

MTA DAILY RIDERSHIP

ANALYSIS & FORECASTING

1- DATA DESCRIPTION

MTA Daily Ridership Data:

Our dataset tracks daily ridership and traffic volumes across multiple NYC transit services starting from March 1, 2020, capturing the impact and recovery from the COVID-19 pandemic.

Included Transit Modes:

- Subways
- Buses
- Long Island Rail Road (LIRR)
- Metro-North Railroad
- Staten Island Railway
- Access-A-Ride (Paratransit)
- Bridges & Tunnels (Vehicle traffic)

Key Metrics (for each mode):

- Total Estimated Ridership/Traffic
- % of Comparable Pre-Pandemic Day
- (shows how current ridership compares to typical pre-COVID levels)

Timeframe:

- Covers 1,706 daily records with no missing data.

2- EXCEL



Data Cleaning in Excel

- **Import CSV File**

via File > Open or Data > Get External Data

- **Convert Date Column to Date Format**

Format Cells > Date

- **Check for Missing Values (non found)**

Use Filters or Conditional Formatting > Blanks

- **Format Numbers**

Comma style for totals

Percentage style for % columns

- **Rename Columns for Clarity**

Data Visualization in Excel

- **Line Chart for Time Trends**

Plot Date vs. Ridership for each mode

- **Pivot Tables for Summaries**

Aggregate by month, mode, or % recovery

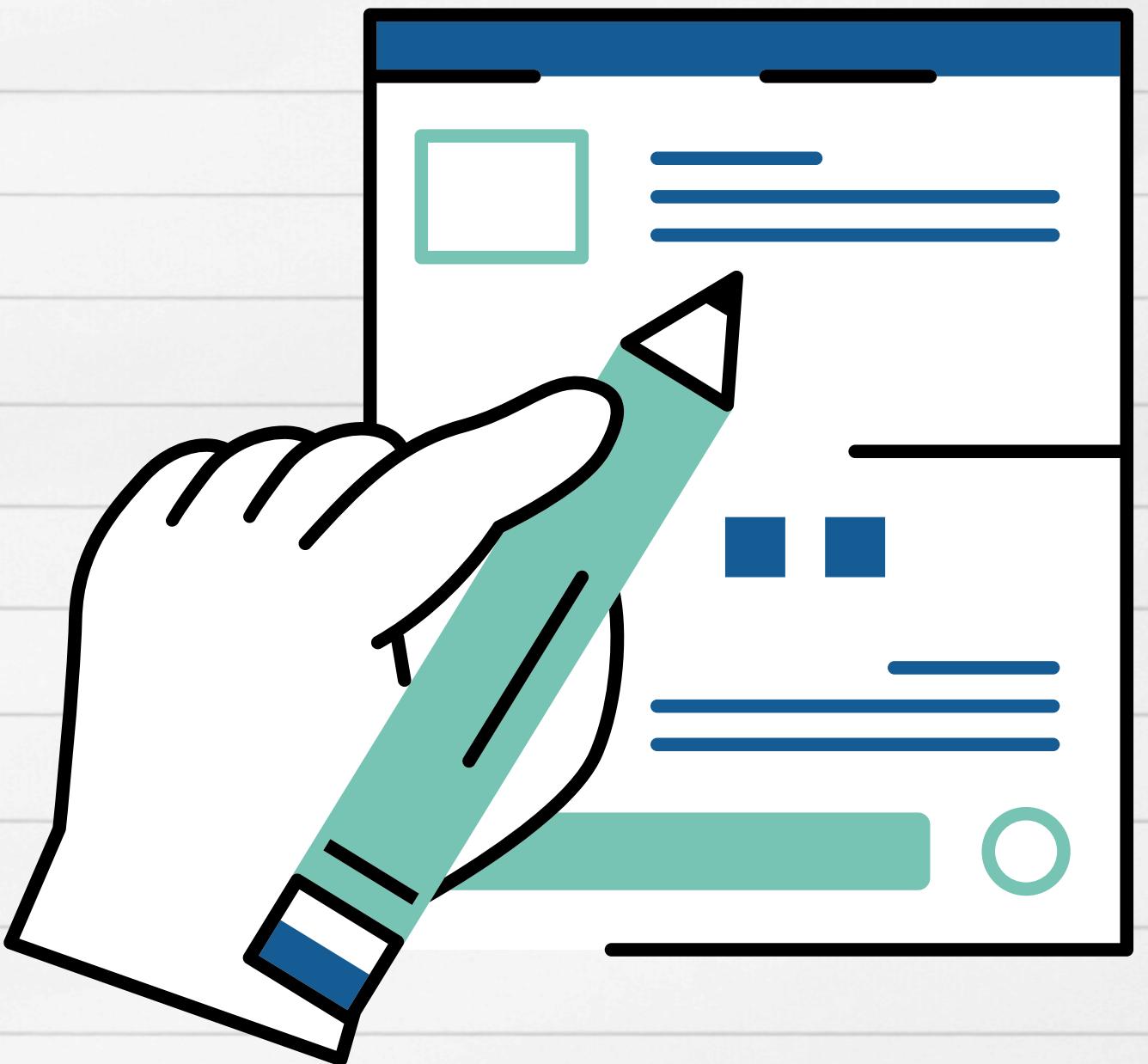
- **Add Slicers or Filters**

Enable interactive exploration

- **Bar/Column Charts**

Compare recovery across transit types

3-SQL

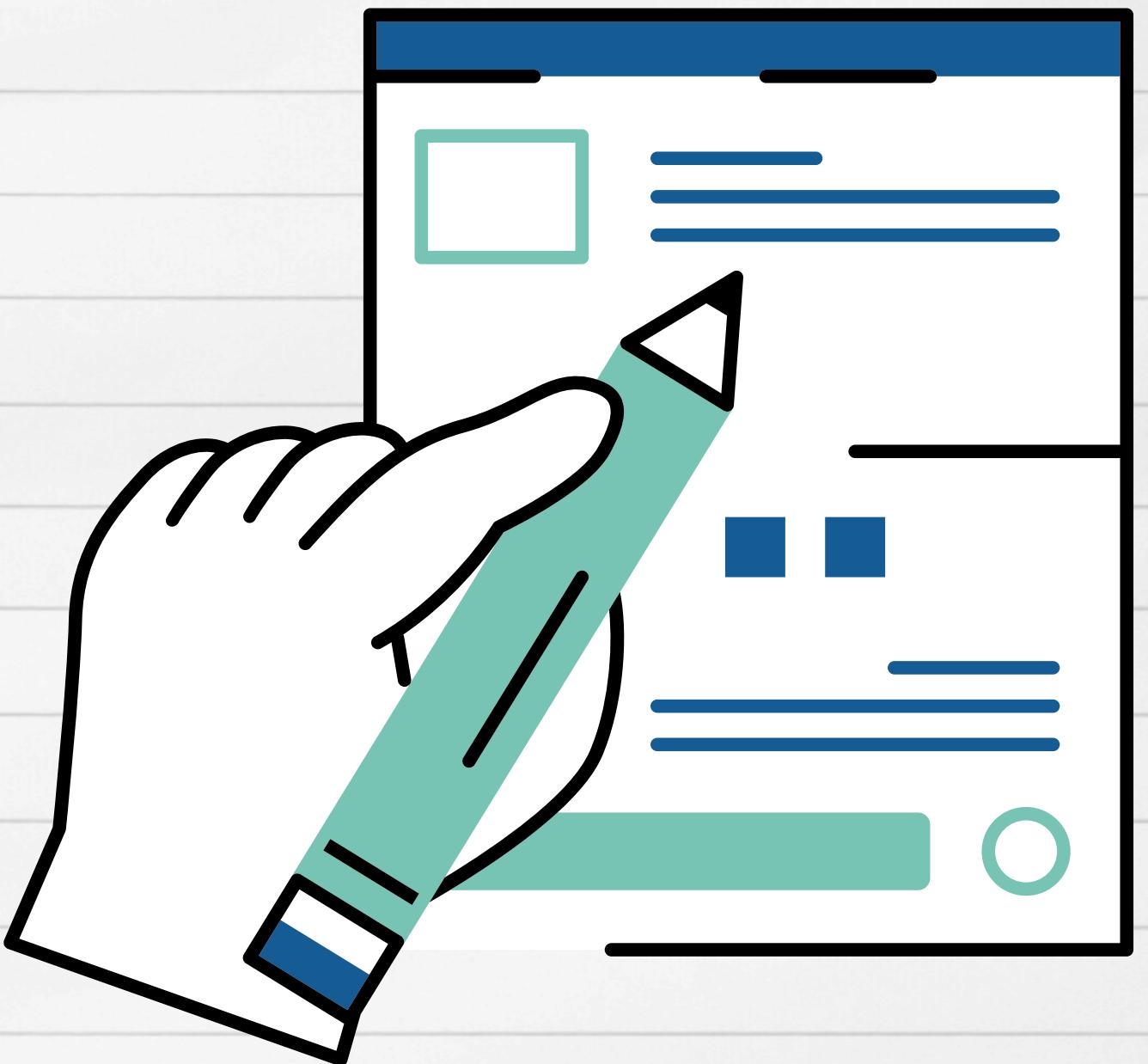


SQL steps

- **Uploaded the csv to** sqliteonline.com
- **ensured that the first row is used as headers having the column names:** founded that it didn't
- **used:** DROP TABLE IF EXISTS your_table_name;
- **Reuploaded the data and inserted first row as header**
- **Viewed the data:** SQL CopyEdit SELECT * FROM your_table_name;
- **Exported the Database:** Go to File → Export DB → Download DB

```
1  SELECT * FROM MTA_Daily_Ridership;
2  DROP TABLE IF EXISTS MTA_Daily_Ridership;
3  |
```

Ч- PYTHON



DATA PREPROCESSING

Goal: Clean and structure the MTA Daily Ridership dataset for analysis.

Steps:

1. Loaded dataset using pandas.
2. Converted date strings to datetime objects.
3.
 - o Created new features:
day_of_week, month, year
is_weekend flag
4. Checked for and handled missing values.
5. Standardized column names.

DATA PREPROCESSING

```
import pandas as pd

# Load and clean dataset
df = pd.read_csv("MTA_Daily_Ridership.csv")
df['date'] = pd.to_datetime(df['date'])
df['day_of_week'] = df['date'].dt.day_name()
df['month'] = df['date'].dt.month_name()
df['year'] = df['date'].dt.year
df['is_weekend'] = df['day_of_week'].isin(['Saturday', 'Sunday'])
```

EXPLORATORY DATA ANALYSIS

Goal: Derive insights by answering key business questions.

INSIGHTS GAINED:

1. Recovery Post-COVID: Bridges and Tunnels rebounded fastest.
2. Weekday vs. Weekend: Weekends showed lower but more stable ridership.
3. Least Drop: Access-A-Ride maintained more consistent usage.

INSIGHTS GAINED:

4. Correlations: Access-A-Ride and Subway showed weak correlation.
5. Seasonality: Ridership dips in winter, peaks in fall/summer.
6. Peak Days/Months: Tuesdays and October had high ridership.

EXPLORATORY DATA ANALYSIS

```
import seaborn as sns
import matplotlib.pyplot as plt

# Example: Correlation heatmap
modes = ['subways', 'buses', 'lIRR', 'metro_north', 'access_a_ride', 'bridges']
plt.figure(figsize=(10, 6))
sns.heatmap(df[modes].corr(), annot=True, cmap="coolwarm")
plt.title("Correlation Between Modes of Transportation")
plt.show()
```

FORECASTING



Goal: Predict ridership for the upcoming month.

Model Used: Linear Regression

Steps:

1. Created total_ridership metric
2. Trained a model on monthly totals
3. Forecasted next month's ridership

Forecast Insight: Estimated steady increase in subway and bus ridership post-pandemic.

FORECASTING



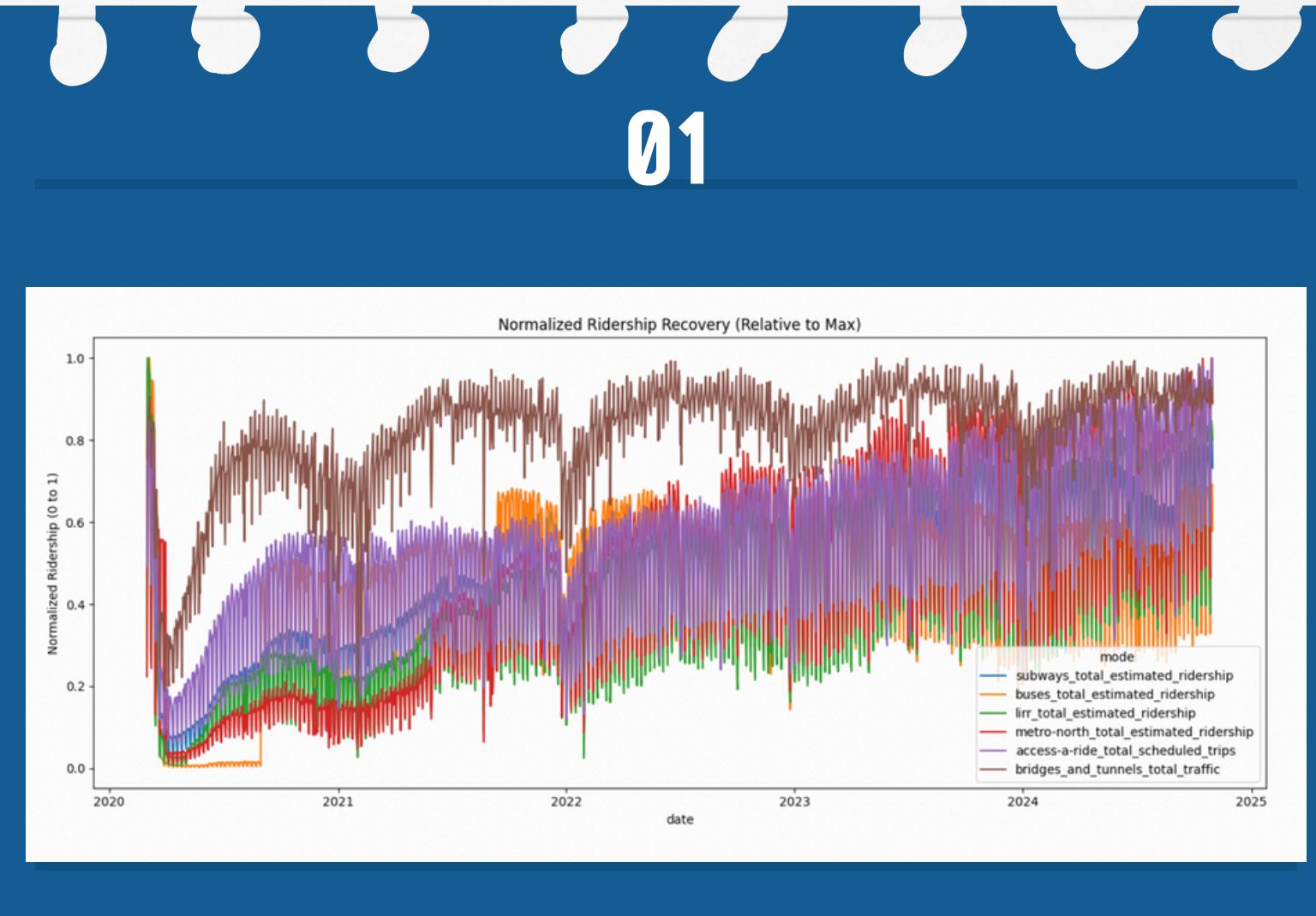
```
from sklearn.linear_model import LinearRegression

monthly_data = df.groupby(['year', 'month']).agg({'subways': 'sum'}).reset_index()
monthly_data['month_num'] = pd.to_datetime(monthly_data['month'], format='%B')
monthly_data['time'] = range(len(monthly_data))

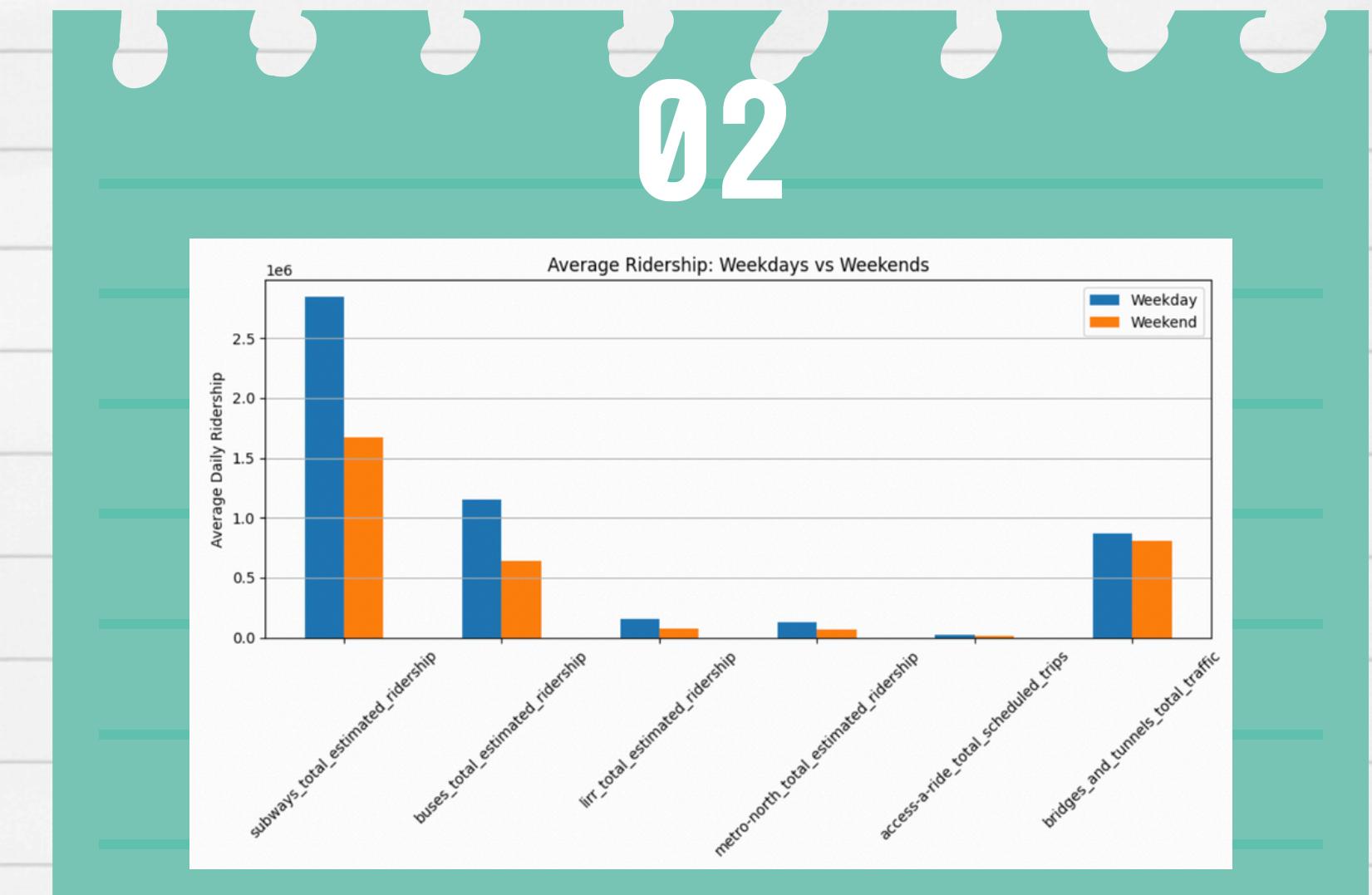
model = LinearRegression()
model.fit(monthly_data[['time']], monthly_data['subways'])
next_month = model.predict([[len(monthly_data)]])
```

SOME OF OUR VISUALIZATIONS

01

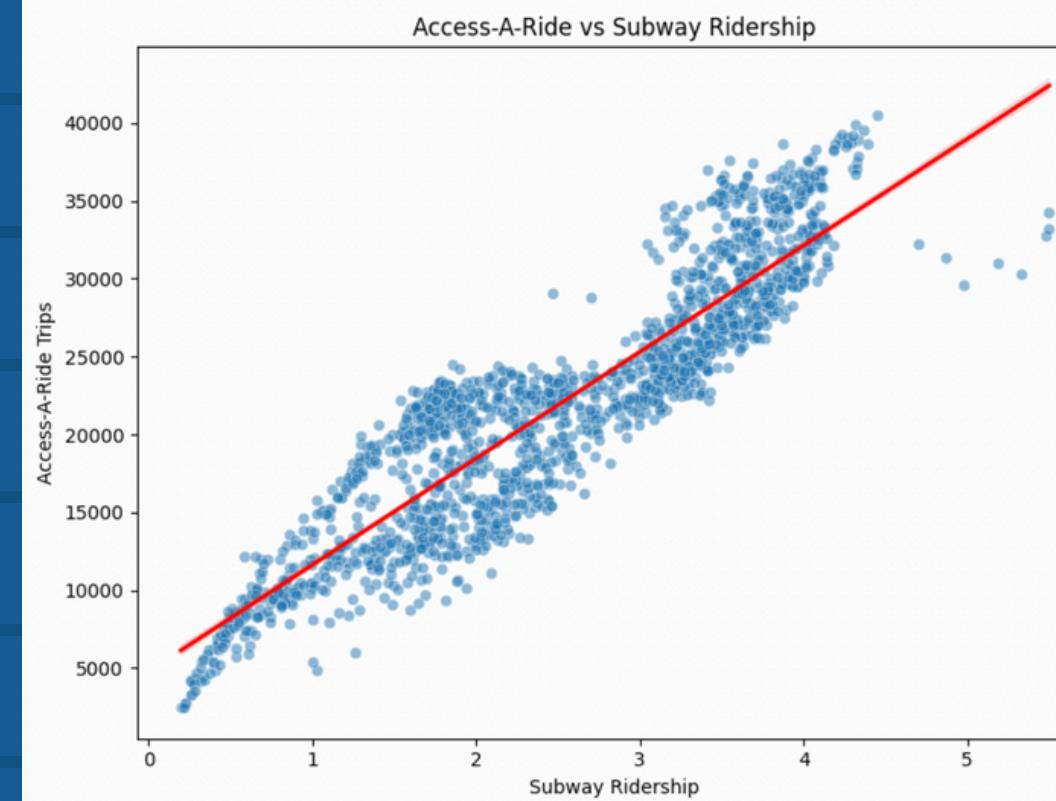


02

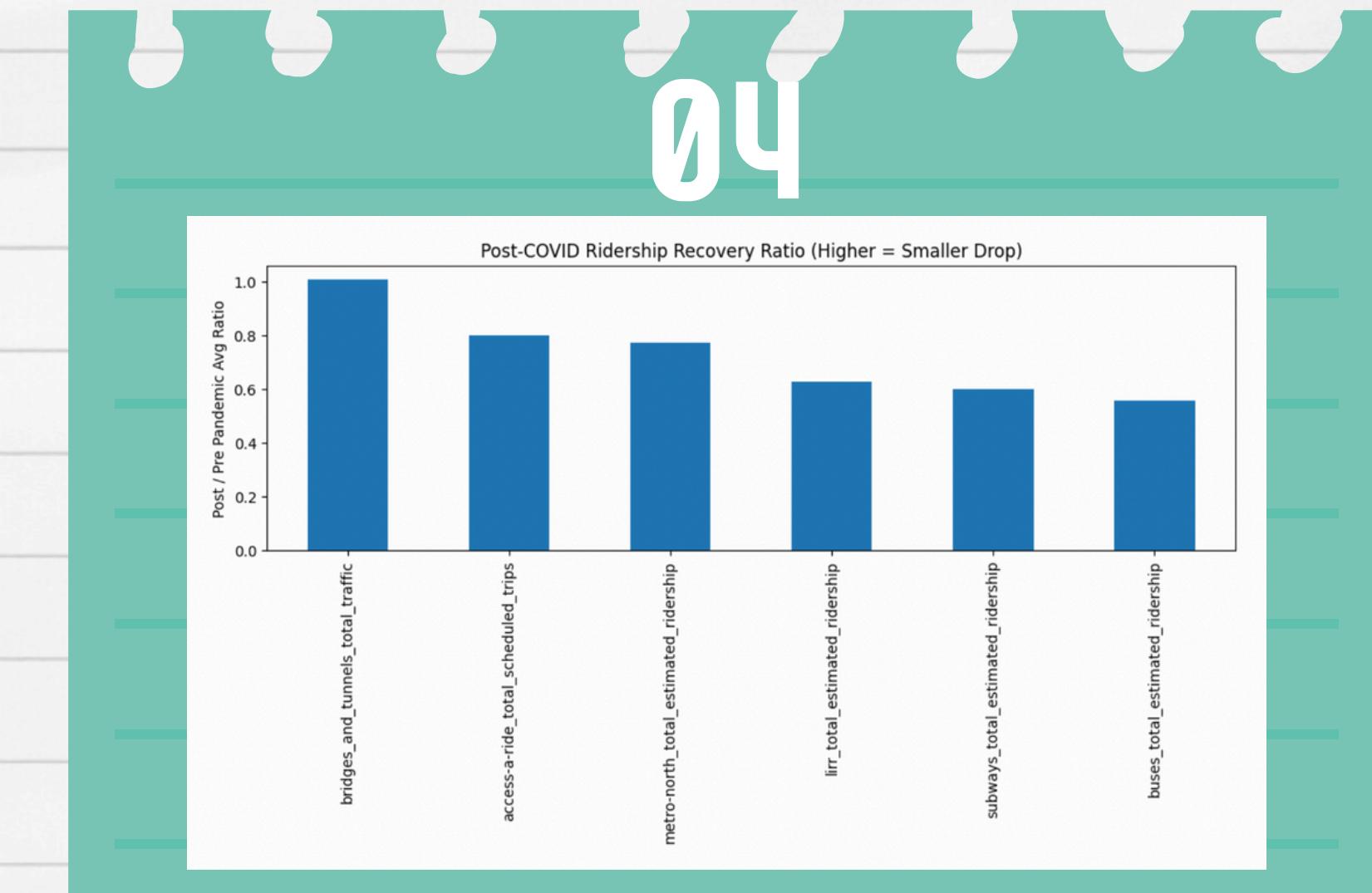


SOME OF OUR VISUALIZATIONS

03

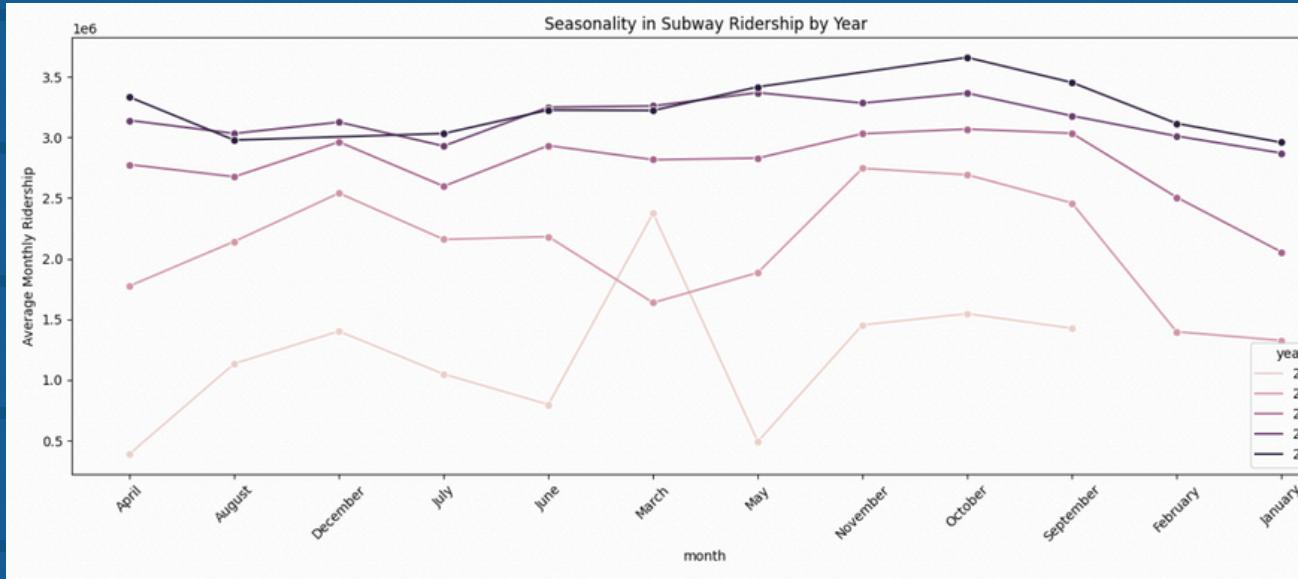


04



SOME OF OUR VISUALIZATIONS

05



06



SUBWAY RIDERSHIP PREDICTION USING GRU MODEL FORECASTING DAILY MTA PASSENGER COUNTS

DATASET DESCRIPTION

Data Overview:

- The dataset contains daily subway ridership counts in New York City
- It covers the period from 2020 to 2023



Data Preprocessing:

- Handled missing values (`df.dropna()`)
- Converted time columns into usable features(`pd.to_datetime()`)
- Applied MinMaxScaler to normalize the data (Used: `MinMaxScaler()` from `sklearn.preprocessing`)

Data Splitting: the data was split into: (Used: `train_test_split()`)

- Training set
- Validation set
- Test set

PROBLEM STATEMENT

- Predicting the number of subway riders each day is a surprisingly smart challenge. In this project, we aim to harness the power of deep learning to forecast daily subway ridership using historical data. Accurate predictions can help the city optimize subway services, ease overcrowding, and create a smoother, more efficient commute for millions of New Yorkers. It's not just about numbers — it's about making everyday travel smarter and more reliable.



MODEL SELECTION

Gated Recurrent Unit (GRU)

we selected GRU as deep learning model for time series forecasting. It's type of RNN(recurrent neural network).

EFFICIENCY

- Has only two gates (update and reset). This design results in fewer trainable parameters, which speeds up the training process and reduces memory.

PERFORMANCE

- GRU delivers strong performance in time series forecasting tasks.

MODEL ARCHITECTURE ARCHITECTURE OF GRU

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import GRU, Dense, Dropout, BatchNormalization
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import EarlyStopping, ReduceLROnPlateau

model = Sequential()

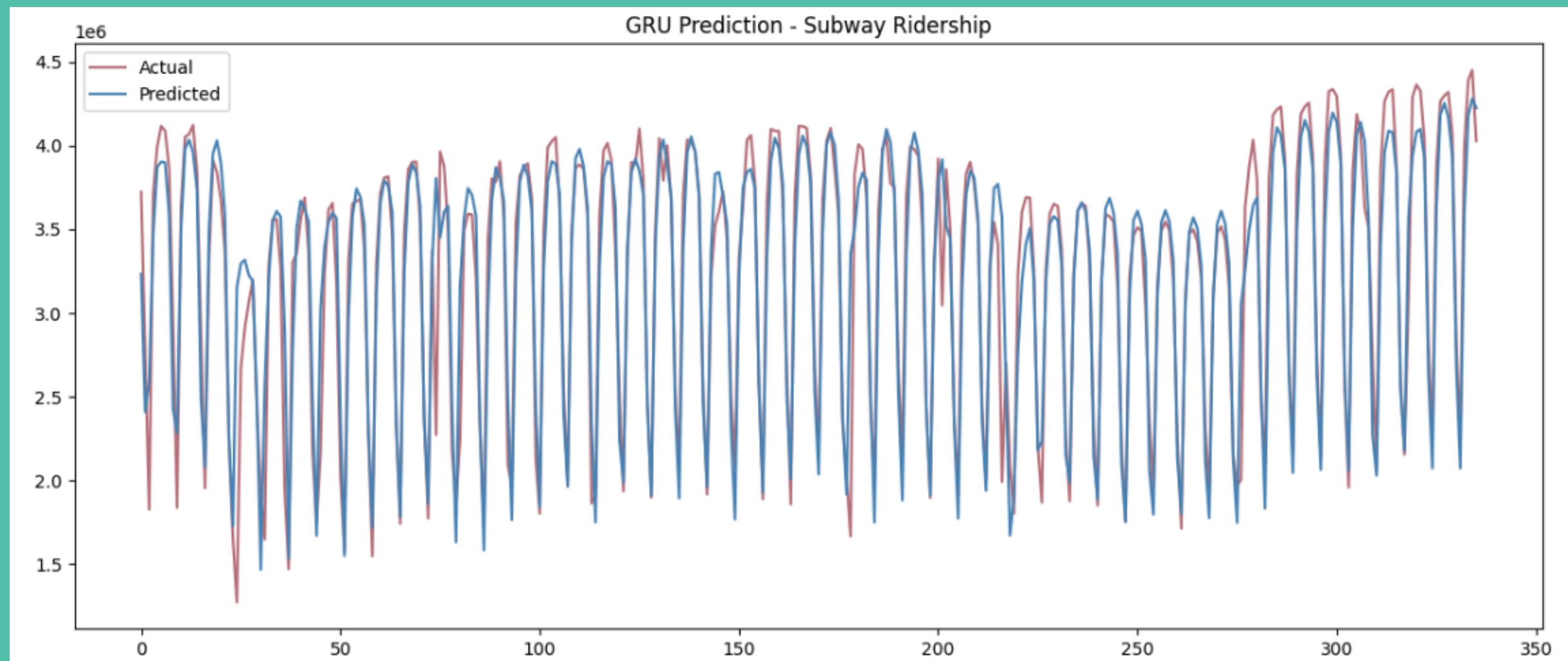
model.add(GRU(128, return_sequences=True, input_shape=(seq_length, 1)))
model.add(Dropout(0.2))
model.add(BatchNormalization())

model.add(GRU(64, return_sequences=False))
model.add(Dropout(0.2))
model.add(BatchNormalization())

model.add(Dense(1))

model.compile(optimizer='adam', loss='mse')
```

RESULTS & EVALUATION



Acutual

PINK LINE

Predicted

BLUE LINE

5-TABLEAU



PROJECT SUMMARY (IN 5 STEPS):

Project Summary (in 5 steps):

- 1. Imported Data:** NYC daily ridership CSV file (Subways, Buses, etc.).
- 2. Cleaned Data:** Selected key columns like date, ridership, and % of pre-pandemic use.
- 3. Trend Analysis:** Created line charts to show how ridership changed over time.
- 4. Impact Comparison:** Used bar charts to compare how each transport mode was affected by COVID.
- 5. Built Dashboard:** Combined everything in an interactive Tableau dashboard with filters.

EXPLORATORY DATA ANALYSIS

Goal: Derive insights by answering key business questions.

INSIGHTS GAINED:

1. Subways were the most affected
Ridership dropped significantly and took longer to recover compared to other modes.
2. Buses recovered faster
Bus usage returned to a higher % of pre-pandemic levels quicker than subways.

INSIGHTS GAINED:

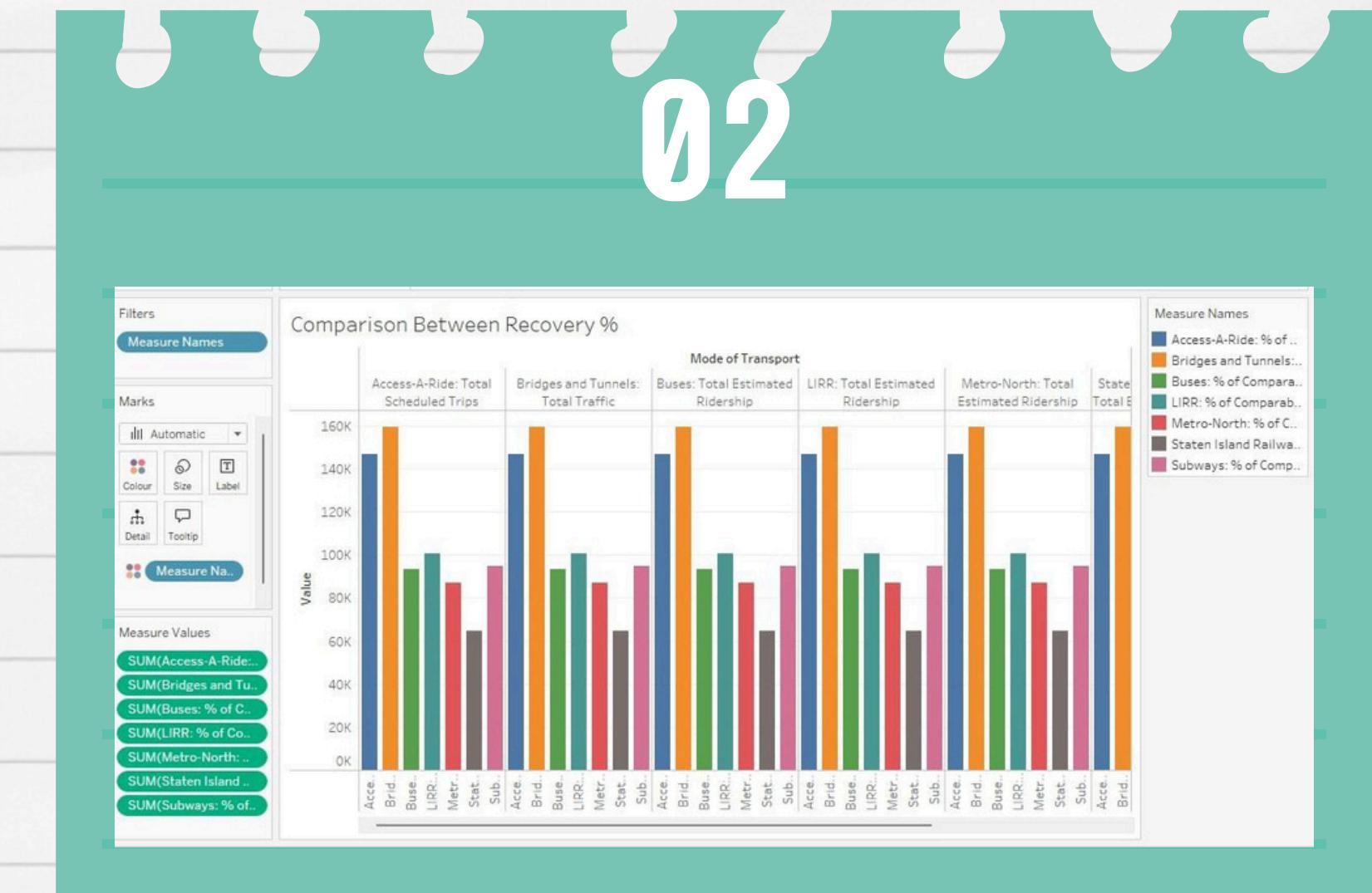
3. LIRR and Metro-North stayed low
These commuter rail systems showed consistently low recovery, indicating continued remote work trends.
4. Pandemic impact peaked around early 2020
Across all modes, the steepest drop in ridership was between March–April 2020.
5. Ridership trends reflect public confidence
Recovery patterns align with reopening phases and public health developments.

SOME OF OUR VISUALIZATIONS

01

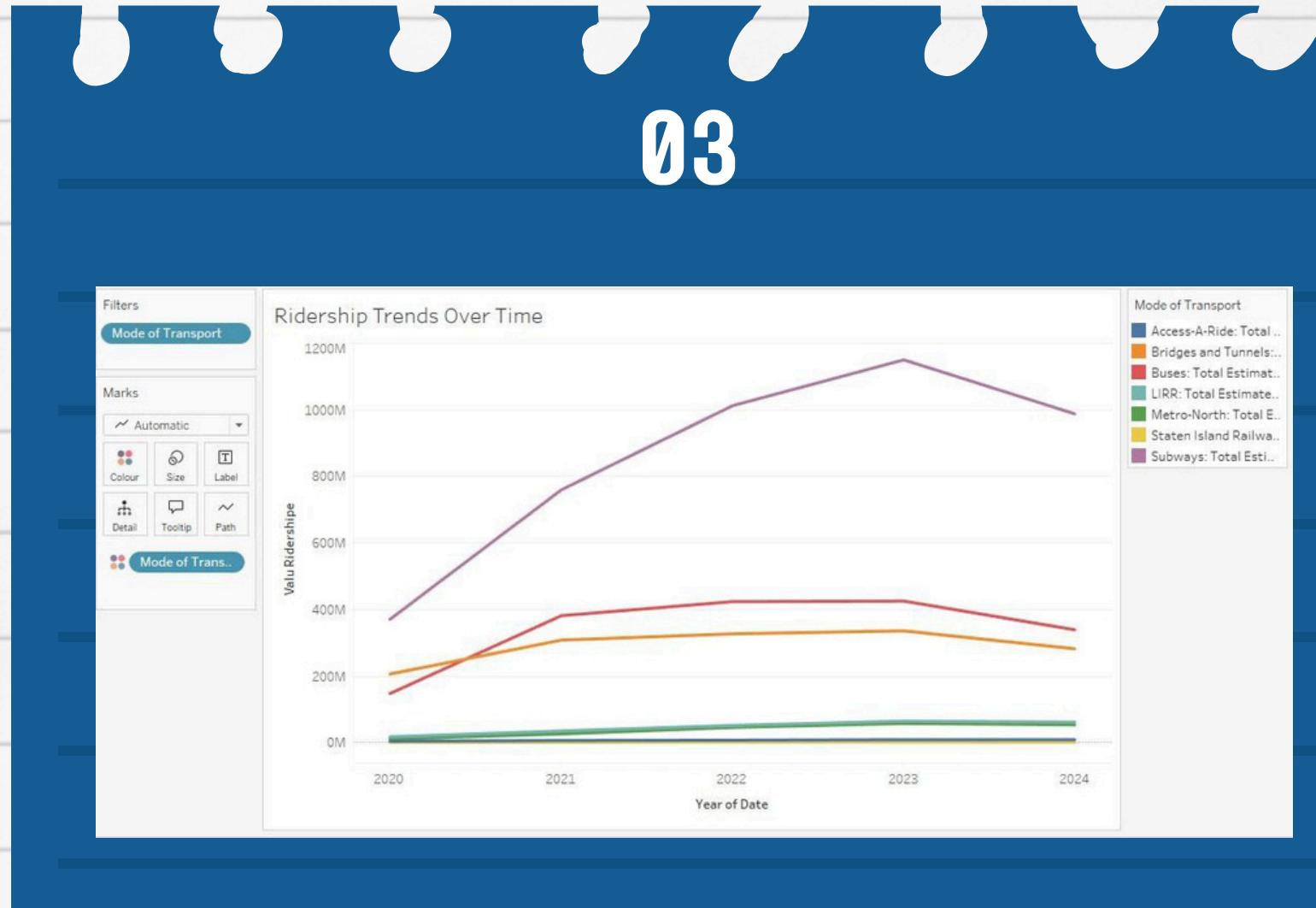


02

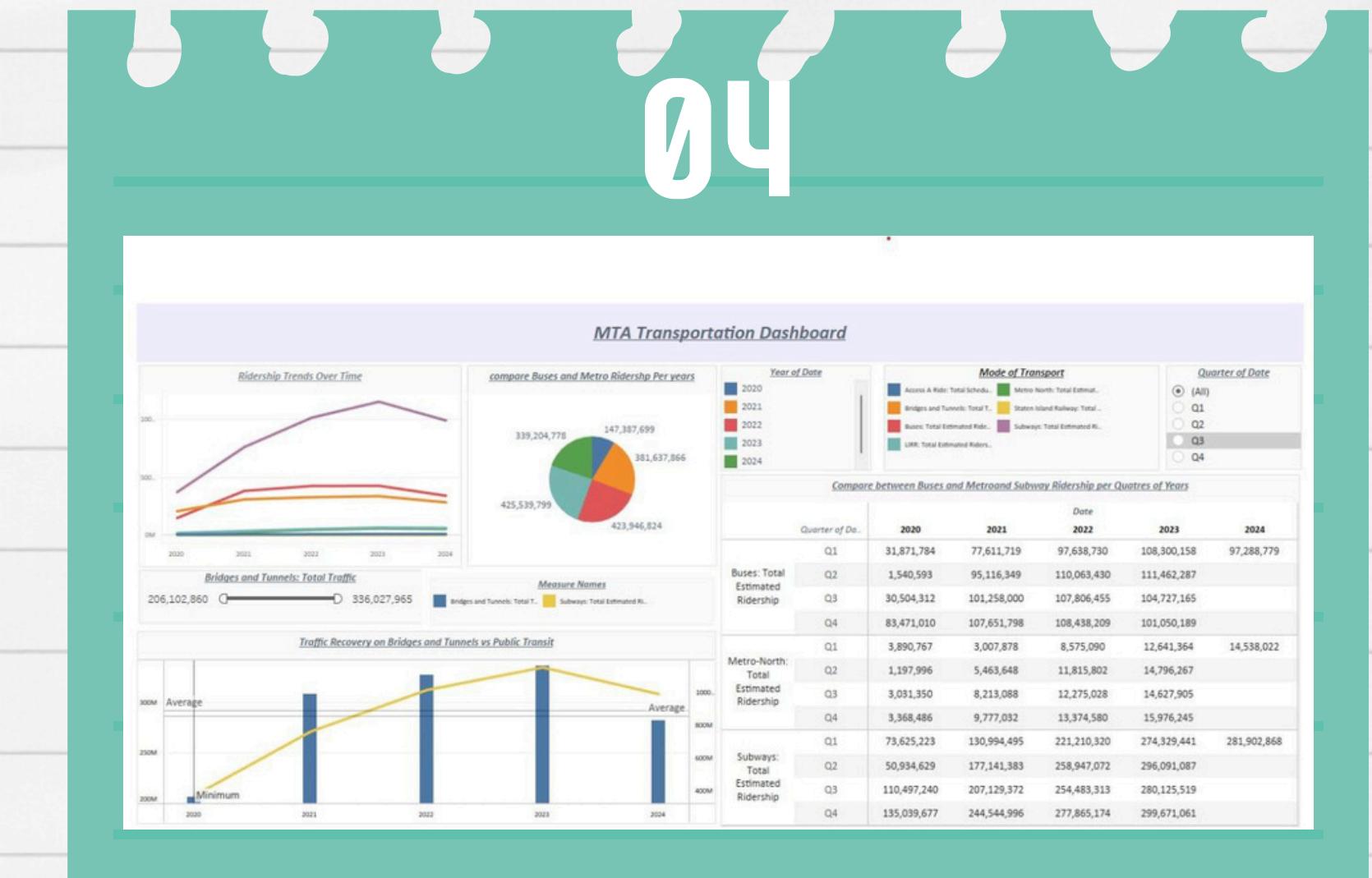


SOME OF OUR VISUALIZATIONS

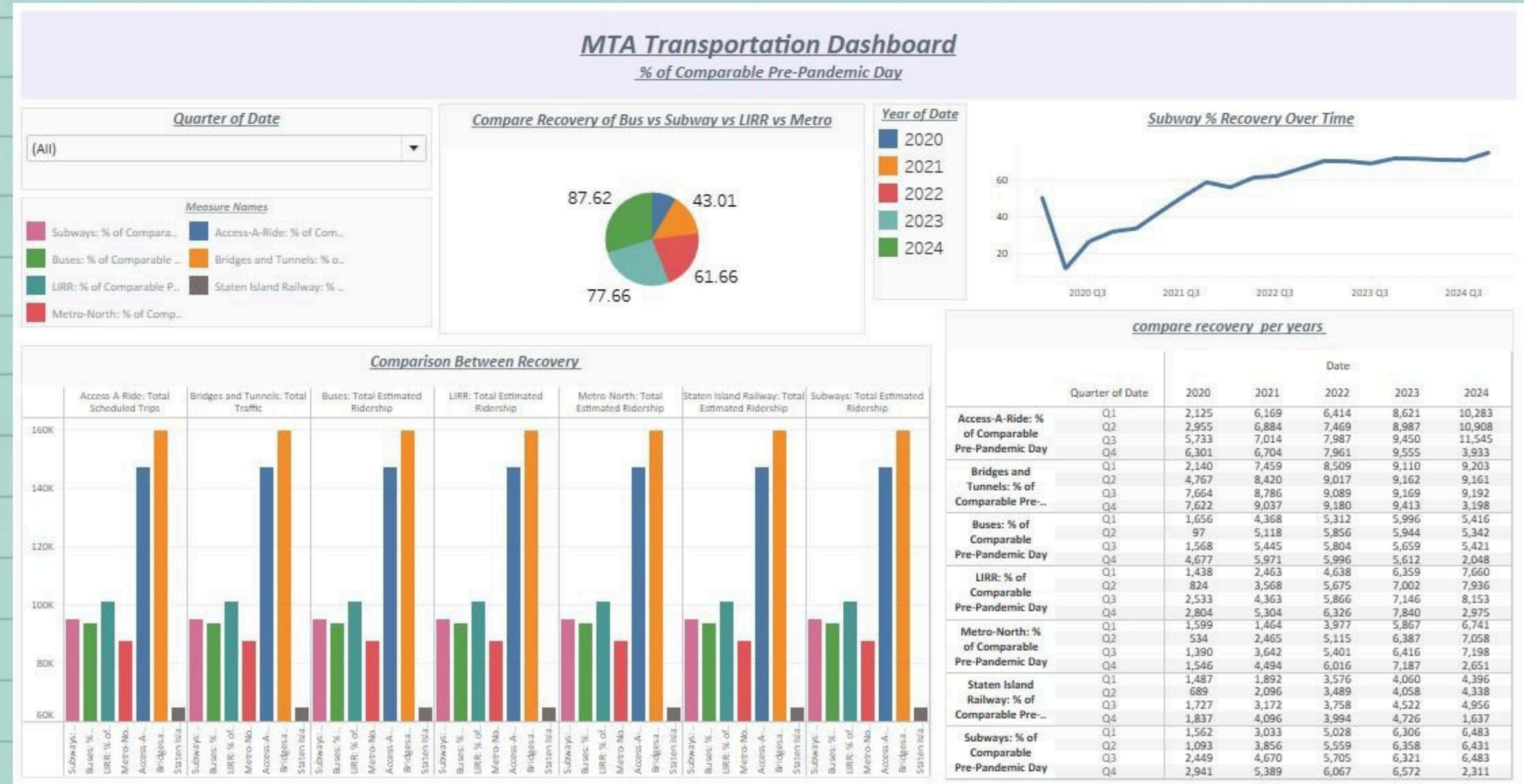
03



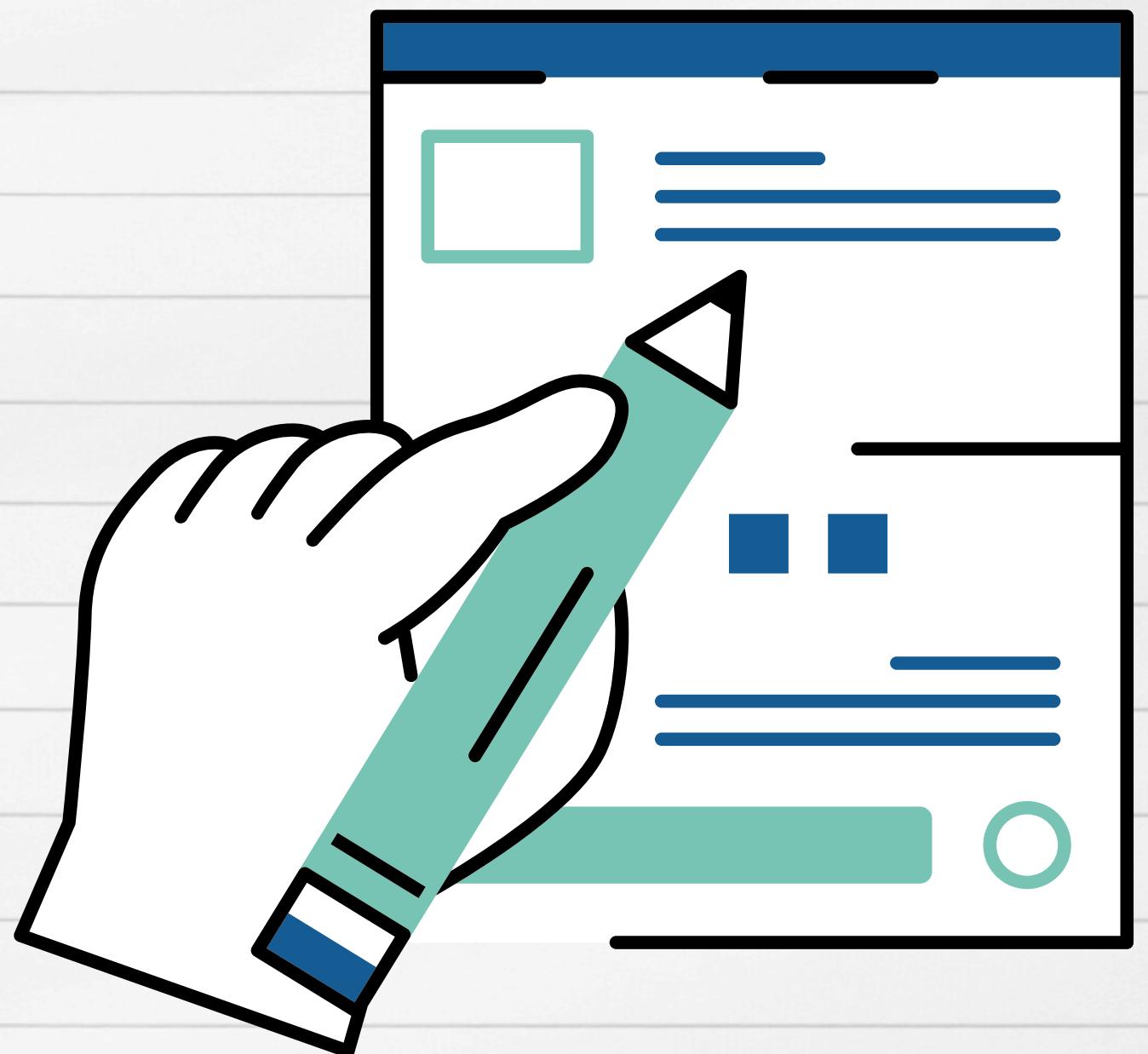
04



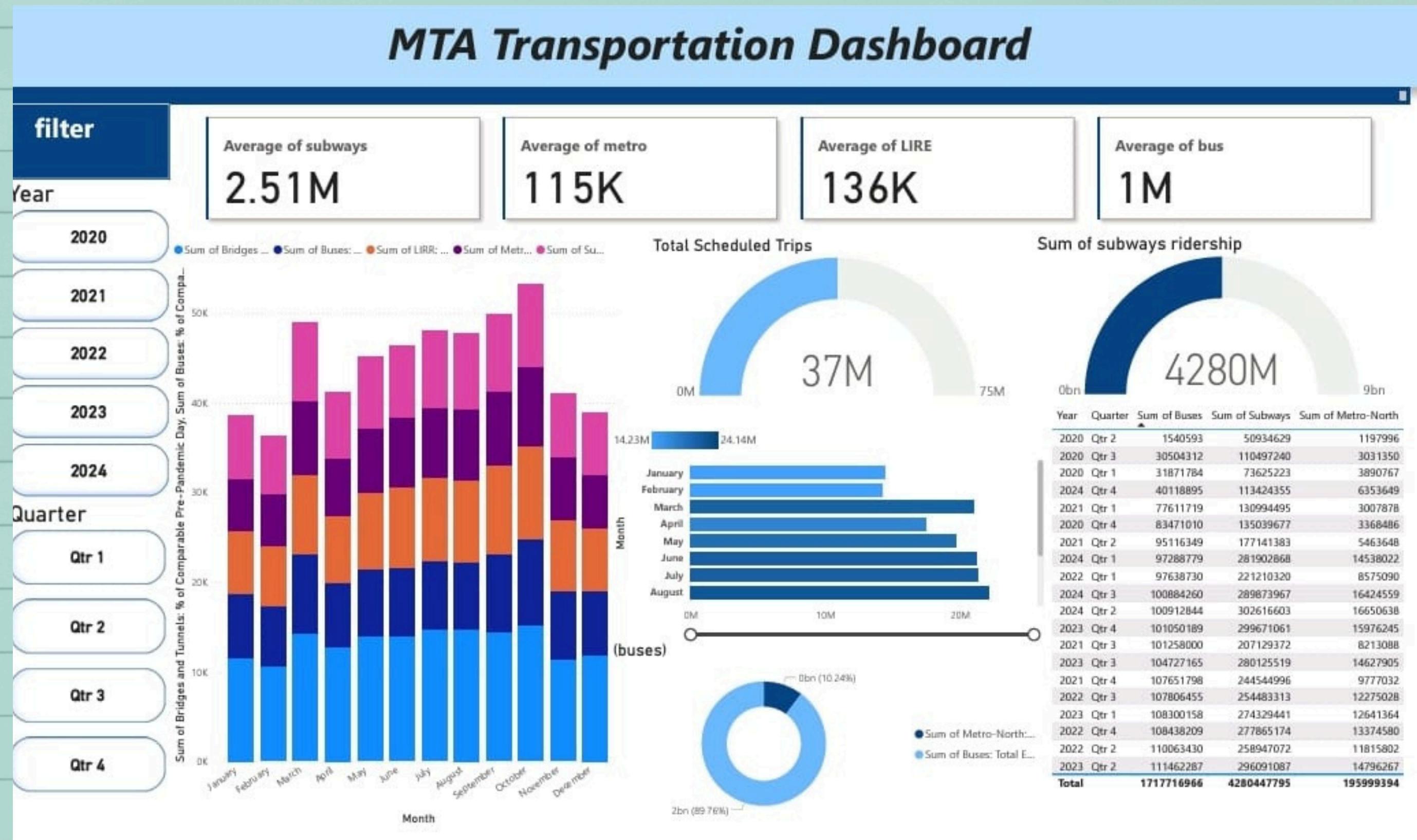
FINAL VISUALIZATIONS



6-POWER BI



FINAL VISUALIZATIONS



**THANK
YOU VERY
MUCH!**