

Apache Tez- Application Data processing Framework

ABHIJIT THAKRE¹

¹School of Informatics and Computing, Bloomington, IN 47408, U.S.A.

²Mechanical Engineer, Nagpur University, 2003

* Corresponding authors: abhijit.thakre@gmail.com

project-000, March 27, 2017

There are lot of advancement in Hadoop framework from Hadoop1.0 to Hadoop 2.0. Hadoop 2.0 is layered architecture provides the YARN layer responsible for the resource management opens up the gate for developing different application engines on top of it. Apache Tez is one of such open source framework build on the top of YARN designed to build data-flow driven runtimes[1]. This paper focusses mostly on introduction to Apache Tez framework. It provides the insight of the architecture used in building Tez. It also tried to cover the technologies using apache Tez as unifying framework and performance improvement achieved due to that. © 2017 <https://creativecommons.org/licenses/>. The authors verify that the text is not plagiarized.

Keywords: Cloud, I524

<https://github.com/cloudmesh/classes/blob/master/docs/source/format/report/report.pdf>

INTRODUCTION

In order to understand Tez as an framework we need to first dig into the history of Hadoop. Hadoop 1.0 have MapReduce as the central execution engine of its application. Any type of problem statement for analysis needs to be restructured to fit it to the map-reduce paradigm. It was also responsible for resource management and resource allocation. With Hadoop 2.0 these responsibilities got divided separately where YARN got the responsibility of general purpose resource management.

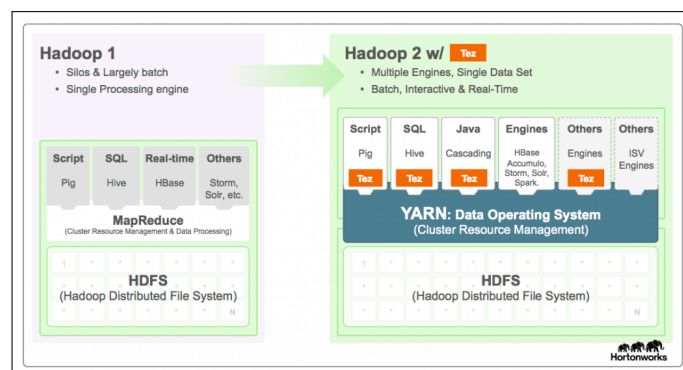


Fig. 1. [?]

Apache Hive, Apache PIG, CASCADING which was initially using Hadoop1 now run on Hadoop2 on Yarn. These listed data processing application including map reduce have certain set of requirement from the hadoop cluster in order to run efficiently.

This is where Tez came in picture. Tez takes care of running the data processing application's efficiency and performance leaving the end user to only concentrate on the business logic.

TEZ TERMINOLOGY

DAG – Direct Acyclic Graphs, it represent overall job. Vertex – Logical step in processing. It contains the details of user logic and dependent environment. Task – There can be multiple task in unit of work that vertex perform. Edge – This represents connection between producer and consumer vertices.

ARCHITECTURE AND IMPLEMENTATION

Tez was designed keeping in mind to address the problems which were not resolved by Hadoop. It was not build from the scratch but on the top of YARN layer to leavaraage the advantages and work that were done for years on Hadoop. So Tez leavaraage the discrete task based compute model, concept of data shuffling in map reduce, resource tentancy and multi-tenancy model and build in security from Hadoop.

Tez focuses mainly on below problem in addition.

- Without Tez, all the algorithm those needs to be executed on clusters needs to somehow translated to map-reduce api. This was impacting the efficiency and performance. However with Tez can naturally map the algorithm to execution engine in cluster.
- Tez provide additionally interface for various application and technology for data source and syncing.
- Performance.

TEZ-API

Tez provides below two API for defining the data processing.

DAG API

This API lets user define the structure for the computation. It lets user define the producer and consumers and how they talk to each other. This class of data processing application is represented as direct acyclic graphs

RUNTIME API

Using this API Layer Tez invokes the user code. This is where the actual code in the task to be executed is defined.

APPLICATION USING TEZ

Below are the product those are updated to be used on TEZ framework to run on YARN.

Apache MapReduce

MapReduce is simple but powerful way of data processing. Tez product comes with build in map-reduce support. The configuration needs to be updated on YARN cluster for map reduce. Tez has inbuilt map processor and reduce processor which provides the respective map reduce functionality.

Apache PIG

PigLatin is the scripting language provided by Apache PIG used to write complex ETL. Tez API can handle the multiple output which is usually the case for the procedural language like PigLatin, helps in keeping the code clean and maintainable.

Apache Hive

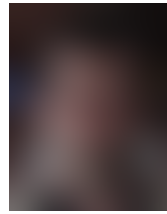
Apache Hive is used to convert the query written in HiveQL to map reduce and execute on hadoop cluster. The HiveQL is translation into map reduce format is often inefficient. Using Tez the trees translate directly into DAGS.

Apache Spark

Apache spark provides scala API for distributed data processing. The output of the spark is DAG of tasks that performance distributed computation. The end output of Apache spark i.e Spark DAG is successfully converted and executed using Tez API.

REFERENCES

AUTHOR BIOGRAPHIES



Abhijit Thakre received his BE (Mechanical) in 2003 from The University of Nagpur.