

---

# Mapping the Pulse of 2024 Traffic

---

**Manasvi Patwa**

Goergen Institute for Data Science  
University of Rochester  
Rochester, NY  
mpatwa@ur.rochester.edu

**Anjaly George**

Goergen Institute for Data Science  
University of Rochester  
Rochester, NY  
ageor16@ur.rochester.edu

## 1 Introduction

In this project, we seek to address the complexities of traffic forecasting, a critical aspect of urban traffic systems. Traditional models fall short in capturing the nonlinear nature and dependencies inherent in traffic flows, impeding effective predictions. Our research uses a Spatio-Temporal Graph Convolutional Network (STGCN) approach, which innovatively applies convolutional structures to graph-structured data, extracting spatio-temporal features simultaneously. By focusing on the interconnectedness and chronological patterns of traffic data for the year 2024, the STGCN model aims to provide unprecedented accuracy in mid- to long-term traffic forecasting.

## 2 Problem Statement

The project aims to leverage the Spatio-Temporal Graph Convolutional Network (STGCN) to analyze and predict traffic patterns using granular, real-time data from California's PeMSD7 dataset. Traditional forecasting models struggle to capture the nonlinear and complex nature of traffic flows due to their inability to account for the strong spatio-temporal dependencies inherent in traffic data. By combining graph convolutional layers for spatial pattern recognition with gated temporal convolutional layers for capturing traffic flow dynamics over time, STGCN provides accurate forecasts across different time horizons. This allows for better-informed real-time traffic management and long-term urban planning, ensuring authorities can anticipate and mitigate congestion effectively.

## 3 Literature Survey

In a study detailed in [1], researchers introduce a Spatio-Temporal Graph Convolutional Network (STGCN) that significantly enhances traffic forecasting capabilities beyond those provided by traditional ARIMA and machine learning techniques. The study effectively demonstrates the STGCN's ability to manage complex spatial-temporal interactions within traffic data. Utilizing the BJER4 and PeMSD7 datasets, which include a variety of traffic metrics. The STGCN framework integrates graph convolutional and gated convolutional layers, effectively harnessing the structured nature of graph-based time series data. This approach results in markedly improved forecasting accuracy across various testing scenarios.

In [2] by Li et al. (2018), introduce a novel framework that combines graph convolution with recurrent neural networks to model traffic as a diffusion process on a network graph. This model, named Diffusion Convolutional Recurrent Neural Network (DCRNN), is designed to capture both spatial and temporal correlations in traffic data effectively. It employs diffusion convolution to address non-Euclidean spatial relationships and recurrent neural networks to model temporal dependencies. The paper demonstrates the DCRNN's superior forecasting ability on real-world traffic datasets, outperforming several baseline and state-of-the-art models. This approach provides significant improvements in traffic forecasting, offering a data-driven solution with practical applications in urban traffic management and planning systems.

Wu and Tan’s 2016 paper [3] proposes an advanced model for predicting short-term traffic flow. They develop a hybrid framework that integrates Convolutional Neural Networks (CNNs) for extracting spatial features from traffic data across road networks, with Recurrent Neural Networks (RNNs) for analyzing temporal patterns in traffic flow over time. This dual approach effectively harnesses the complex interplay between spatial and temporal factors in traffic data, aiming to deliver more accurate short-term forecasts. The effectiveness of their model is underscored by its performance in experiments, which suggests significant improvements over traditional models that consider spatial and temporal components separately.

These studies are highly relevant to this project because they highlight the importance of combining spatial and temporal features to accurately model traffic patterns. Each study explores unique ways to capture complex traffic dependencies across road networks, ultimately achieving superior forecasting accuracy compared to traditional models. By understanding these frameworks, we can refine the design of our own STGCN model and tailor it to the specific challenges of the PeMSD7 dataset, ensuring more effective traffic management and planning solutions.

## **4 Dataset**

In our traffic analysis project, we utilize the PeMSD7 dataset from California’s Performance Measurement System (PeMS) [4], focusing on data from March and February 2024. This dataset is collected at 5-minute intervals, capturing detailed traffic metrics such as speed, flow, and density, and providing 288 data points per node per day from sensor stations across California’s highways. We model the traffic networks as graphs where nodes represent sensor stations or road segments, and edges depict connectivity. This granular, graph-based structure facilitates the application of spatial-temporal graph convolutional networks, enabling comprehensive analysis and predictive modeling of traffic patterns under various conditions.

## **5 Methodology**

### **5.1 Data Preprocessing**

In our traffic data preprocessing workflow, we carefully filtered out the relevant columns from a larger dataset that included numerous features. The columns we selected were 'Timestamp', 'Station', 'Freeway No', 'Direction of Travel', 'Samples', 'Total Flow', 'Avg Occupancy', and 'Avg Speed'. This selection process ensured that we focused on the most relevant data attributes necessary for our analysis.

To address missing values within our dataset, we applied a linear interpolation method, which efficiently estimated missing data points. This step was crucial in maintaining a continuous and complete dataset, enhancing its reliability for detailed traffic pattern analysis. We also applied Z-score normalization on the data to scale the numerical features.

We further enriched our primary traffic data by integrating auxiliary metadata from a secondary dataset. This metadata included crucial geographical details such as the latitude and longitude of each station. We merged this geographical metadata with our main dataset using a right join on the 'Station' column, significantly enhancing the dataset’s value for spatial analyses. This enriched dataset now provided a comprehensive foundation for our subsequent traffic flow modeling and analysis efforts.

### **5.2 Why we chose to work with undirected graphs?**

In the analysis of traffic data using graph-based approaches, the utilization of undirected graphs simplifies the modeling of traffic networks considerably. This simplification is particularly advantageous when representing bidirectional traffic flows on roads, as it allows each road segment to be modeled without regard to the directionality of traffic, thus reducing model complexity and computational overhead. Furthermore, this approach supports the generalization of traffic patterns, capturing symmetrical traffic conditions effectively across both directions of a road, which is beneficial for aggregate data analysis. From an algorithmic perspective, the use of undirected graphs enhances the efficiency of graph-based algorithms, particularly those used in machine learning frameworks

like Graph Convolutional Networks (GCNs), where simplifying the computation of the Laplacian matrix is crucial. This adaptation facilitates a more straightforward application of spectral graph theory, which is instrumental in leveraging spatial features of traffic data for predictive modeling. This methodological choice, therefore, not only streamlines the analytical processes but also ensures compatibility with advanced analytical tools, making it a practical approach for large-scale traffic studies.

### 5.3 Constructing graphs

In constructing the traffic network graph for our analysis, each node within the graph represents a sensor station, capturing traffic data from various locations across the highway system. The edges between these nodes signify the distance between stations, providing a measure of connectivity that influences traffic flow between segments. To quantify the relationship between nodes, we compute the adjacency matrix  $\mathbf{W}$ , where each element  $w_{ij}$  is calculated using the formula:

$$w_{ij} = \begin{cases} \exp\left(-\frac{d_{ij}^2}{\sigma^2}\right), & i \neq j \text{ and } \exp\left(-\frac{d_{ij}^2}{\sigma^2}\right) \geq \epsilon \\ 0, & \text{otherwise} \end{cases}$$

Here,  $d_{ij}$  represents the Euclidean distance between stations  $i$  and  $j$ ,  $\sigma$  is a predefined threshold that controls the decay rate of the weight with respect to distance, and  $\epsilon$  is a small threshold value to maintain sparsity in the adjacency matrix. This formulation ensures that closer nodes have a stronger connection, reflecting higher interaction potential in terms of traffic flow, while distant nodes have minimal or no influence on each other.

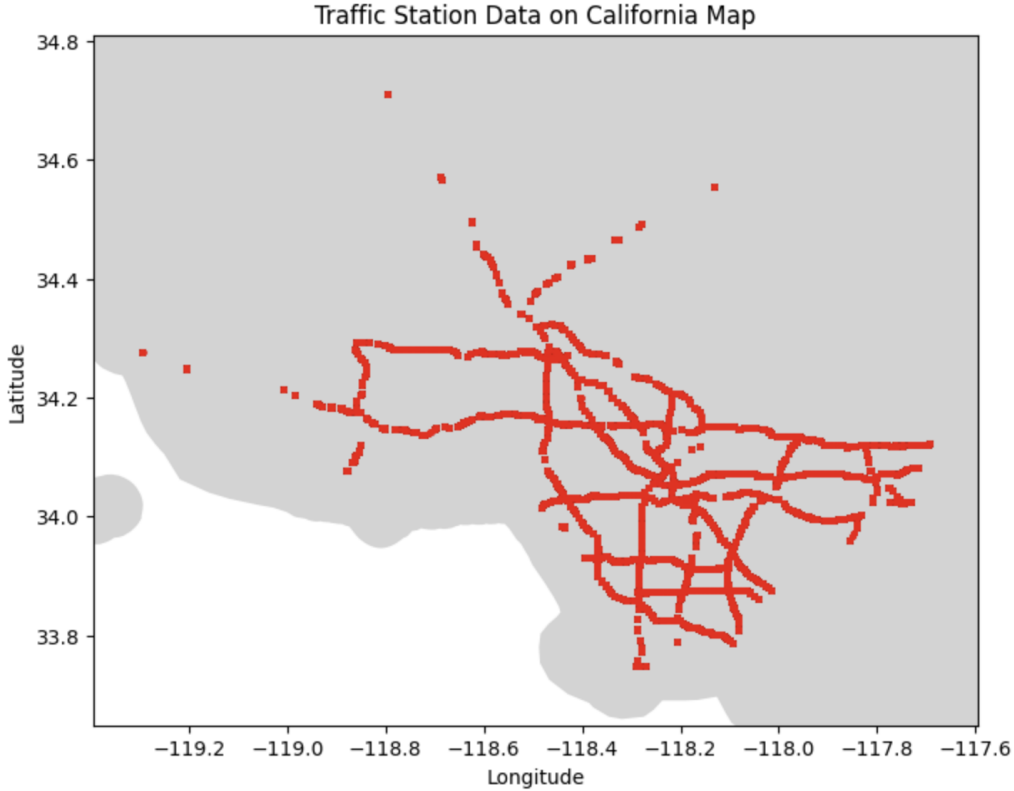


Figure 1: Traffic Stations on California Map

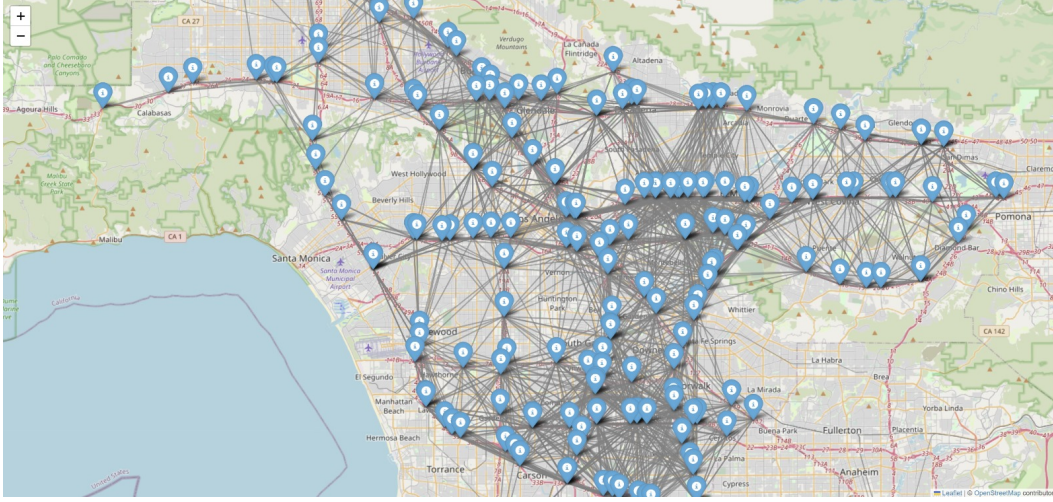


Figure 2: California Traffic Station Network Graph

This graph in Figure 2 visualizes the intricate network of 190 traffic monitoring stations in California, derived from the PeMSD dataset. Each station is represented as a node, while the connections (edges) between nodes indicate spatial relationships. The weights of the edges were calculated using a Gaussian kernel formula to represent geographic distances between stations accurately. The map provides a clear picture of the connectivity between these traffic stations, showing the spatial structure and geographical distribution of California's highway network.

The distance between each station pair was calculated using the Haversine formula. A Gaussian kernel was then applied to generate the weights between stations, filtering connections based on a predefined distance threshold. An adjacency matrix was constructed to represent the connectivity, which was then used to build the graph. An adjacency matrix was constructed to represent the connectivity, which was then used to build the graph.

## 5.4 Exploratory Data Analysis

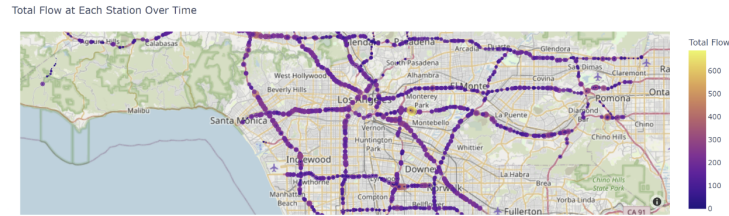


Figure 3: Total Flow at Each Station at 12AM

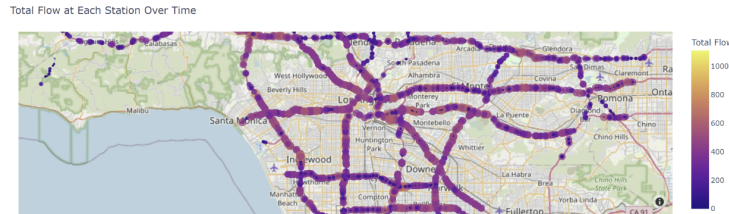


Figure 4: Total Flow at Each Station at 6AM

Figures 3 and 4 depict the total traffic flow at California's stations, revealing a clear daily pattern. Traffic is lowest at midnight, as expected, since most people are at home. It begins to rise from 6 AM, peaking during the hours when most people are commuting to work or school, and then starts to taper off after 9 PM, coinciding with the end of typical business hours and lessening social activities.

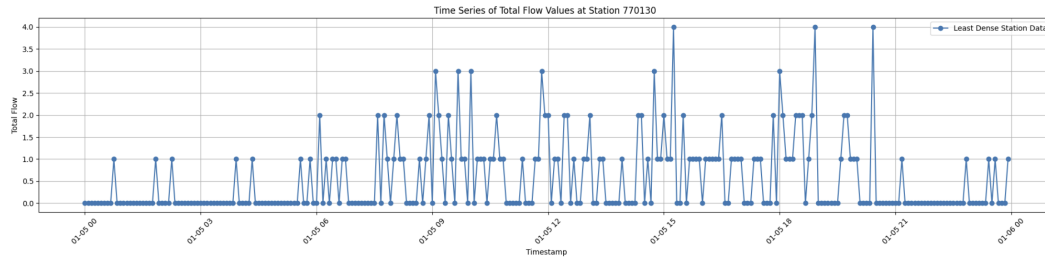


Figure 5: Time Series of Total Flow Values at Station 718308 which is the station having highest traffic flow

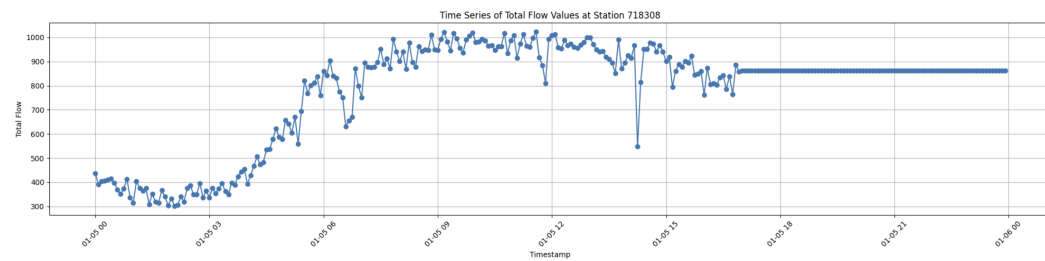


Figure 6: Time Series of Total Flow Values at Station 770130 which is the station having lowest traffic flow

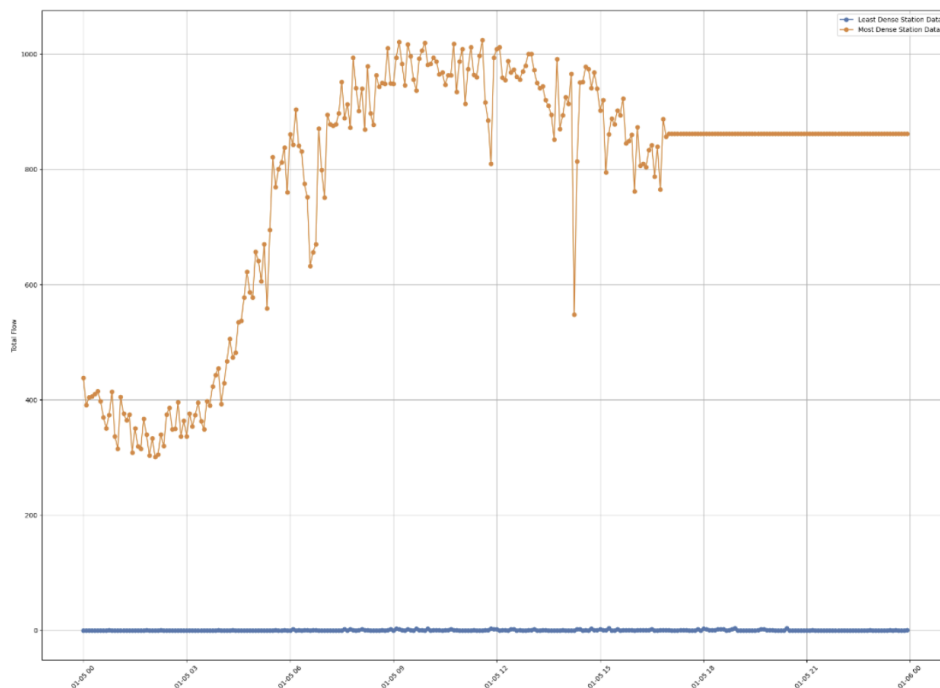


Figure 7: Time Series of Total Flow Values at Most Dense and Least Dense Stations

In Figures 5, 6, and 7, we observe the time series graphs representing the traffic flow of the least and most dense stations, alongside their overlapped data, respectively. Figure 4 illustrates a station with minimal traffic flow, suggesting it could either be an area with intrinsically low traffic or a data collection point with potential recording issues. Figure 5 exhibits the contrasting pattern of the busiest station, indicating consistent traffic throughout the day with identifiable peak hours. The overlay presented in Figure 6 offers a comparative perspective, highlighting the vast differences in traffic density between the two stations.

These visualizations, particularly the contrast between traffic densities, will be pivotal in understanding the range of traffic patterns across the network. When analyzed over a two-month period, such insights can pinpoint key areas for traffic management interventions, inform infrastructure development, and facilitate targeted congestion mitigation strategies. They offer a granular understanding of station-specific behaviors, which is essential for tailoring traffic solutions and predicting potential bottlenecks.

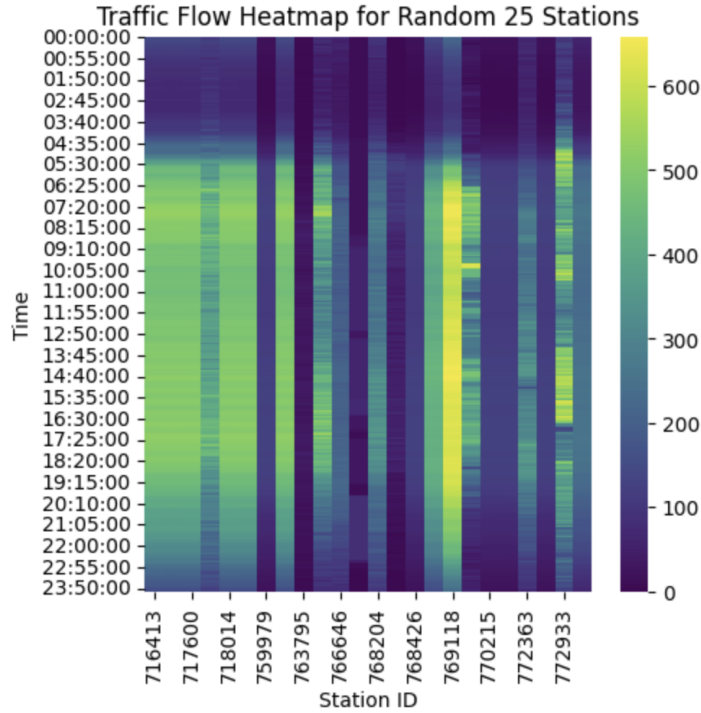


Figure 8: Traffic Flow Heatmap for randomly selected 25 stations

This heatmap in Figure 9 provides a succinct visual summary of traffic flows for 25 random stations over a 24-hour period, with varying colors indicating traffic volume at 5-minute intervals. Lighter shades during certain periods hint at potential rush hours, while continuous dark or light patterns across stations may pinpoint key traffic nodes or underutilized areas, respectively. The visual disparity among stations and times offers insights into daily and location-specific traffic trends, which are essential for identifying congestion points and anomalies that could affect traffic dynamics. Analyzing such patterns over a two-month period will be critical for detecting consistent traffic behaviors, informing infrastructure planning, and optimizing traffic management strategies to improve overall commuter experience.

## 5.5 STGCN Model Architecture

The STGCN model as seen in Figure 8 is constructed to process data that has both spatial and temporal components, making it particularly useful for tasks such as traffic forecasting, where the time and location are key factors.

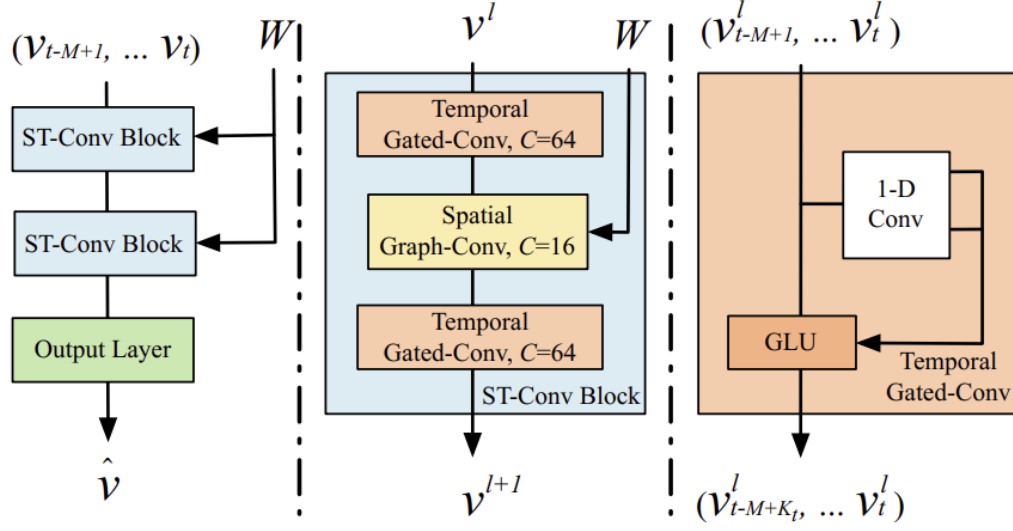


Figure 9: STGCN Model Architecture

### Structure and Components

#### 1. Input Layer:

- The input to the model consists of time-series data from multiple sensors or nodes, represented as  $V_{t-M+1}$  to  $V_t$ , where  $M$  is the number of time steps considered.

#### 2. Spatio-Temporal Convolutional Blocks (ST-Conv Blocks):

- Each ST-Conv Block is structured in a “sandwich” configuration, comprising two temporal gated convolutional layers with a spatial graph convolutional layer positioned between them.
- Temporal gated convolutions handle the time-series aspect, effectively learning patterns across different time steps.
- The spatial graph convolution targets the spatial dimension, learning the inter-node (or sensor) relationships and dependencies.
- This combination allows the network to simultaneously learn both the temporal progression and spatial distribution of the data.

#### 3. Residual Connections and Bottleneck Strategy:

- Inside each ST-Conv Block, residual connections help in avoiding the vanishing gradient problem by allowing gradients to flow through alternate pathways.
- The bottleneck strategy reduces the dimensionality internally within the block, enhancing computational efficiency and focusing the network’s learning capacity.

#### 4. Output Layer:

- Following the ST-Conv blocks, the fully-connected output layer produces the final prediction, utilizing the features extracted by the preceding blocks.

#### 5. Workflow:

- The model processes data through successive ST-Conv blocks, each refining the feature representation.
- The output layer then maps these refined features to the desired output, such as predicting future values in a time series.

The STGCN model leverages its unique design to effectively process and predict patterns in data that varies across both time and space. Its architecture is tailored to integrate and analyze complex interactions within data, making it a powerful tool for forecasting and other related tasks.

## 5.6 Using STGCN on the dataset

In the project, we implement the Spatio-Temporal Graph Convolutional Network (STGCN) model to analyze our traffic dataset, following the methodology outlined in the referenced study. Our objective was to leverage this model to capture both spatial and temporal dependencies within the traffic data effectively. By constructing a graph where nodes represent traffic stations and edges depict the connectivity between these stations, we apply the STGCN model to predict traffic flow dynamics over time. The model utilizes convolutional layers to process the spatial information and gated convolutional layers for the temporal data, ensuring that both aspects are integrated seamlessly. Our output from this model includes accurate predictions of traffic conditions, such as vehicle flows and speeds, at various times and locations across the network. These predictions are crucial for real-time traffic management and long-term urban planning, aiming to enhance the understanding of traffic patterns and make informed data-driven decisions in traffic system optimization.

## 6 Results

The model was trained to forecast traffic flow for different time horizons (15, 30, and 45 minutes) using historical data from our training set. The training configuration included 34 days of data for training, 4 days for validation, and 4 days for testing. Key parameters were set as follows: spatial kernel size ( $k_s$ ) of 3, temporal kernel size ( $k_t$ ) of 3, a historical window size ( $n_{his}$ ) of 12, and a prediction horizon ( $n_{pred}$ ) of 9. Training was conducted over 35 epochs with a batch size of 50, using the RMSProp optimizer with a learning rate of 0.001. We used an inference mode of 'merge,' saving results every 10 epochs.

During the training phase, the model demonstrated strong predictive capabilities. For the shortest interval (15 minutes), it achieved relatively low errors across all performance metrics, indicating high accuracy and efficiency, as seen in Table 1. As we extended the forecast horizon to 30 and 45 minutes, the model maintained its predictive strength, showing only slight increases in error values due to the challenges of longer-term prediction.

In testing, the model continued to perform well, providing consistent and accurate predictions across all forecast intervals. The minimal discrepancy between training and testing results indicates a high level of generalization, proving that the STGCN model effectively captures spatial and temporal patterns.

These results confirm the efficacy of STGCN in traffic forecasting, offering timely and reliable predictions that can significantly aid urban traffic management. By understanding the spatio-temporal dependencies in traffic networks, this model provides practical insights that support real-time decision-making and long-term planning.

Table 1: Accuracy metrics

(a) Training				(b) Testing			
	MAPE(%)	MAE	RMSE		MAPE(%)	MAE	RMSE
15 mins	1.414	0.067	0.120	15 mins	1.287	0.068	0.121
30 mins	1.858	0.080	0.136	30 mins	1.511	0.080	0.138
45 mins	1.929	0.090	0.050	45 mins	0.365	0.091	0.152

## 7 Conclusion and Future Work

In this report, we have demonstrated the efficacy of the Spatio-Temporal Graph Convolutional Network (STGCN) model in accurately forecasting traffic patterns across 190 sensor stations using historical traffic data. The model combines graph convolutional layers for learning spatial patterns with gated temporal convolutional layers to capture traffic flow dynamics over time. The training results for short-term (15 minutes) and mid-term (30 and 45 minutes) horizons highlighted the model's robustness, achieving consistent and accurate predictions. The minimal discrepancy between training and testing results validated a high level of generalization, confirming that the model is well-suited to real-time traffic management and long-term planning.



## 7.1 Future Work and Scope

1. **Weather Data Integration:** Incorporate weather data to improve the model’s resilience and predictive accuracy. By including meteorological conditions, the model can better account for external factors affecting traffic patterns.
2. **Comparative Study:** Compare STGCN with newer graph models, such as Graph Attention Networks (GATs) and Dynamic Graph CNNs. This study will evaluate each model’s accuracy and computational efficiency to identify optimal architectures for traffic forecasting.
3. **Scalability and Relevance:** Adapt the model to evolving traffic patterns and data sources to ensure continued relevance. By dynamically updating the graph structure, the model can better handle changes in road networks and travel behavior.
4. **Refining Graph Structures:** Explore dynamic graph structures or real-time graph connectivity to better reflect evolving traffic conditions. Learning graph connectivity dynamically would adapt the model to changing travel behaviors.
5. **Advanced Neural Network Architectures:** Applying attention mechanisms, graph neural networks, or transformer models may further refine predictive capabilities by capturing spatial and temporal

## 8 Acknowledgement

We would like to express our sincere gratitude to Professor Gonzalo Mateos for his invaluable guidance and support throughout this project.

We also thank the California Department of Transportation for providing the PeMSD7 dataset, which was crucial to our work.

## References

- [1] Yu, B., Yin, H., & Zhu, Z. (2018). Spatio-Temporal Graph Convolutional Networks: A Deep Learning Framework for Traffic Forecasting. *IJCAI*.
- [2] Li, Y., Yu, R., Shahabi, C., & Liu, Y. (2018). Diffusion Convolutional Recurrent Neural Network: Data-driven Traffic Forecasting. *ICLR*.
- [3] Wu, Y., & Tan, H. (2016). Short-term Traffic Flow Forecasting with Spatial-Temporal Correlation in a Hybrid Deep Learning Framework. *arXiv:1612.01022*.
- [4] California Department of Transportation. (n.d.). Performance Measurement System (PeMS). Retrieved from <https://pems.dot.ca.gov/>.