# NYC Real-Time Traffic Speed Analysis
## Overspeed Rate and Jam Rate Analysis

Zhengda Liu, Eva Song, Jingwei Zhang, Shengyao Ye, Tianle Qiu

*2024 Fall STAT605 Group 2*

# 1 Introduction

This report examines NYC real-time traffic data, focusing on overspeed and congestion rates. We calculated these rates and applied clustering to uncover traffic patterns. Persistent congestion on streets like 11th Avenue in Manhattan highlights shortcomings in traffic management. Results are presented through interactive visualizations, offering insights for urban traffic planning and policy improvements.

# 2 Body

## 2.1 Data Description

We used the NYC Real-Time Traffic Speed Data from Kaggle, a 28.36GB CSV file with 13 columns. Each record represents the average speed of vehicles traveling between endpoints on a link, defined as a sequence of latitude and longitude points.
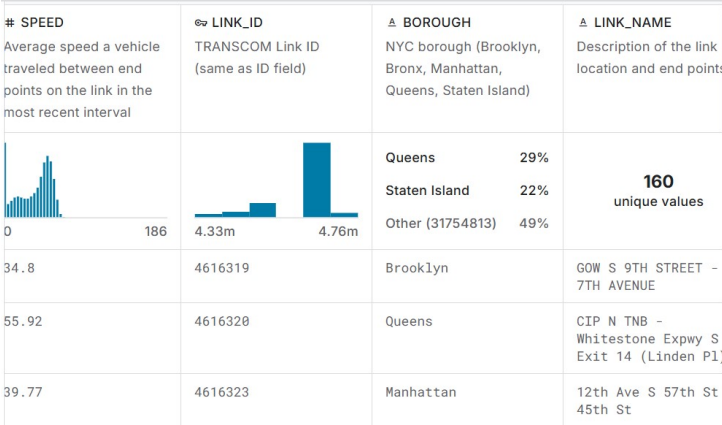


Figure 1: Dataset Overview

Since the Kaggle dataset is not updated, we identified an updated version on NYC Open-Data. While this version was not analyzed in this study, it is noted for future use.

## 2.2 Statistical Computation

All computations were performed using CHTC, as outlined below. The tasks include **A. computing summary statistics** (overspeed rate and jam rate) for each link, and **B. splitting data by months** for map visualization.
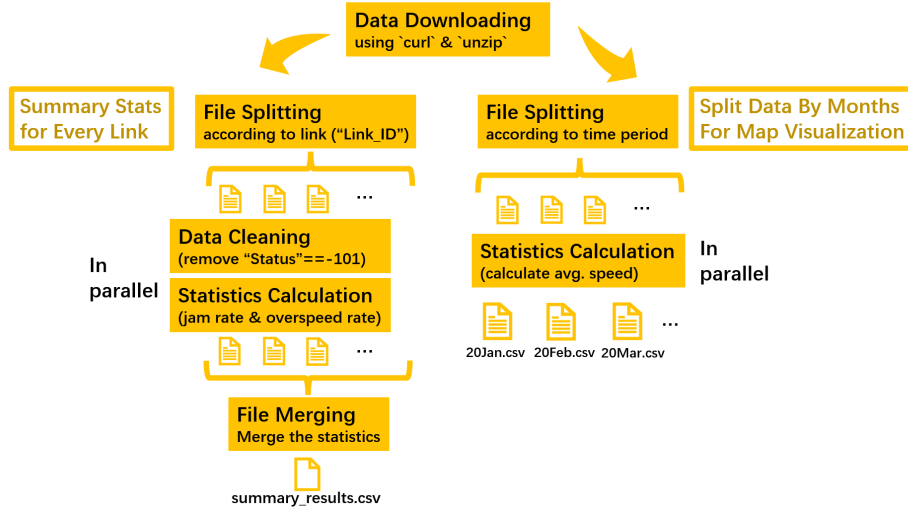
Figure 2: Computation Flowchart for Tasks A and B

### 2.2.1 Data Acquisition

Kaggle data was downloaded using `curl` and decompressed with `unzip`. After acquisition, two DAG files were executed to process tasks A and B.

### 2.2.2 Data Splitting and Cleaning

Data was split by key fields: `LinkID` for task A and month for task B. Records with `status` = -101, indicating zero speed, were removed as they distort summary statistics.

### 2.2.3 Computation

For task A, overspeed rate and jam rate were computed as:

$$OverspeedRate = \frac{\sum_i I(s_i > 60)}{m}, \quad JamRate = \frac{\sum_i I(s_i < 10)}{m},$$

where $s_i$ is the speed of the $i$-th record, and $m$ is the total number of records. The results were merged into a single CSV for visualization.

For task B, average speeds for each link per minute were computed. The outputs included multiple files, each representing link speeds for a specific month.

## 2.3 Result

All the codes and results can be seen in our repository. Below is the command to clone it:

Listing 1: Git Clone Command

```
git clone https://github.com/JustinDs0205/STAT605_Final.git
```

### 2.3.1 Graphical summaries

We developed interactive tools to visualize traffic patterns:

**Overspeed and Jam Rate Visualization**

View the visualization here.

**Interactive Traffic Map**

Users can select streets and date ranges to visualize traffic speeds using an interactive map powered by Leaflet. Explore the map here.

**Cluster Analysis**

Overspeed rates and jam rates for clustering road sections with similar traffic patterns. This helps authorities identify and manage roads with comparable issues more effectively.
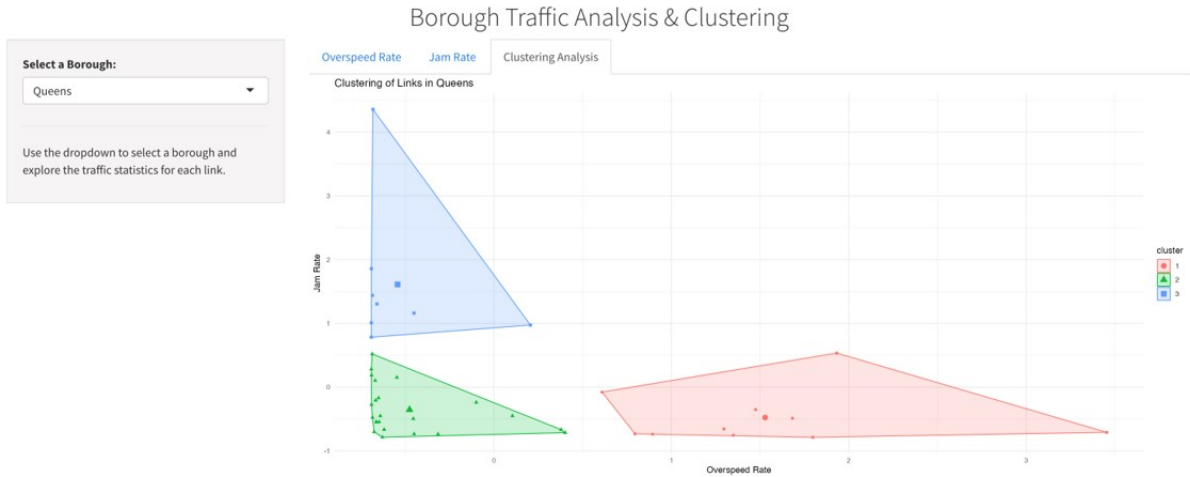


Figure 3: Cluster Analysis Result

Table 1: Cluster Meansfor Overspeed Rate and Jam Rate

| Cluster | Overspeed Rate | Jam Rate |
|---------|----------------|----------|
| 1 | 0.099275982 | 0.01440353 |
| 2 | 0.009634136 | 0.01927453 |
| 3 | 0.006628807 | 0.09819441 |

Cluster 1 consists of road segments with high overspeed rates and low congestion, requiring enhanced traffic monitoring to address speeding. Cluster 2 includes stable road segments with minimal speeding or congestion issues. Cluster 3 comprises road segments with high congestion and low overspeed rates, suggesting the need for road optimization to reduce congestion.

### 2.3.2 Statistical computation

The number of jobs, the typical job time, memory, and disk space required are described as follows.

Table 2: Job Resource Usage

| Job | Time | Memory | Disk |
|---|---|---|---|
| getKdata.sub | 1 hour 2 mins | 1GB | 40GB |
| A_pre.sub | 1 hour 17 mins | 1GB | 40GB |
| A_para.sub | 58 mins | 1GB | 40GB |
| B_para.sub | 1 hour 35 mins | 1GB | 40GB |

### 2.3.3 News Validation

Streets with high congestion rates in Manhattan, like 11th Avenue, reveal issues in policy implementation and enforcement. A Reddit post and a CBS News report highlight ongoing traffic challenges, with drivers ignoring lane discipline. Our visualization confirms that 11th Avenue has the highest congestion rate, underscoring the need for improved traffic management and stricter policy enforcement.
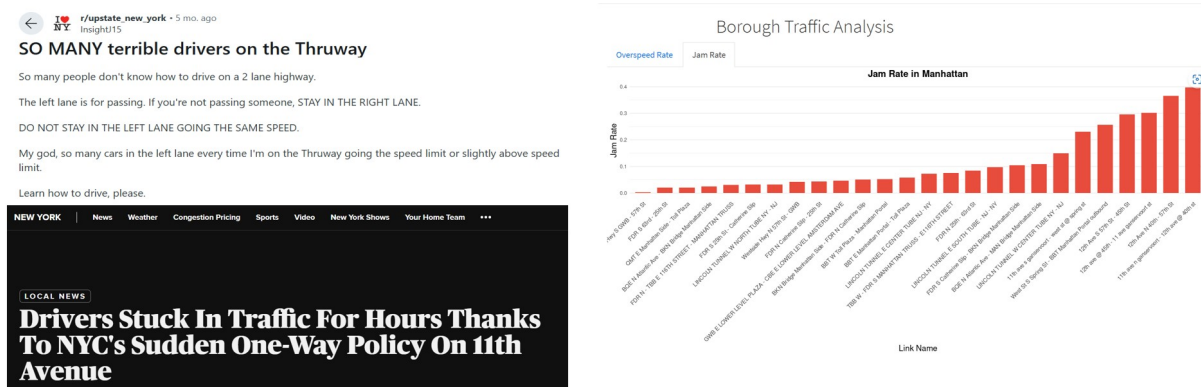


Figure 4: Validation for news and plots

### 2.3.4 Weaknesses

During data cleaning, we did not remove data related to special occasions, such as farmers markets or festival parades. Users can only know the overall historical congestion level of a road, but cannot get real-time information.

# 3 Conclusion

This report analyzes NYC traffic speed data, focusing on overspeed and congestion rates across streets. High congestion on streets like 11th Avenue in Manhattan highlights inefficiencies in traffic management and policy enforcement.

Clustering analysis identified patterns among streets with similar traffic characteristics, providing valuable insights into traffic behavior. The findings suggest that targeted measures, such as enhanced lane management and stricter enforcement, could mitigate congestion and overspeeding.

Future research could integrate additional factors, such as weather, to achieve a more comprehensive analysis. Predictive models may also enable real-time traffic monitoring using NYC OpenData.

# 4 Contributions

| Member | Proposal | Coding | Presentation | Report | PPT |
|---|---|---|---|---|---|
| Zhengda Liu | 1 | 1 | 1 | 0.6 | 0.4 |
| Eva Song | 0.6 | 1 | 1 | 1 | 0.2 |
| Jingwei Zhang | 1 | 0 | 1 | 0 | 1 |
| Shengyao Ye | 0 | 0.2 | 1 | 0 | 1 |
| Tianle Qiu | 0 | 0 | 1 | 1 | 0 |

Table 3: Group Member Contribution Chart

- In the chart above, 1 = full contribution, 0.1–0.9 = partial contribution, 0 = no contribution.