# 311 Complaint: Street Condition

# &

# Motor Vehicle Collisions Crashes

## Narrative Description:

The first dataset that was selected contains 311 service requests from 2010 to present in New York City. According to the 311 Service Requests data, in 2017, street condition was the 5th most common complaint issue with a total of 93,270 complaints for the year alone. Therefore, this specific complaint type was chosen to explore whether vehicle accidents are more likely to occur depending on road conditions.

Street condition complaints could be key indicators in predicting the likelihood of collisions and essential in preventing future crashes. As a result, a second dataset has been selected that contains the details on motor vehicle collision events in New York City. By merging these two datasets, we hope to find correlations between complaints of road conditions, which included potholes, failed street repairs, and more, and vehicle collisions.

For instance, with the information that we have, we hope to easily identify which locations should often receive inspections to assess the road conditions due to the higher volume of complaints. As well as what time of the year would streets have hazardous road conditions? For instance, during the fourth quarter of the year where there are more snow storms, streets would be more damaged than the rest of the year. By having access to this information, it would be easier for us to prepare and improve the road conditions and thus, reduce vehicle accidents drastically.

**Data Sources**

1. 311 Service Requests

2. Motor Vehicle Collisions Crashes

**Extracted Sample Data:**

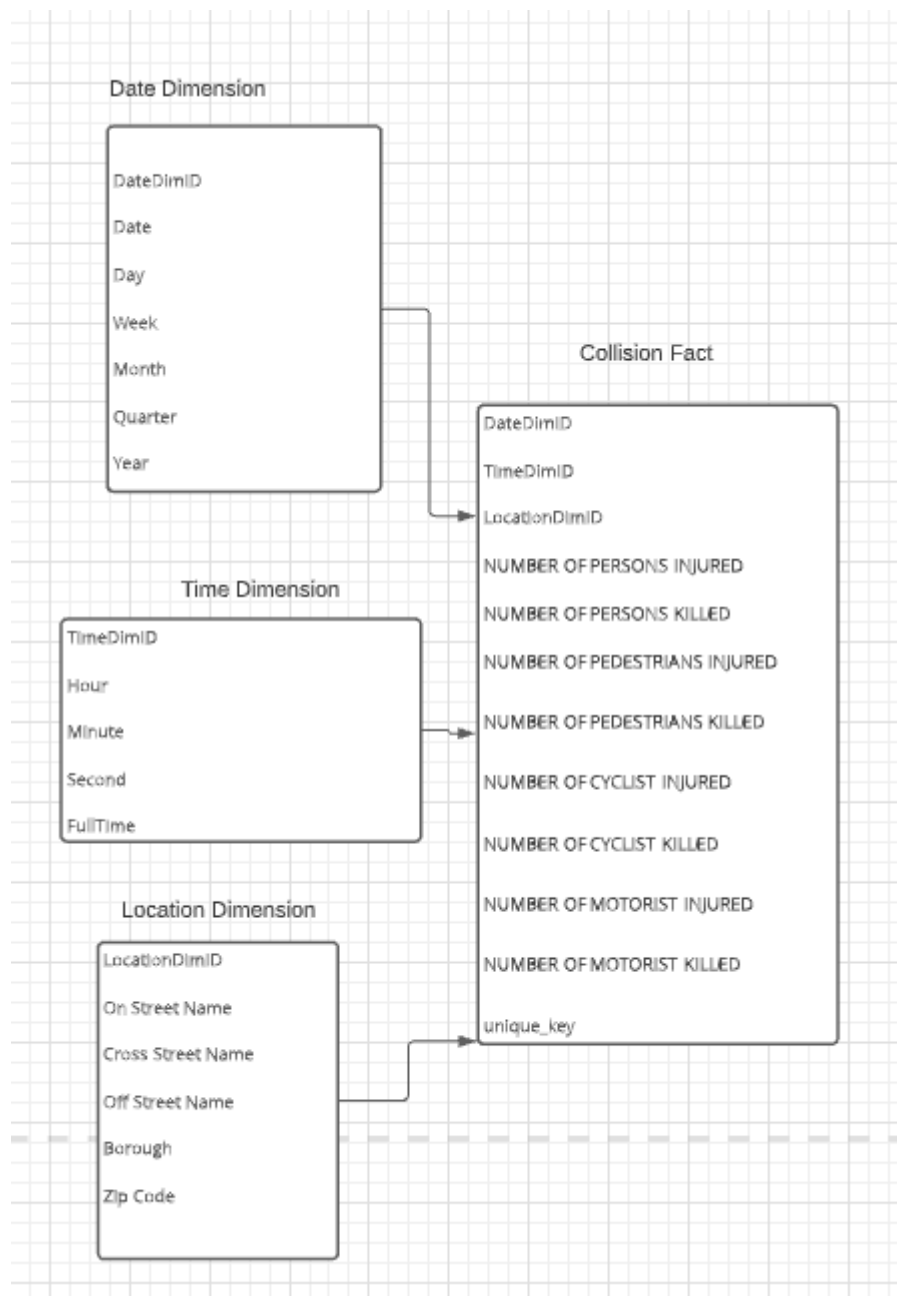1. 311 Service Requests

2. Motor Vehicle Collisions

**KPIs**

1. Street Condition Complaints and Collisions Count by Zip Code

2. Number of Persons Killed By Borough

3. Number of Persons Injured by Borough

4. Number of Deaths by Road Condition Description

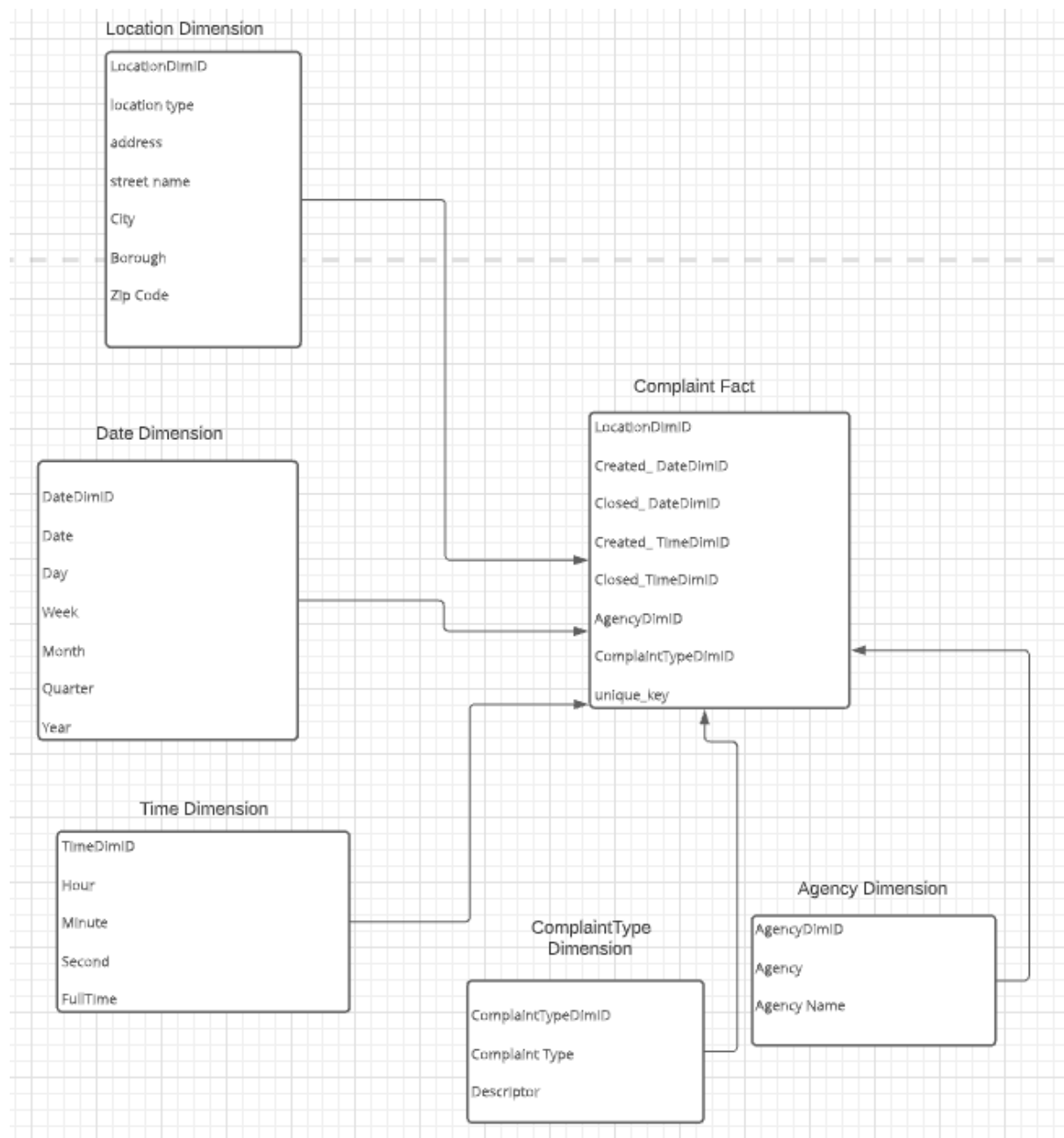5. Occurrences of Road Condition Complaints Based on Collisions

**Dimensional Model:** Link to LucidChart Dimensional Model

Some KPIs involve frequency monitoring based on time. Therefore, transaction grain has been selected for every record in the fact table for each of the transactions. The only attribute that would most likely change is location; however, infrequently. Thus, type 3 slowly changing dimensions has been selected.

## Collision Fact



**Date Dimension**
- DateDimID
- Date
- Day
- Week
- Month
- Quarter
- Year

**Time Dimension**
- TimeDimID
- Hour
- Minute
- Second
- FullTime

**Location Dimension**
- LocationDimID
- On Street Name
- Cross Street Name
- Off Street Name
- Borough
- Zip Code

**Collision Fact**
- DateDimID
- TimeDimID
- LocationDimID
- NUMBER OF PERSONS INJURED
- NUMBER OF PERSONS KILLED
- NUMBER OF PEDESTRIANS INJURED
- NUMBER OF PEDESTRIANS KILLED
- NUMBER OF CYCLIST INJURED
- NUMBER OF CYCLIST KILLED
- NUMBER OF MOTORIST INJURED
- NUMBER OF MOTORIST KILLED
- unique_key

# Complaint Fact



## Location Dimension

- LocationDimID
- location type
- address
- street name
- City
- Borough
- Zip Code

## Date Dimension

- DateDimID
- Date
- Day
- Week
- Month
- Quarter
- Year

## Time Dimension

- TimeDimID
- Hour
- Minute
- Second
- FullTime

## Complaint Fact

- LocationDimID
- Created_ DateDimID
- Closed_ DateDimID
- Created_ TimeDimID
- Closed_TimeDimID
- AgencyDimID
- ComplaintTypeDimID
- unique_key

## ComplaintType Dimension

- ComplaintTypeDimID
- Complaint Type
- Descriptor

## Agency Dimension

- AgencyDimID
- Agency
- Agency Name

**Final Dimensional Schema (Integrated Model)**

The Collision and Complaint fact tables share Date, Time and Location dimensions.

The Complaint fact table also contains Agency and ComplaintType dimensions.

**Date Dimension**
- DateDimID
- Date
- Day
- Week
- Month
- Quarter
- Year

**Time Dimension**
- TimeDimID
- Hour
- Minute
- Second
- FullTime

**Agency Dimension**
- AgencyDimID
- Agency
- Agency Name

**Complaint Fact**
- LocationDimID
- Created_DateDimID
- Closed_DateDimID
- Created_TimeDimID
- Closed_TimeDimID
- AgencyDimID
- ComplaintTypeDimID
- unique_key

**Collision Fact**
- DateDimID
- TimeDimID
- LocationDimID
- NUMBER OF PERSONS INJURED
- NUMBER OF PERSONS KILLED
- NUMBER OF PEDESTRIANS INJURED
- NUMBER OF PEDESTRIANS KILLED
- NUMBER OF CYCLIST INJURED
- NUMBER OF CYCLIST KILLED
- NUMBER OF MOTORIST INJURED
- NUMBER OF MOTORIST KILLED
- unique_key

**Location Dimension**
- LocationDimID
- Borough
- Zip Code

**ComplaintType Dimension**
- ComplaintTypeDimID
- Complaint Type
- Descriptor

**Data Profiling:** Data Profiling

A Pandas module was used for profiling. This was a quick exploratory analysis step we took before proceeding. These multiple visualisations saved a lot of time and gave us an understanding of the distribution of each variable.

**ETL Programming For Dimensions and Fact Tables:**

**Lineage Graph**

With the ETL tool DBT, we created our dimension models and fact tables. The diagram shows the relationship between each dimension to the fact tables.

**Complaint Type Dimension**

A surrogate key is created and the complaint type and descriptor columns are created. It will also display these columns in the dimension. Our "where" clause only takes records for the specific complaint type, "Street Condition".

## Location Dimension

The datasets are integrated by zip code and borough. After unioning the separate location tables for the datasets, the location dimension was created with a surrogate key.

**Time Dimension**

A CSV file was integrated that contains time_dim_id, fulltime, hours, minutes, and seconds. It was uploaded to BigQuery and selected all records from this dataset to create the time dimension.

**Date Dimension**

The dbt utils package was installed to construct a date dimension with the dimension ID, full date, the number sequence of day (out of 365 days) in the year, the day of the week, the number sequence of week (out of 52 weeks) the day falls in in the year, the month of the year, the quarter in the year, and the year.

Code Source: gitlab-data/analytics/date

**Agency Dimension**

Agency name and agency as the surrogate key from the 311 Service Requests dataset.

It will also display these columns with the original values in the dimension.

**Collision Fact**

We have selected the dimensions ID, the degenerated dimension, and the columns from the Collision dataset that will add valuable information to create the KPIs. The datasets are joined with the related dimensions (Date, Location, and Time) on the shared columns in order to display the table.

**Complaint Fact**

We have selected the dimension IDs and the degenerated dimension that are critical for this fact table. The dataset are joined with the related dimensions (Date, Location, Agency, Complaint_Type, and Time) on the shared columns in order to display the table.



**Visualisations of KPIs: Presentation: Real Time Insights**

- Data Lake: Google BigQuery

- Data Warehouse: Google BigQuery

- ETL Tools: DBT

- BI Application: Tableau

- Programming Language: SQL & Python

- Diagramming Application: LucidChart