# Analysis of the Factors that Affect Bus Delay in Toronto

## Team Members

Jaydenn Chang    N01511476
Kaiyan Chen      N01489178
Mira Philip      N01495720
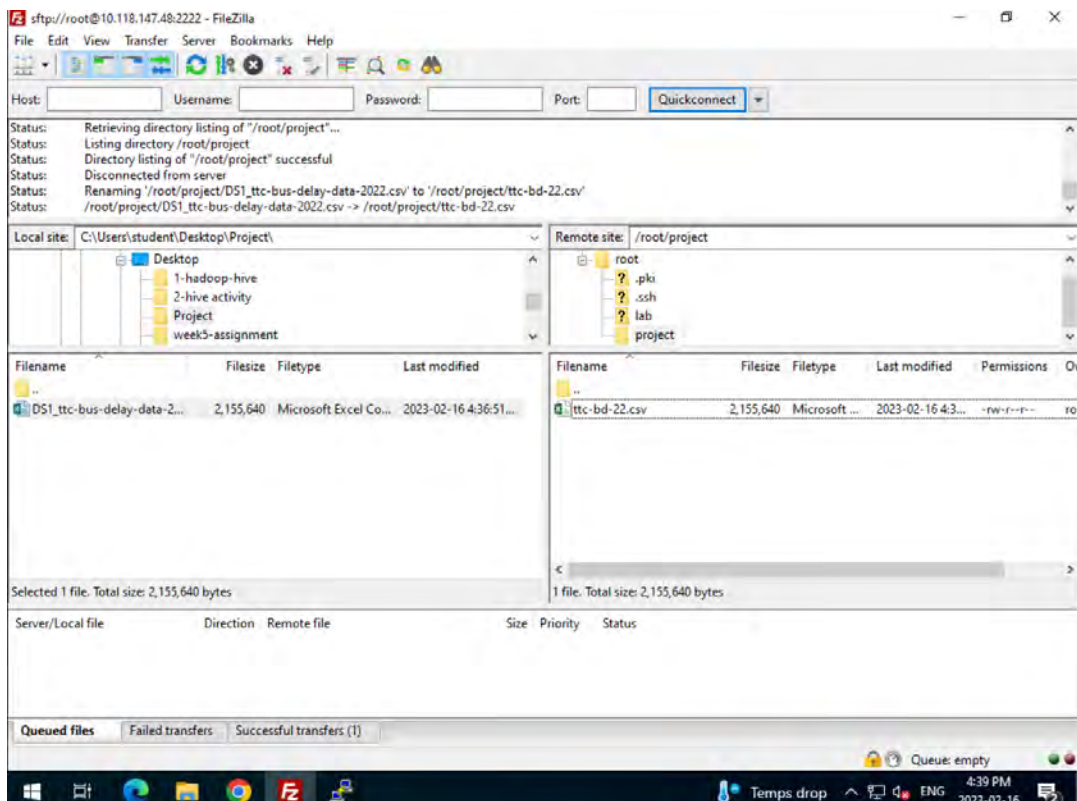Simul Bista      N01489966

## Problem Statement:

The project aims to analyze the factors that affect bus delay in Toronto based on the dataset that is sourced from the Toronto Transit Commission (TTC), which provides detailed records of bus delays across the city. It focuses on four key factors, namely bus route, days of the week, time, and incident type. The findings of the project to can be used to inform policy and decision-makers to improve the reliability and efficiency of the public transit system.

## Steps

1. **HDFS**
   a. First, we create a directory in the sandbox called project i.e., /root/project and using FileZilla, we upload our dataset into the sandbox.

b. We create a project directory in Hadoop using the command:
   hadoop fs -mkdir /user/root/project

```
[root@sandbox-hdp project]# hadoop fs -mkdir /user/root/project
[root@sandbox-hdp project]# hadoop fs -ls /user/root
Found 7 items
drwx------    - root root          0 2023-02-11 12:00 /user/root/.Trash
drwxr-xr-x    - root root          0 2023-01-30 21:06 /user/root/.hiveJars
drwx------    - root root          0 2023-01-23 20:37 /user/root/.staging
drwxr-xr-x    - root root          0 2023-02-10 20:47 /user/root/data
drwxr-xr-x    - root root          0 2023-02-13 20:49 /user/root/lab
drwxr-xr-x    - root root          0 2023-02-16 21:42 /user/root/project
drwxr-xr-x    - root root          0 2023-02-10 20:51 /user/root/twitter
[root@sandbox-hdp project]#
```

c. We generally give read and write access permission using chmod to the folder, however, we already being root, we don't do that as of now.

d. We then load the data from the sandbox to HDFS using the command:
   hadoop fs -put /root/project/ttc-bd-22.csv /user/root/project

```
[root@sandbox-hdp project]# hadoop fs -put /root/project/ttc-bd-22.csv /user/root/project
[root@sandbox-hdp project]# hadoop fs -ls /user/root
Found 7 items
drwx------    - root root          0 2023-02-11 12:00 /user/root/.Trash
drwxr-xr-x    - root root          0 2023-01-30 21:06 /user/root/.hiveJars
drwx------    - root root          0 2023-01-23 20:37 /user/root/.staging
drwxr-xr-x    - root root          0 2023-02-10 20:47 /user/root/data
drwxr-xr-x    - root root          0 2023-02-13 20:49 /user/root/lab
drwxr-xr-x    - root root          0 2023-02-16 21:43 /user/root/project
drwxr-xr-x    - root root          0 2023-02-10 20:51 /user/root/twitter
[root@sandbox-hdp project]# hadoop fs -ls /user/root/project
Found 1 items
-rw-r--r--    1 root root    2155640 2023-02-16 21:43 /user/root/project/ttc-bd-22.csv
[root@sandbox-hdp project]#
```

2. **HIVE**
   a. Access hive using the command:
      Hive

   b. Now, we create a database called ttc using the following command:

      CREATE DATABASE ttc;

   c. We can verify that the database has been created using:

      SHOW DATABASES;

d. Before creating the table, we must make sure that we are inside the correct database since by default a default database is selected. We do this using:
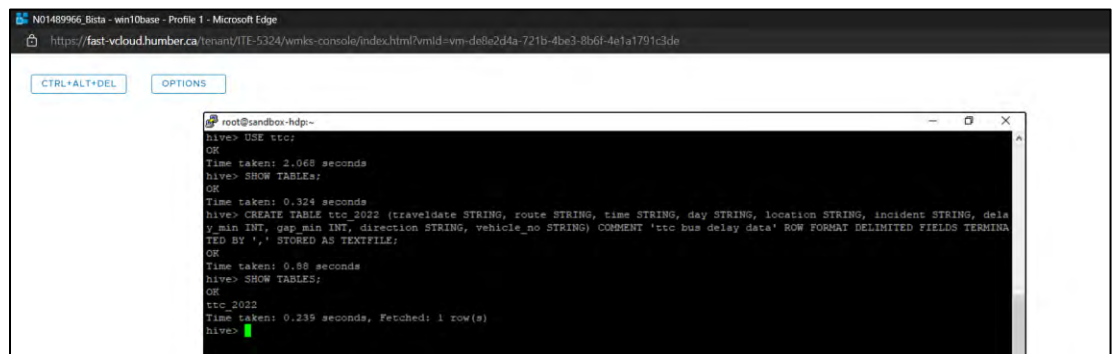
USE ttc;

e. We then create the table in which we want to load the data from our dataset which is in hdfs. The command is:

CREATE EXTERNAL TABLE ttc_bus_delay_2022 (traveldate STRING,route STRING, time TIMESTAMP, day STRING, location STRING, incident STRING, delay_min INT, gap_min INT, vehicle_no STRING);

f. We can check if the table has been created using the command:

SHOW TABLES IN ttc;



g. Now its time to load the data into the table, we do that using:

LOAD DATA INPATH '/user/root/project/ttc-bd-22.csv' OVERWRITE INTO TABLE ttc_bus_delay_2022;

h. Use a select query to view the first couple of records that has been loaded into the table:

SELECT * FROM ttc_bus_delay_2022 LIMIT 10;

**3. Zeppelin (Spark)**

    a. We create the data frame called ttc from the given dataset (csv file)

b. We then clean the data - remove null columns, perform some column renaming and then format the data in the time column to show the hours only (truncating the mins).



c. Next, we check the data (a couple of columns that we just cleaned) to verify that everything is as planned. And then we convert the data frame to a temp view to start visualizing and gain some meaningful information from the representation.

d. **Visualization 1:** We figured out the time of the day in which most bus delays happened. The result shows that around 2-5pm was most likely for the bus delay to happen.

e. **Visualization 2:** We figure out the day of the week in which most bus delays occurred. The result shows the day of the week doesn't show significant difference.

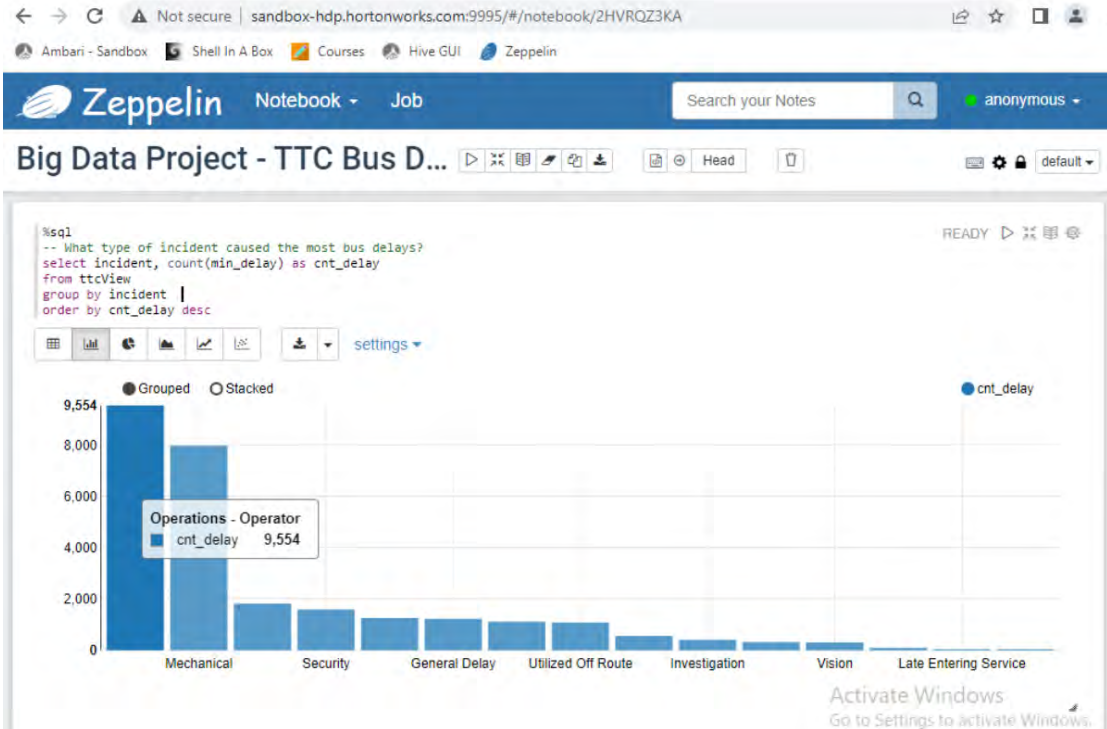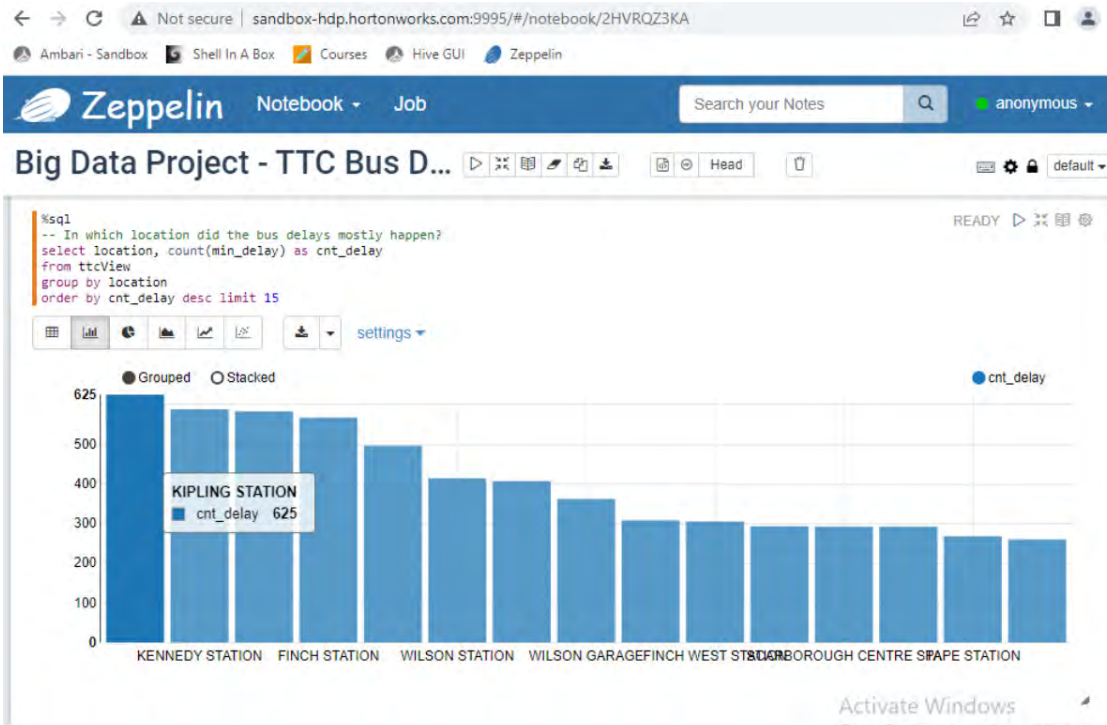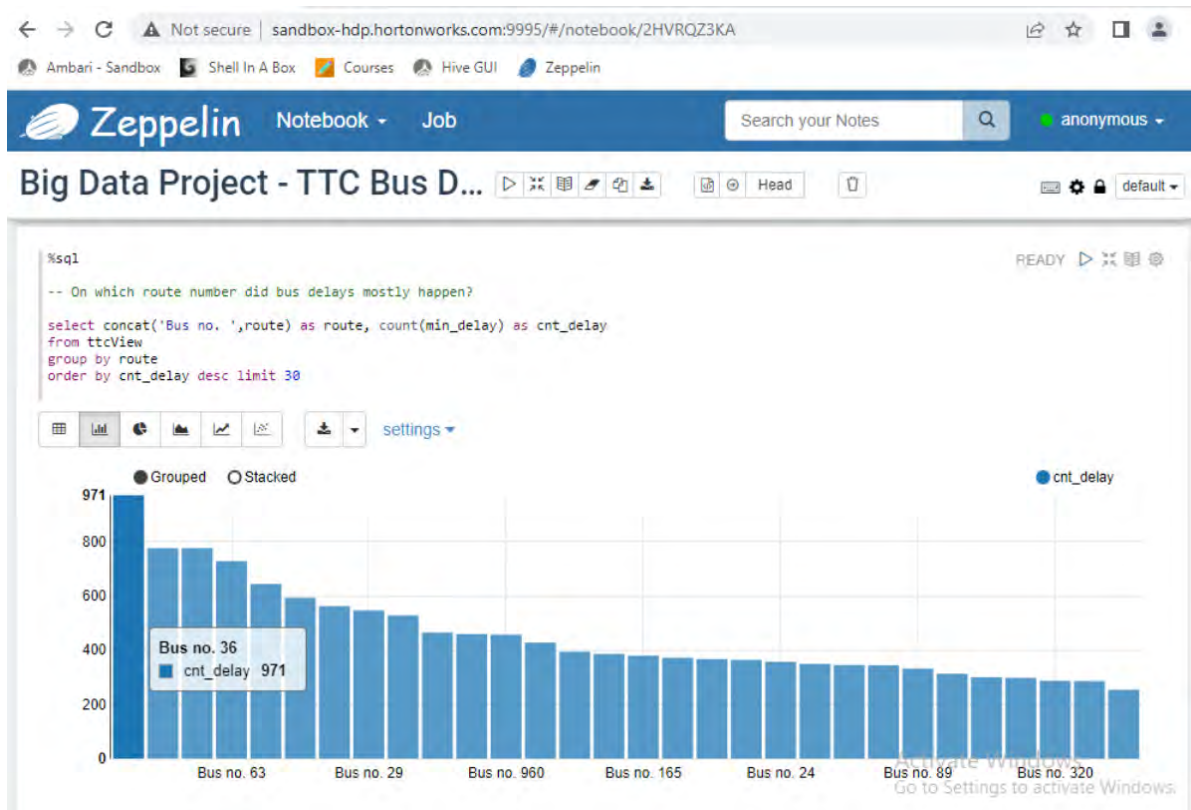f. **Visualization 3:** Next, we figure out the incident type which caused the most bus delays. The result shows that "Operations/Operators" and the "Mechanical" are the top reason behind bus delays.

g. **Visualization 4:** We figured out the location(station) where the most bus delays happened. The result shows that at Kipling station and Kennedy station are two major stations where the bus delay happened. They happened both to be the end station of the subway lines.

h. **Visualization 5:** We noticed on which routes the most bus delay happened. The result shows that its Bus Route 36.



**Conclusion**:

Bus delays mostly happened in the afternoon around 2-5pm at the end station of subway line 2 with the majority cause categories of "Operations/Operators" and "Mechanical." We can focus on these causes and dive deeper into the each reason to put together a further conclusion and possibly solutions to the bus delays.