.
- Compatibility with HiveQL and Hive MetaStore
- Simple data flow and easy maintenance.

# Benefits of Apache Tajo

Apache Tajo offers the following benefits:

- Easy to use
- Simplified architecture
- Cost-based query optimization
- Vectorized query execution plan
- Fast delivery
- Simple I/O mechanism and supports various type of storage.
- Fault tolerance

# Use Cases of Apache Tajo

The following are some of the use cases of Apache Tajo:

## Data warehousing and analysis

Korea's SK Telecom firm ran Tajo against 1.7 terabytes worth of data and found it could complete queries with greater speed than either Hive or Impala.

## Data discovery

The Korean music streaming service Melon uses Tajo for analytical processing. Tajo executes ETL (extract-transform-load process) jobs 1.5 to 10 times faster than Hive.

## Log analysis

Bluehole Studio, a Korean based company developed TERA — a fantasy multiplayer online game. The company uses Tajo for game log analysis and finding principal causes of service quality interrupts.

# Storage and Data Formats

Apache Tajo supports the following data formats:

- JSON
- Text file(CSV)
- Parquet
- Sequence File
- AVRO
- Protocol Buffer
- Apache Orc

Tajo supports the following storage formats:

- HDFS
- JDBC
- Amazon S3
- Apache HBase
- Elasticsearch

The following illustration depicts the architecture of Apache Tajo.



The following table describes each of the components in detail.

| Component | Description |
| --- | --- |
| Client | **Client** submits the SQL statements to the Tajo Master to get the result. |
| Master | Master is the main daemon. It is responsible for query planning and is the coordinator for workers. |
| Catalog server | Maintains the table and index descriptions. It is embedded in the Master daemon. The catalog server uses Apache Derby as the storage layer and connects via JDBC client. |

10

| | |
|---|---|
| Worker | Master node assigns task to worker nodes. TajoWorker processes data. As the number of TajoWorkers increases, the processing capacity also increases linearly. |
| Query Master | Tajo master assigns query to the Query Master. The Query Master is responsible for controlling a distributed execution plan. It launches the TaskRunner and schedules tasks to TaskRunner. The main role of the Query Master is to monitor the running tasks and report them to the Master node. |
| Node Managers | Manages the resource of the worker node. It decides on allocating requests to the node. |
| TaskRunner | Acts as a local query execution engine. It is used to run and monitor query process. The TaskRunner processes one task at a time.<br><br>It has the following three main attributes:<br><br>• Logical plan - An execution block which created the task.<br>• A fragment - an input path, an offset range, and schema.<br>• Fetches URIs |
| Query Executor | It is used to execute a query. |
| Storage service | Connects the underlying data storage to Tajo. |

## Workflow

Tajo uses Hadoop Distributed File System (HDFS) as the storage layer and has its own query execution engine instead of the MapReduce framework. A Tajo cluster consists of one master node and a number of workers across cluster nodes.

The master is mainly responsible for query planning and the coordinator for workers. The master divides a query into small tasks and assigns to workers. Each worker has a local query engine that executes a directed acyclic graph of physical operators.

In addition, Tajo can control distributed data flow more flexible than that of MapReduce and supports indexing techniques.

The web-based interface of Tajo has the following capabilities:

• Option to find how the submitted queries are planned

11

End of ebook preview

If you liked what you saw…

Buy it from our store @ https://store.tutorialspoint.com