# Table of Contents

# Introduction to the Project

## A. Objective and Scope

The primary objective of this project is to leverage the advanced tools and techniques offered by the Hadoop Distributed Platform (HDP) for the processing and analysis of a comprehensive dataset. Students will utilize these tools to create analytical reports and visually compelling presentations that can serve as valuable resources for decision-makers.

Specifically, the project aims to conduct an in-depth analysis of data related to truck fleets to gain profound insights into the risk factors associated with trucking operations. Large truck accidents remain a leading cause of injuries and deaths in the United States. The overarching goal is to enhance safety within the trucking industry, addressing a critical issue that has led to numerous injuries and fatalities nationwide. The primary objective is to identify dangerous commercial truck drivers across the country.

## B. Significance of Cloudera Big Data Ecosystems

Cloudera's Big Data ecosystem plays a pivotal role in this project by providing a robust and scalable platform for data management and analysis. This ecosystem demonstrates its practical relevance by serving as the technological foundation for addressing real-world challenges within the trucking industry.

The project showcases how Cloudera's Big Data solutions can be effectively applied to improve driver safety, ensure regulatory compliance, and proficiently manage risks in the dynamic trucking sector.

# Project Phases

# Environment Setup and Data Integration

## 1) Installation of Tableau Drivers

This phase involves the installation of essential Tableau drivers, including JDBC and ODBC connectors. These drivers are crucial as they enable seamless connectivity between the project environment and the Hadoop File System (HDFS), ensuring efficient data access and processing.

## 2) Compilation of a Geolocation Dataset

In this stage, students will create the Data Definition Language (DDL) for the geolocation table in Hive (see step#3 below) and then proceed to load the geolocation dataset into this table. This dataset includes crucial information such as latitude, longitude, city, state, and essential truck fleet data, such as total mileage, average mileage, and fuel consumption. Additionally, students will create the necessary geolocation data set "geolocation.csv."

## 3) Create the geolocation Table in Hive

To prepare the geolocation data for analysis, students will create the geolocation table in Hive using the following Data Definition Language (DDL) command:

```
CREATE TABLE geolocation (
  truckid STRING,
  driverid STRING,
  event STRING,
  latitude DOUBLE,
  longitude DOUBLE,
  city STRING,
  state STRING,
  velocity BIGINT,
  event_ind BIGINT,
  idling_ind BIGINT
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
TBLPROPERTIES
("skip.header.line.count"="1");
```

This table will serve as the repository for geolocation data within the Hive environment, ensuring it is structured and organized for further use and analysis.

## 4) Load the Geolocation Dataset

With the 'geolocation' table in Hive, students will proceed to load the geolocation dataset, which includes the "geolocation.csv" file, into 'geolocation' Hive table. Loading the dataset into the Hive table is a critical step in preparing the data for subsequent analysis.

By completing this task, students will have the necessary geolocation data ready within the **geolocation** table in Hive, setting the stage for in-depth exploration and visualization.

## 5) Loading Data into Hadoop File System (HDFS)

In this task, students will transfer the input "geolocation.csv" file into Cloudera VM by utilizing the system copy command (scp). The students will use the concepts they have learned to upload the data into the HIVE instance by utilizing the "LOAD IN PATH" command.

## 6) Integration of Hive Table with Tableau

The integration of the Hive table with Tableau is a pivotal step in the project's workflow. This integration ensures that the data stored in HDFS becomes visible and accessible within the Tableau environment, facilitating comprehensive data analysis, exploration, and visualization.

To achieve this integration, students will identify the correct driver, install the ODBC driver for Impala, and load the tables into Tableau via ETL transactions. This will enable the creation of the necessary charts and visualizations needed for the project.

## 7) Steps for installing Impala ODBC Driver

a) Go to the  https://www.cloudera.com/downloads.html
b) Download Impala ODBC Drive only.

If you are using Tableau as an analytical tool, follow the below steps to connect to HDFS.
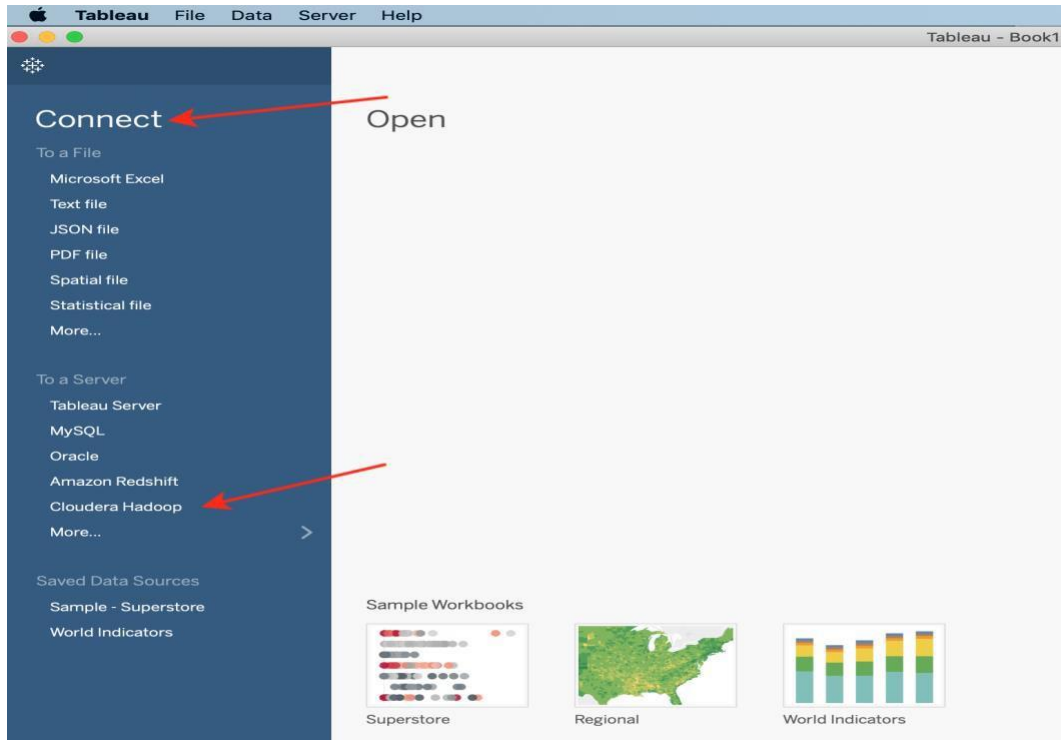


**Database Drivers**
The Cloudera ODBC and JDBC Drivers for Hive and Impala enable your enterprise users to access Hadoop data through Business Intelligence (BI) applications with ODBC/JDBC support.
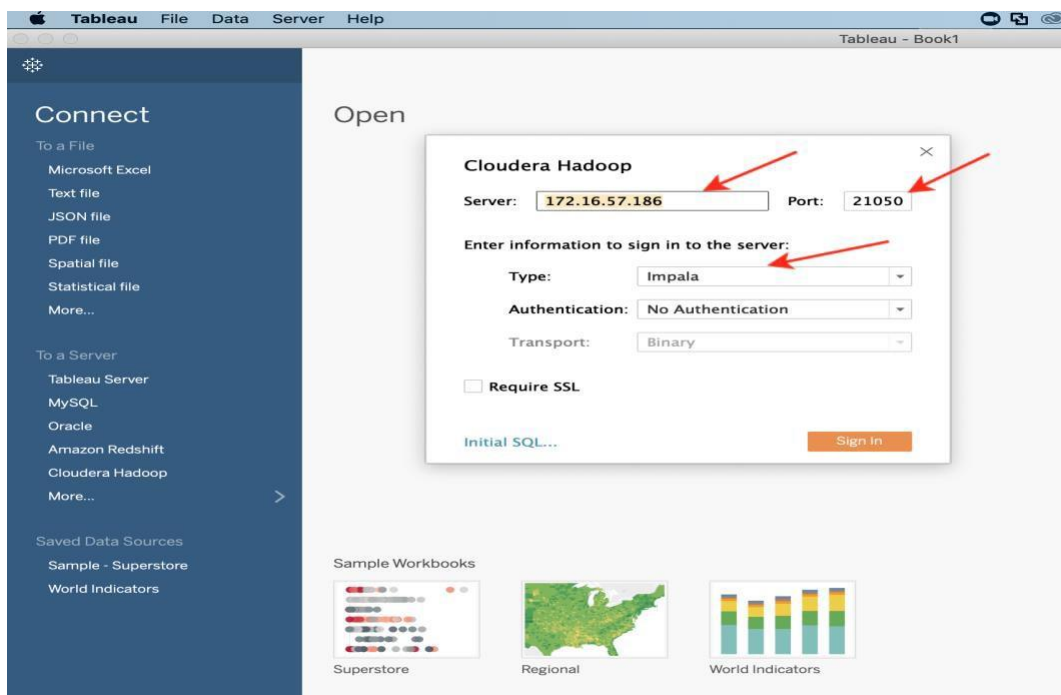
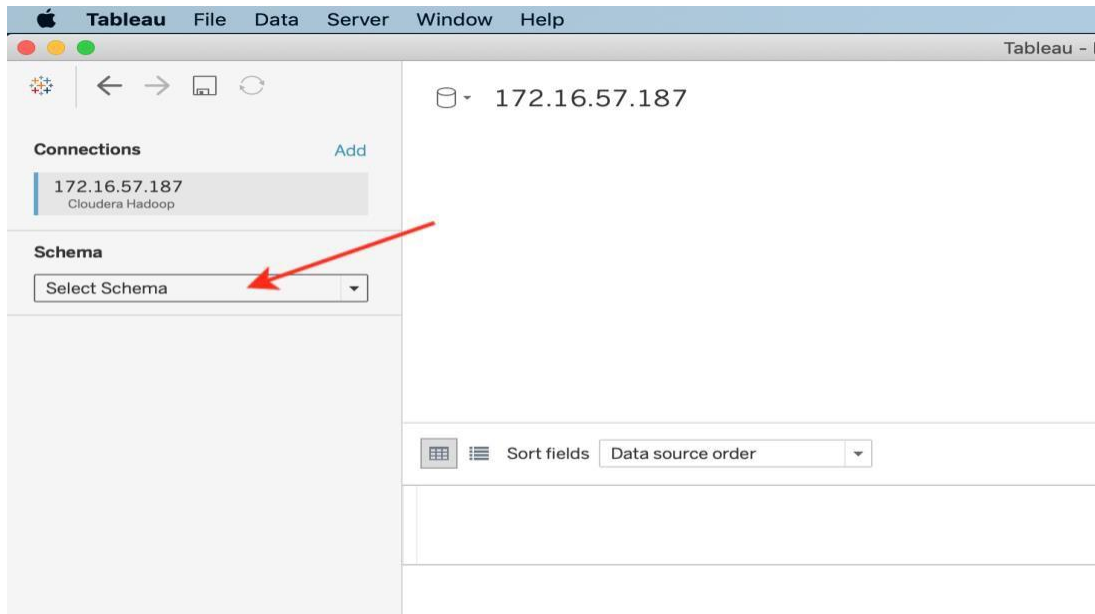**Hive ODBC Driver Downloads** >
**Hive JDBC Driver Downloads** >
**Impala ODBC Driver Downloads** >
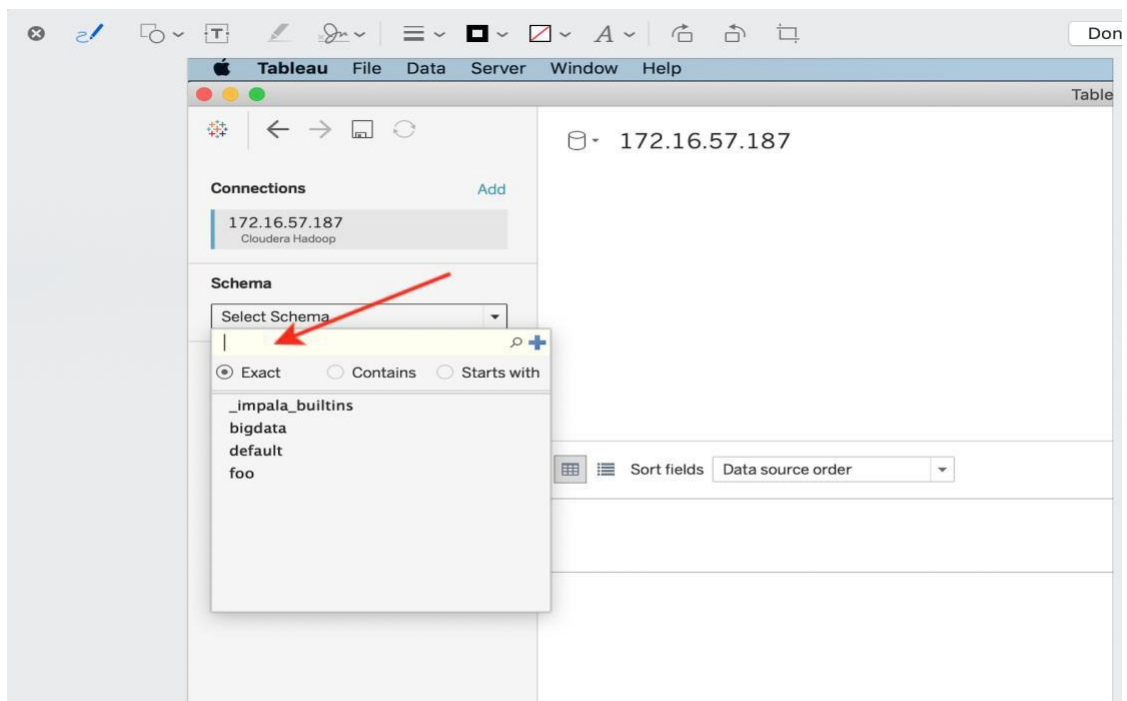**Impala JDBC Driver Downloads** >

From your Cloudera Image, run the command line 'hostname -I' to locate the IP address associated with your Instance. Copy the IP address and inserted into server bar as it shows below. Make sure you are using the port# 21050 and the Type = Impala.
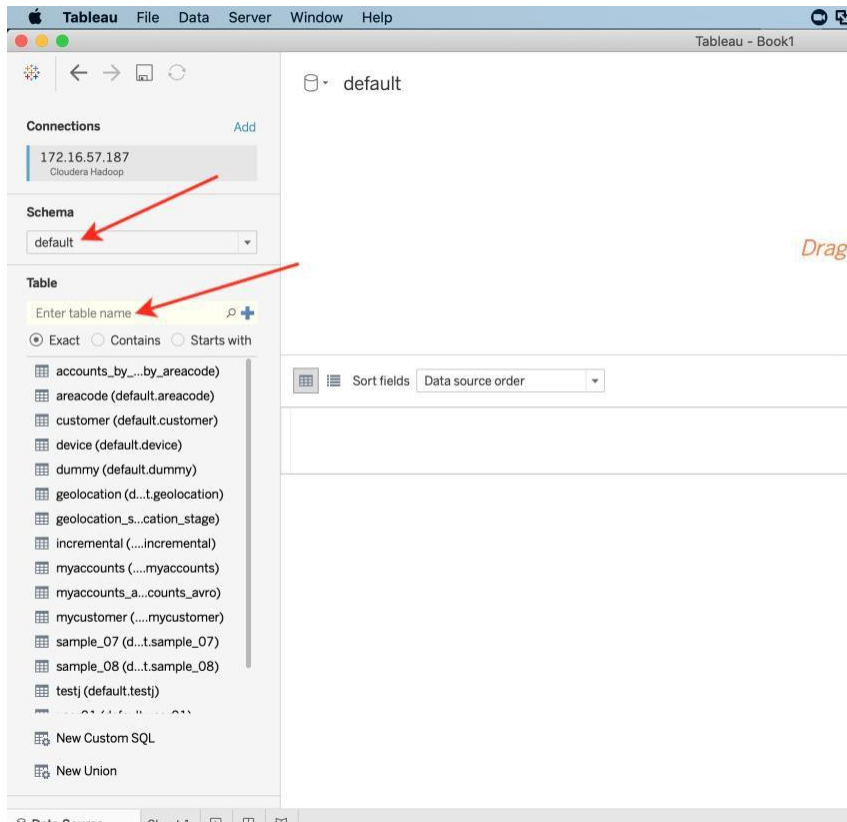
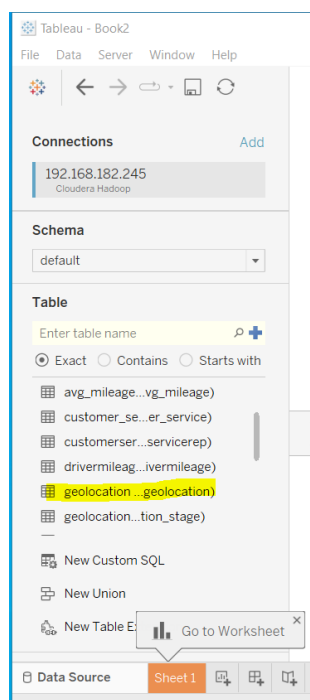Place the cursor at the Schema Bar and hit Enter.

Place the cursor at the Table bar and hit Enter. Then you will start getting the Hive/Impala tables.



Drag geolocation table to the middle of the sheet and click on "Update now "as it shows below.

You should have geolocation data migrated into Tableau as it shows below.