# Spatiotemporal Taxi Demand Forecasting with MS-RNN and GNN in NYC

Tianqi Wang

April 2025

# Acknowledgements

First of all, I would like to express my sincere gratitude to my Dissertation supervisor, **Professor Ligang He**, for his meticulous guidance and selfless help in the process of topic selection, research design, and thesis writing. When I encountered difficulties, he patiently pointed me in the right direction and offered valuable academic advice.

I am also grateful to the senior members of Professor He's research group for their support and inspiration in data processing, experimental tuning, model implementation, and code debugging: **Yan Qian**, **Da Yan**, **Xuan Huang**, **Weitong Liao**, **Stephen Xu**, **Ashley Au**, **Xiao Qin**, **Zhihao Dai**, **Yiming Wei**, **Xuanyu Liu**, **Yuchen Liu**, **Yi Su**, and **Peng Jiang**.

My thanks go to the open source community and the developers of tools such as PyTorch, PyTorch Geometric, Pandas, and NumPy, whose generously shared code and documentation made this research far more efficient.

Finally, I would like to thank my family and my girlfriend, Yuanting Shi, for their understanding and support throughout my academic journey. Their encouragement has allowed me to focus on my research and continue moving forward.

Tianqi Wang
2025–04–20

**Abstract**

Accurate taxi passenger demand prediction is essential for efficient urban transportation management. We introduce a novel framework that fuses a Multi-Scale RNN (GRU for 1–24 h, LSTM for 1–7 days, Transformer for monthly seasonality) with GraphSAGE to jointly model temporal trends and spatial dependencies in New York City taxi data. GraphSAGE employs neighbor sampling to dynamically adapt to evolving origin–destination flows while keeping computational costs low. A rolling-window incremental training strategy updates only the latest data segments, enabling real-time responsiveness without full retraining. In benchmark comparisons against CSTN and STGCN, our approach reduces MAE to 19.8 and RMSE to 31.2, while cutting overall runtime by approximately 20%. This scalable, efficient model offers actionable demand forecasts and lays the groundwork for future extensions incorporating exogenous inputs (e.g., weather, events) and finer-grained spatial modeling at street or intersection levels.

# Contents

# 1  Introduction

## 1.1  Background and Context

Rapid urbanization and population growth have intensified urban mobility challenges globally, manifesting in increased traffic congestion, inefficient management of public transportation systems, and suboptimal allocation of resources [1]. Among various transportation modes, taxis play a significant role in urban mobility due to their flexibility and responsiveness to demand [15]. However, unpredictable fluctuations in taxi passenger demand frequently lead to inefficient dispatching, operational inefficiencies, and reduced passenger satisfaction. This demand–supply mismatch poses substantial challenges during distinct periods. During peak hours, events, or inclement weather conditions, passenger demand significantly exceeds taxi availability, resulting in prolonged waiting times, elevated frustration among passengers, and reduced overall satisfaction [19]. Conversely, during periods of low passenger demand, taxis frequently circulate empty, incurring unnecessary operational costs such as fuel and maintenance. This underutilization not only impacts driver earnings but also contributes to increased urban pollution and road congestion [28].

Complex traffic flow exhibits spatial correlations, such as interactions between neighboring regions, and temporal dependencies, including variations between peak and off–peak periods [25]. To address these intricacies, this research leverages a combination of Multi–Scale Recurrent Neural Networks (MS–RNN) and Graph Neural Networks (GNNs). Specifically, the MS–RNN framework employs Gated Recurrent Units (GRU) for short–term temporal patterns, Long Short–Term Memory (LSTM) networks for mid–term cyclic variations, and Transformer models for long–term periodic trends. This multi–scale temporal modeling ensures robust handling of demand dynamics across different horizons.

To capture spatial dependencies, we integrate GraphSAGE—a scalable GNN architecture capable of dynamically modeling relationships between taxi zones. By employing neighborhood sampling, GraphSAGE aggregates information from adjacent regions, accurately representing spatial interactions while significantly reducing computational complexity [12]. The resulting MS–RNN + Graph-SAGE framework addresses both spatial and temporal complexities in taxi demand forecasting, facilitating improved real–time dispatch decisions and optimizing urban mobility. Ultimately, this approach aims to alleviate economic losses, reduce emissions, and enhance urban quality of life by improving taxi utilization.

## 1.2  Motivation

Inaccurate forecasting of taxi demand leads to mismatches between supply and demand, with negative effects at multiple levels. For passengers, unreliable predictions cause longer wait times and increased uncertainty, lowering satisfaction and driving them toward alternative modes of transport [27]. Taxi drivers—especially less experienced ones—struggle to position themselves effi-

ciently without accurate forecasts, resulting in idle time, reduced earnings, and higher turnover rates, which further degrade service quality [27]. From an urban management perspective, poor forecasts exacerbate vehicle imbalance and empty cruising, fueling congestion and greenhouse–gas emissions. Ineffective dispatching creates traffic bottlenecks, undermining overall transportation efficiency [27].

Emerging advances in big data analytics and artificial intelligence enable more accurate demand prediction through deep learning. Traditional statistical methods often fail to capture complex spatiotemporal dynamics, whereas integrated RNN–GNN approaches have shown strong promise in modeling both temporal dependencies and spatial correlations inherent in urban taxi data.

**Keywords:** Multi–scale RNN; Graph Neural Network; GraphSAGE; Taxi Demand Forecasting; Spatiotemporal Data Modeling;Smart Mobility; Urban Traffic Management

## 1.3 Project Aim and Objectives

The primary aim of this research is to develop a robust, scalable, and computationally efficient spatiotemporal forecasting model for accurately predicting urban taxi passenger demand. Addressing the intricate patterns in both time and space domains is critical for optimizing taxi dispatch operations, improving service quality, and mitigating urban congestion. To achieve this aim, we set the following detailed objectives:

1. **Design a Multi Scale Recurrent Neural Network (MS RNN):** We will architect a multi branch recurrent network that simultaneously captures temporal patterns at different scales:

    - **GRU branch for short term (1–24 h) dynamics:** This branch models immediate, hour to hour fluctuations driven by localized factors such as weather changes or short lived public events.

    - **LSTM branch for mid term (1–7 day) cyclic patterns:** Capturing daily and weekly rhythms, this branch learns recurrent patterns like weekday commuter peaks versus weekend leisure demand.

    - **Transformer branch for long term (monthly) trends:** Employing a self attention mechanism, this branch detects broader seasonal variations, such as monthly tourism surges or policy driven changes in mobility behavior.

2. **Integrate GraphSAGE for Dynamic Spatial Modeling:** We will construct a dynamic graph where nodes represent taxi zones and edges are weighted by recent origin–destination flows. GraphSAGE will inductively sample and aggregate neighborhood information, allowing the model to dynamically adapt to shifting spatial dependencies without retraining from scratch.

3. **Implement Incremental Rolling Training for Real Time Adaptability:** Rather than retraining the entire model on cumulative data, we will adopt a rolling window approach that updates only on newly available data. This minimizes computational cost, enables daily model refreshing, and ensures the forecasting system remains responsive to evolving urban patterns.

4. **Visualize Results through Heat Maps and Dispatch Recommendations:** We will develop intuitive visualizations of predicted demand distributions across the city and derive actionable dispatch suggestions based on spatial and temporal demand imbalances. These outputs will provide valuable insights for taxi operators and urban planners seeking to optimize fleet allocation and reduce service gaps.

## 1.4 Contributions of the Project

This study makes several significant contributions that advance both the methodological landscape and the practical deployment of spatiotemporal forecasting systems:

- **Hybrid Spatiotemporal Architecture:** We propose a novel hybrid model that combines a Multi Scale Recurrent Neural Network (MS RNN) with a Graph Neural Network (GraphSAGE). By fusing GRU, LSTM, and Transformer architectures within the MS RNN, the model captures multi horizon temporal patterns more comprehensively than single scale models. The addition of GraphSAGE enables dynamic spatial refinement, further boosting predictive accuracy. Extensive experiments demonstrate that this hybrid design significantly outperforms classical methods such as pure GRU, LSTM, Transformer, and established baselines like CSTN and STGCN in both MAE and MSE metrics.

- **Efficient Rolling Training Strategy:** We implement an incremental training mechanism based on a sliding temporal window. Each day, the model fine tunes itself using only the newest 24 hours of data, leveraging warm start optimization without reinitializing parameters. This approach not only minimizes retraining overhead but also maintains model stability and generalization across time. Empirical results show that daily incremental updates complete in under two minutes, making real time deployment feasible even on moderate hardware setups (e.g., i7 CPUs and mid range GPUs).

- **Spatial Refinement via GraphSAGE:** Our use of GraphSAGE for post processing refinement represents a key technical innovation. By applying learned aggregation functions over dynamically sampled neighborhoods, the GNN corrects region specific biases in MS RNN predictions, smoothing inconsistencies while preserving local anomalies where appropriate. Experimental analysis shows that this refinement step yields substantial

additional gains, with a 12.7% reduction in MAE and a 25.5% reduction in MSE compared to MS RNN outputs alone.

- **Operational Insights for Real World Impact:** Beyond quantitative improvements, our framework provides practical outputs that are immediately useful for stakeholders. Predicted heat maps highlight spatial demand hotspots in advance, allowing fleet operators to proactively position taxis in high demand zones. Dispatch recommendations based on predicted imbalances help optimize vehicle distribution, reduce passenger wait times, and potentially lower city wide congestion levels by minimizing inefficient cruising behavior.

- **Benchmarking and Extensibility for Future Research:** We provide thorough benchmarking results comparing our approach against leading spatiotemporal models, such as CSTN and STGCN, under standardized evaluation protocols. Furthermore, the system is designed with extensibility in mind: future researchers can easily incorporate additional data sources (e.g., weather, events) or experiment with more advanced GNN variants (e.g., dynamic graphs, EvolveGCN). The modularity and reproducibility of our codebase lay a robust foundation for ongoing research in fine grained urban forecasting and dynamic graph learning.

Overall, this project delivers not only a technically sophisticated solution to a complex forecasting problem but also produces actionable outputs with clear relevance to real world urban transportation systems. By advancing both methodological innovation and practical deployment readiness, the proposed framework represents a meaningful step forward toward intelligent, data driven mobility management in smart cities.

# 2    Literature Review

## 2.1    Taxi Demand Prediction

### 2.1.1    Historical Methods and Traditional Forecasting Techniques

To build efficient urban transportation management, resource allocation, and improved passenger satisfaction, we require a accurate prediction of taxi passenger demand [27]. Historically, there are diverse methods that have been used to forecast taxi demand, evolving significantly from traditional statistical techniques to advanced machine learning and deep learning approaches. This section provides a valuable review for both historical methods and modern computational approaches, highlighting their strengths, limitations, and applicability to current taxi demand forecasting scenarios.

Because of interpretability and computational simplicity, early taxi demand forecasting efforts mainly relied on classical statistical methods. Common techniques included moving averages, exponential smoothing, and autoregressive integrated moving average (ARIMA) models[**?**]. Moving averages is one of the

earliest and simplest forecasting techniques, it calculate future taxi demand based on historical averages over specified periods. Although It's easy to implement, this method assumes demand stability and neglects temporal variations and trend shifts, which limiting its predictive power.[4]

Exponential smoothing is an enhancement over simple moving averages, assigns exponentially decreasing weights to older data points. This approach captures recent demand patterns better, improving short term forecasting accuracy. Despite these improvements, exponential smoothing models still struggle to account for complex, non linear demand fluctuations that are typical in dynamic urban environments.

ARIMA models emerged as an advancement over simpler techniques, explicitly modeling autocorrelations within historical demand data[4]. ARIMA effectively captures linear temporal dependencies by integrating moving average and autoregressive components. Seasonal ARIMA (SARIMA) is a variant of ARIMA that has been further developed to handle periodic patterns (e.g., daily or weekly cycles) commonly seen in taxi demand or traffic data[23]. Although these models demonstrated improved predictive capabilities, they inherently assume linearity and stationarity of time series data. Such assumptions frequently fail to represent real world taxi demand accurately, particularly during unusual events or abrupt demand shifts.

Other classical statistical approaches, such as regression models, have been employed to correlate taxi demand with external variables (e.g., weather conditions, public events, and economic indicators)[19]. Regression based methods enable a more detailed understanding of factors influencing taxi demand but often require extensive feature engineering and face challenges with multicollinearity, variable selection, and capturing non linear relationships[19].

### 2.1.2   Machine Learning and Deep Learning based Approaches

The advent of machine learning (ML) techniques introduced significant advancements in taxi demand prediction. These approaches overcome many limitations associated with traditional statistical methods by efficiently handling complex, non linear, and high dimensional datasets. Among these, supervised learning methods such as Support Vector Machines (SVM) [6], Random Forests (RF), ridge regression model (RRM), and combination forecasting model (CFM) (Liu et al., 2020) have demonstrated substantial improvements in predictive accuracy.

Support Vector Machine (SVM) models is recognized for their strong generalization capabilities, manage complex relationships within data by mapping input features into high dimensional spaces.[13] However, the practical deployment of SVM for large scale demand forecasting remains constrained by issues of computational complexity and sensitivity to hyperparameter tuning(Rafael Gomes Mantovani et al., 2015). In response, ensemble techniques such as Random Forests have become popular due to their robustness, interpretability, and ability to effectively capture nonlinear interactions among variables(Auret and Aldrich, 2012). Similarly, Ridge Regression Models (RRM) have emerged as viable al-

ternatives, leveraging regularization techniques to mitigate multicollinearity and enhance predictive stability, particularly in scenarios with numerous correlated features[3]. Moreover, Combination Forecasting Models (CFM), which integrate predictions from multiple distinct forecasting methodologies, have demonstrated superior performance by harnessing the complementary strengths of individual models, leading to improvements in forecast accuracy, reliability, and robustness.

In recent years, deep learning methods have prominently emerged in taxi demand prediction tasks, driven by their exceptional ability to capture intricate spatiotemporal dependencies within large scale datasets[16]. Recurrent Neural Networks (RNNs), particularly Long Short Term Memory (LSTM) and Gated Recurrent Units (GRU) have gained widespread use owing to their effectiveness in modeling complex sequential temporal patterns[2]. These architectures successfully address the vanishing gradient problem inherent in traditional RNNs, enabling them to capture both short and mid term temporal dependencies efficiently[2]. Furthermore, these models have demonstrated notable success in scenarios involving periodic fluctuations, such as daily and weekly cycles in taxi demand.

Convolutional Neural Networks (CNNs) have also found robust applications in demand forecasting, particularly when spatial dependencies and local geographical features are significant[7]. By effectively extracting spatial features, CNNs facilitate the modeling of intricate spatial relationships. Hybrid CNN–LSTM architectures, in particular, combine the spatial modeling strengths of CNNs with the temporal modeling capabilities of LSTMs, yielding highly accurate predictions in complex spatiotemporal scenarios[14].

Graph Neural Networks (GNNs), encompassing Graph Convolutional Networks (GCN), Graph Attention Networks (GAT), and GraphSAGE, represent cutting edge methodologies in capturing spatial relationships within urban regions. These models efficiently represent spatial structures by leveraging graph based data representations, enabling dynamic and accurate modeling of interactions between taxi zones[18]. GraphSAGE, specifically, enhances scalability and flexibility by employing neighborhood sampling techniques, thus effectively capturing relevant local information while maintaining computational efficiency. The integration of GNNs with multi scale RNN frameworks further enhances the modeling capability, accurately addressing both spatial interactions and temporal dynamics in urban taxi demand[18].

Moreover, Transformer based models, originally designed for natural language processing, have recently shown promising potential for taxi demand prediction[10]. Their self attention mechanisms uniquely capture long range temporal dependencies and global contextual information within datasets. This capability allows Transformers to model complex temporal correlations more effectively than traditional sequential neural networks, especially for extended forecasting horizons[10]. The combination of Transformers with other deep learning architectures presents opportunities for even more sophisticated predictive models, potentially revolutionizing future forecasting methodologies in urban transportation management.

## 2.2 Recurrent Neural Networks (RNN)

Recurrent Neural Networks (RNNs) are a class of neural architectures specifically designed to process sequential data[20]. Unlike traditional feedforward networks, RNNs possess internal loops that allow them to maintain a hidden state vector, which evolves over time as new elements of a sequence are observed. This hidden state serves as a form of memory, enabling the network to capture temporal dependencies between sequential inputs.

The mathematical formulation of an RNN at time step $t$ can be described as:

$$h_t = \sigma(W_{xh}x_t + W_{hh}h_{t1} + b_h)$$

$$y_t = W_{hy}h_t + b_y$$

where $h_t$ is the hidden state, $x_t$ is the input at time $t$, $y_t$ is the output, $W$ are weight matrices, $b$ are biases, and $\sigma$ is a nonlinear activation function such as tanh or ReLU.

The recurrent structure allows RNNs to leverage both the current input and the historical context encoded in $h_{t1}$. This makes RNNs particularly suitable for tasks where the order of inputs is crucial, including language modeling, financial time series forecasting, and transportation demand prediction.

In the context of taxi demand forecasting, RNNs offer a natural framework for learning patterns over time, such as daily commuting cycles, weekend fluctuations, and special event impacts[24]. By training on historical passenger pick up data, RNNs can infer the temporal dynamics underlying demand variations and generalize to predict future demand patterns.

Despite their conceptual simplicity and broad applicability, traditional RNNs suffer from notable challenges when modeling long sequences. Chief among these are the vanishing gradient and exploding gradient problems, which impede the network's ability to learn long range dependencies during backpropagation through time (BPTT). In response to these limitations, more advanced architectures such as Long Short Term Memory (LSTM) networks and Gated Recurrent Units (GRU) were developed.

### 2.2.1 Long Short Term Memory (LSTM) and Gated Recurrent Units (GRU)

Traditional RNNs often struggle to capture long term dependencies due to the inherent nature of gradient propagation across many time steps. Specifically, during backpropagation, gradients can decay exponentially (vanishing gradients) or grow uncontrollably (exploding gradients), making learning unstable or ineffective. To address these challenges, specialized RNN variants were proposed, most notably Long Short Term Memory (LSTM) networks and Gated Recurrent Units (GRU)[2].

**Long Short Term Memory (LSTM)**    LSTM networks, introduced by Hochreiter and Schmidhuber (1997), fundamentally modify the RNN structure by introducing a memory cell $c_t$ alongside the hidden state $h_t$. The cell state acts as a conduit for preserving information across long sequences, regulated by three gating mechanisms: the input gate, forget gate, and output gate[20].

The core computations in an LSTM cell are as follows:

$$f_t = \sigma(W_f \cdot [h_{t1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t1}, x_t] + b_c)$$

$$c_t = f_t \odot c_{t1} + i_t \odot \tilde{c}_t$$

$$o_t = \sigma(W_o \cdot [h_{t1}, x_t] + b_o)$$

$$h_t = o_t \odot \tanh(c_t)$$

where $\sigma$ denotes the sigmoid function and $\odot$ denotes element wise multiplication.

Through these gates:

- The **forget gate** decides which past information should be discarded.

- The **input gate** controls what new information is incorporated.

- The **output gate** determines the exposure of internal memory to the next time step.

This sophisticated gating allows LSTMs to maintain and manipulate information over arbitrarily long horizons, making them especially suitable for sequence modeling tasks where long term dependencies are critical.

**Gated Recurrent Unit (GRU)**    The GRU simplifies the LSTM architecture while retaining similar capabilities. A GRU combines the forget and input gates into a single update gate and merges the cell state and hidden state into a single entity. This simplification leads to faster computation and easier parameter optimization without significantly compromising modeling power[17].

The GRU updates are given by:

$$z_t = \sigma(W_z x_t + U_z h_{t1})$$

$$r_t = \sigma(W_r x_t + U_r h_{t1})$$

$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t1}))$$

$$h_t = (1 z_t) \odot h_{t1} + z_t \odot \tilde{h}_t$$

Here:

- The **update gate** $z_t$ balances the previous hidden state and the candidate hidden state.

- The **reset gate** $r_t$ controls how much of the previous state to forget.

**Comparison Between LSTM and GRU**

- LSTM is more expressive and better suited for very long sequences with complex dependencies[?, ?].

- GRU is computationally lighter, leading to faster convergence and suitability for scenarios with limited computational budgets or where real time inference is required[2].

**Application to Taxi Demand Forecasting** In urban mobility forecasting, both LSTM and GRU have been successfully applied to capture recurrent patterns such as morning/evening rush hours, weekend effects, and holiday season anomalies. Their flexibility in handling variable length sequences and resilience against noisy data make them highly valuable for modeling the inherently stochastic nature of taxi demand.

### 2.2.2 Transformer Models

Although RNN variants like LSTM and GRU excel at modeling sequential data, they still process sequences step by step, making them less parallelizable and prone to difficulties when handling very long sequences. To overcome these limitations, Transformer architectures were introduced by Vaswani et al. (2017)[21], fundamentally changing the landscape of sequence modeling.

**Transformer Fundamentals** Transformers rely entirely on attention mechanisms to model dependencies between sequence elements, discarding the need for recurrence altogether. The key innovation is the **self attention** mechanism, where each position in a sequence attends to all other positions, allowing the model to capture long range dependencies without regard to the order of processing.

Scaled dot product attention computes attention scores as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

where $Q$, $K$, and $V$ are the query, key, and value matrices respectively, and $d_k$ is the dimension of the key vectors.

**Multi Head Attention**  To allow the model to jointly attend to information from different representation subspaces, Transformers employ multi head attention, running several self attention operations in parallel and concatenating their outputs.

**Positional Encoding**  Since Transformers lack recurrence, they inject positional information into the input embeddings via learned or fixed positional encodings, enabling the model to preserve sequence order information.

**Advantages Over RNNs**

- Fully parallelizable sequence processing, leading to faster training times.

- Direct modeling of arbitrarily distant dependencies without degradation.

- Superior scalability to very long sequences (e.g., thousands of steps).

**Transformer in Taxi Demand Forecasting**  For taxi demand forecasting, Transformer based models can capture long term seasonal patterns (e.g., monthly dips, holiday surges) that span beyond the memory limits of conventional RNNs. By observing demand sequences holistically, the model can associate distant time steps and make globally coherent predictions.

**Limitations**  Despite their strengths, Transformers can struggle with capturing fine grained local patterns unless augmented with local attention windows or hybrid architectures. For tasks heavily reliant on short term signals, RNN components may still outperform.

**Hybrid Architectures**  Given the complementary strengths of RNNs and Transformers, many state of the art models adopt hybrid designs, combining GRU/LSTM branches for immediate signals and Transformer branches for long term trend modeling. This multi scale approach, used in our framework, leverages the best of both paradigms.

### 2.2.3 Multi scale RNN Methods

Taxi demand prediction typically involves analyzing temporal dependencies occurring at different scales—hourly fluctuations, daily cycles, weekly patterns, and even monthly seasonal variations. Multi scale RNN (MS RNN) methods have been developed to address the complexities associated with modeling such diverse temporal dynamics within a single integrated framework.

The MS RNN framework typically combines different RNN architectures, such as GRU, LSTM, and Transformers, each dedicated to capturing specific temporal scales. GRU models are usually applied to capture short term dependencies (hourly fluctuations), as their computational efficiency enables rapid processing and accurate predictions for immediate demand scenarios. LSTM models are employed for mid term patterns (daily and weekly cycles) due to their robust memory retention capabilities, enabling accurate modeling of recurring demand trends.

Transformer models within MS RNN frameworks are utilized to capture long term temporal dependencies, such as monthly or seasonal variations. Their powerful self attention mechanism effectively captures global context, enhancing the accuracy of long range forecasting. By integrating these diverse RNN architectures within a cohesive framework, the MS RNN approach achieves comprehensive temporal modeling, significantly outperforming single scale approaches in predictive accuracy and robustness.

Additionally, MS RNN methods often incorporate adaptive mechanisms, such as dynamic weighting of different temporal scales, allowing the model to flexibly emphasize the most relevant temporal features under varying circumstances. These adaptive techniques further enhance the model's predictive performance and reliability across different forecasting horizons. Furthermore, extensive experimentation has demonstrated the practical applicability of MS RNN methods in real world scenarios, showing superior performance in accurately forecasting taxi demand compared to traditional single scale methods. As a result, MS RNN frameworks have become increasingly popular and valuable for complex urban transportation management applications.

Moreover, feature fusion layers are typically implemented within MS RNN architectures, combining the distinct temporal features extracted by GRU, LSTM, and Transformer components. This integration is commonly achieved through fully connected layers, enabling the model to learn complex interactions between different temporal scales effectively. The fusion of multi scale features results in significantly improved forecasting performance, making MS RNN frameworks highly suitable for practical taxi demand prediction scenarios.

## 2.3 Graph Neural Networks (GNN)

Graph Neural Networks (GNNs) have emerged as powerful tools for effectively analyzing and modeling data structured in graphs. Unlike traditional neural network architectures designed primarily for grid like data (e.g., images in CNNs) or sequential data (e.g., text in RNNs), GNNs are specialized to exploit the rela-

tionships and interactions inherent in graph structured data. Graphs naturally arise in a variety of real world domains, including social networks, molecular chemistry, recommendation systems, knowledge graphs, and transportation networks[18].

The fundamental principle of GNNs is to iteratively update the representation of a node by aggregating information from its local neighbors, thereby capturing both node level features and the broader relational context within the graph. This neighbor aware feature learning allows GNNs to reason over complex topologies and relational patterns that are otherwise difficult to model using conventional neural architectures.

In this section, we elaborate on three prominent and representative GNN architectures: Graph Convolutional Networks (GCN), Graph Attention Networks (GAT), and Graph Sample and Aggregate (GraphSAGE). Each model offers distinct mechanisms for information propagation and aggregation, with unique strengths and limitations.

### 2.3.1 Graph Convolutional Networks (GCN)

Graph Convolutional Networks (GCNs), introduced by Kipf and Welling (2016) [8], represent one of the foundational approaches for extending convolutional neural networks to irregular, non Euclidean graph structures. GCNs aim to perform localized, spectral convolutions by operating in the graph frequency domain

The basic layer wise propagation rule for a GCN is given by:

$$H^{(l+1)} = \sigma \left( \tilde{D}^{1/2} \tilde{A} \tilde{D}^{1/2} H^{(l)} W^{(l)} \right),$$

where:

- $H^{(l)}$ is the matrix of node features at layer $l$,

- $\tilde{A} = A + I$ is the adjacency matrix with added self loops,

- $\tilde{D}$ is the diagonal node degree matrix of $\tilde{A}$,

- $W^{(l)}$ is the learnable weight matrix,

- $\sigma(\cdot)$ is an activation function, typically ReLU.

The symmetric normalization $\tilde{D}^{1/2} \tilde{A} \tilde{D}^{1/2}$ ensures numerical stability and balances contributions from neighbors with varying degrees.

**Advantages of GCNs:**

- Efficient and simple to implement.

- Good at semi supervised node classification tasks.

- Provides smooth node embeddings by averaging features over local neighborhoods.

**Limitations of GCNs:**

- Equal weighting of neighbors may be suboptimal when neighbor importance is heterogeneous.

- Scalability issues arise for very large graphs, as full adjacency matrices must be processed.

- Over smoothing problem: with deeper layers, node embeddings tend to become indistinguishable.

**Applications:** GCNs have been widely used in citation network classification, protein protein interaction modeling, and recommendation systems where relational structures are relatively stable and small to medium scale.

### 2.3.2   Graph Attention Networks (GAT)

Graph Attention Networks (GATs), introduced by Veličković et al. (2018) [22], extend the GNN paradigm by introducing attention mechanisms into neighborhood aggregation. Unlike GCNs, which uniformly aggregate all neighboring features, GATs allow a node to learn dynamic weights for each of its neighbors based on learned attention scores[22].

The core idea involves computing attention coefficients $e_{ij}$ between node $i$ and each neighbor $j$:

$$e_{ij} = \text{LeakyReLU}\left(\mathbf{a}^\top \left[W\mathbf{h}_i \,\|\, W\mathbf{h}_j\right]\right),$$

where $\|$ denotes concatenation, $W$ is a learnable linear transformation, and $\mathbf{a}$ is a learnable attention vector.

The normalized attention coefficients are then computed via a softmax:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in \mathcal{N}(i)} \exp(e_{ik})},$$

and used to perform weighted aggregation:

$$\mathbf{h}_i' = \sigma\left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij} W\mathbf{h}_j\right).$$

**Advantages of GATs:**

- Nodes can focus selectively on the most relevant neighbors.

- Better handles graphs with highly varying neighborhood structures.

- Supports inductive settings (generalization to unseen nodes).

**Limitations of GATs:**

- Computationally intensive due to pairwise attention score calculations.

- Less stable on large scale, sparse graphs unless sampling methods are employed.

**Applications:** GATs excel in tasks requiring fine grained relational reasoning, such as social network analysis, fraud detection, and knowledge graph completion, where neighbor importance varies substantially.

### 2.3.3 GraphSAGE

GraphSAGE (Graph Sample and Aggregate), introduced by Hamilton et al. (2017) [5], is designed to address scalability and generalization challenges in graph learning. Unlike GCNs and GATs, which assume access to the entire graph during training (transductive learning), GraphSAGE enables inductive learning by generating embeddings for unseen nodes via sampling and aggregating local neighborhoods[12].

The GraphSAGE framework consists of the following steps:

1. **Sampling**: For each node, a fixed size subset of its neighbors is sampled.

2. **Aggregation**: Neighbor embeddings are aggregated using a function such as mean, LSTM aggregator, or pooling.

3. **Update**: The node's own embedding is updated based on its current embedding and the aggregated neighbor information.

Mathematically, the neighborhood aggregation at layer $k$ is:

$$\mathbf{h}_i^{(k)} = \sigma \left( W^{(k)} \cdot \text{AGGREGATE} \left( \{\mathbf{h}_i^{(k1)}\} \cup \{\mathbf{h}_j^{(k1)}, \forall j \in \mathcal{N}(i)\} \right) \right),$$

where $\sigma$ is a nonlinearity, and $W^{(k)}$ are learnable weights.

**Advantages of GraphSAGE:**

- Supports inductive generalization to new, unseen nodes and graphs.

- Enables mini batch training via neighborhood sampling, improving scalability.

- Flexible aggregator choices allow adaptation to different graph structures.

**Limitations of GraphSAGE:**

- Sampling introduces variance; training stability depends on sampling quality.

- Limited receptive field unless multiple aggregation layers are stacked.

**Applications:** GraphSAGE is widely used in web scale recommendation systems (e.g., Pinterest, Twitter), dynamic social networks, and real time fraud detection platforms, where new nodes and edges are continuously added.

### 2.3.4 Comparative Analysis

While all three architectures—GCN, GAT, and GraphSAGE—are designed to learn effective node representations by exploiting graph structures, they differ fundamentally in their aggregation strategies and application scenarios.

GCNs are effective for small to medium static graphs where neighbor importance is relatively uniform. GATs shine in graphs with heterogeneous or highly variable neighborhood structures by learning attention based importance. GraphSAGE offers the best scalability and generalization to unseen nodes, making it the architecture of choice for very large or dynamically evolving graphs.

In practice, the choice among GCN, GAT, and GraphSAGE depends on the specific characteristics of the target graph (e.g., size, density, homogeneity) and the computational constraints of the deployment environment.

### 2.3.5 Comparative Summary

In summary, GCN, GAT, and GraphSAGE represent distinct yet complementary approaches to Graph Neural Network architectures. GCNs provide foundational insights into graph convolutional operations, efficiently capturing structural information through uniform neighbor aggregation. GAT extends this concept by incorporating dynamic attention mechanisms, providing enhanced representational capabilities through adaptive neighbor weighting. GraphSAGE further improves on scalability and inductive capabilities, enabling effective embedding generation in large scale and evolving graphs.

Collectively, these architectures provide a comprehensive toolkit for addressing various graph based challenges, including node classification, link prediction, and dynamic network analysis, thereby significantly advancing the capability to manage and analyze complex relational data in diverse real world scenarios.

### 2.3.6 Application of GNNs in Traffic Prediction

Graph Neural Networks (GNNs) have recently demonstrated remarkable success in modeling urban transportation systems, particularly for traffic prediction tasks. Traffic systems inherently exhibit complex spatial relationships characterized by the interconnectivity of road networks, intersections, and urban regions. Traditional methods often overlook or insufficiently capture these spatial correlations, resulting in limited predictive accuracy. GNNs, by contrast, explicitly exploit the graph based nature of transportation networks, offering substantial improvements in capturing spatial dependencies.

In traffic prediction, transportation networks can naturally be modeled as graphs where nodes represent distinct urban regions, intersections, or road segments, and edges represent the physical or functional connections between these

elements. GNNs effectively learn the interactions and spatial dependencies between nodes by aggregating information from neighboring nodes. This capability allows GNN models to accurately capture local and global patterns in traffic flow, significantly improving predictive accuracy compared to conventional methods.

Graph Convolutional Networks (GCNs) were among the first GNN architectures adopted for traffic prediction tasks due to their effectiveness in aggregating neighborhood information. Researchers have leveraged GCNs to predict short term traffic conditions, including traffic speed, volume, and congestion levels. GCN models have shown significant predictive advantages over traditional statistical and machine learning approaches, especially by effectively modeling recurrent congestion patterns and identifying critical nodes within urban transportation networks.

Graph Attention Networks (GATs) further enhance traffic prediction by introducing attention mechanisms to assign differentiated importance to neighboring nodes dynamically. Traffic conditions at specific intersections or road segments often depend disproportionately on certain critical neighbors, such as major intersections or heavily congested nodes. GAT's ability to prioritize relevant neighbors during information aggregation makes it particularly suitable for accurately modeling such nuanced spatial interactions. Studies employing GAT architectures have successfully demonstrated improved predictive accuracy for tasks like congestion forecasting, accident risk assessment, and traffic incident prediction.

GraphSAGE offers substantial advantages for large scale and dynamic traffic prediction scenarios, owing to its inductive learning capabilities and computational efficiency. Urban transportation networks frequently evolve due to changes such as road maintenance, construction, or temporary road closures. GraphSAGE efficiently adapts to such evolving scenarios by generating embeddings for newly introduced nodes without requiring extensive retraining. Moreover, its neighborhood sampling and mini batch training enable GraphSAGE models to handle large urban networks with thousands or even tens of thousands of nodes, maintaining scalability and responsiveness crucial for real time traffic prediction systems.

In practical applications, hybrid models combining GNNs with temporal modeling frameworks, such as recurrent neural networks (RNN) and transformer architectures, have become increasingly popular. These hybrid approaches effectively capture both spatial and temporal dependencies inherent in traffic data. Specifically, integrating GNNs for spatial modeling with multi scale RNN models enables comprehensive spatiotemporal forecasting, accurately predicting traffic conditions across various time horizons, from immediate future conditions to longer term patterns.

Applications of GNN based traffic prediction models extend beyond mere congestion forecasting. They support dynamic traffic routing, optimization of transportation infrastructure, management of urban public transit systems, and enhancement of intelligent transportation systems (ITS). These predictive models empower transportation agencies and urban planners to make informed

decisions, optimize resource allocation, and implement proactive measures to mitigate congestion, thereby significantly enhancing urban mobility and sustainability.

## 2.4   Summary and Gap Identification

This chapter has extensively reviewed diverse methodologies in taxi demand prediction, ranging from traditional statistical models to advanced neural network frameworks. Historical methods such as moving averages, exponential smoothing, and ARIMA provided foundational insights but demonstrated significant limitations in capturing complex, non linear temporal dynamics and spatial dependencies in urban taxi demand data.

Advances in machine learning significantly enhanced predictive capabilities through techniques like Support Vector Machines (SVM), Random Forests (RF), and Gradient Boosting Machines (GBM). These methods offered substantial improvements in robustness and accuracy, successfully capturing intricate non linear relationships. However, challenges remained in scalability, computational efficiency, and adaptability to large, dynamic datasets.

Deep learning techniques, particularly Recurrent Neural Networks (RNN), Convolutional Neural Networks (CNN), and Transformer architectures, marked significant breakthroughs in capturing sophisticated temporal and spatial patterns. LSTM and GRU networks effectively addressed temporal modeling challenges associated with sequential data. CNN architectures excelled in spatial feature extraction, while Transformer models enhanced long range temporal forecasting. Multi scale RNN frameworks further integrated these methods, providing comprehensive temporal modeling solutions.

Graph Neural Networks (GNNs) introduced innovative capabilities to model spatial dependencies inherent in urban transportation systems effectively. Architectures like GCN, GAT, and GraphSAGE demonstrated superior performance in capturing spatial interactions, significantly improving predictive accuracy. GraphSAGE, in particular, introduced essential inductive learning capabilities, offering efficient scalability and adaptability crucial for dynamic traffic prediction scenarios.

Despite these advancements, several research gaps remain:

- Effectively integrating incremental and online learning paradigms with complex neural network models (e.g., multi scale RNNs and GNNs) remains under explored.

- Scalability and computational efficiency of advanced architectures, particularly Transformers and hybrid models, need improvement for real time, large scale deployment.

- Systematic incorporation of external factors (weather, events, socioeconomic data, urban policies) into predictive models could boost accuracy and reliability.

- Robustness to noisy, incomplete, or erroneous real world transportation data requires further development to handle uncertainty and improve practical applicability.

Addressing these identified gaps will advance taxi demand forecasting methodologies, significantly improving their practical effectiveness and reliability in complex, dynamic urban transportation systems.

# 3 Methods

## 3.1 Multi scale RNN Model

The Multi scale Recurrent Neural Network (MS RNN) model lies at the core of our framework for predicting hourly taxi demand across New York City's urban regions. Its design is motivated by the observation that urban mobility patterns exhibit temporal dependencies at multiple time scales: immediate (within the last few hours), periodic (daily and weekly), and long term seasonal trends (monthly or more). The MS RNN framework is explicitly designed to model and fuse information from all these temporal levels.

This section provides an in depth explanation of each sub module in the MS RNN architecture—how they are implemented in code and why they are used. We also explain how the multi scale representations are fused and how this architecture enables rolling, incremental training without full retraining.

### 3.1.1 Architecture Overview

The MultiScaleModel is meticulously designed as a subclass of the nn.Module class in the PyTorch framework, enabling a highly modular, extensible, and efficient implementation. Its architecture is strategically constructed to capture the complex multiscale temporal dynamics observed in urban mobility demand data, which exhibit fluctuations at varying periodicities ranging from hourly noise to monthly seasonal trends.

To accurately model short to medium term temporal dependencies, the model incorporates two dedicated Long Short Term Memory (LSTM) branches. The first LSTM branch is configured to extract daily (1 day) temporal patterns, capturing regular variations such as morning and evening rush hours. The second LSTM branch is specialized for weekly (1 week) periodic structures, which are critical for modeling behavioral trends that vary across different days of the week, such as weekday commuting patterns versus weekend leisure travel. Each LSTM branch processes its respective time scaled input independently, enabling the network to learn distinct temporal features at these two important scales.

To address longer term dependencies that span over several weeks or months, the MultiScaleModel integrates a Transformer based branch. The Transformer, known for its self attention mechanisms and superior performance in modeling long range dependencies, is specifically tasked with capturing monthly (1 month) periodic trends. This branch allows the model to incorporate knowledge about

seasonal variations, special events, or other low frequency phenomena that significantly influence passenger demand but may not be detectable through short term recurrence alone.

Following the independent extraction of temporal features from the LSTM and Transformer branches, a feature fusion layer is employed to integrate these multiscale representations into a unified latent vector. This fusion layer performs a weighted combination and nonlinear transformation of the outputs from the different branches, ensuring that complementary temporal information across multiple scales is effectively aggregated while preserving their distinctive contributions.

To further refine the fused temporal features and to capture the short term, high frequency fluctuations that occur at the hourly level, the model incorporates an additional Gated Recurrent Unit (GRU) layer. The GRU is chosen for its relatively lightweight architecture compared to the LSTM, making it well suited for modeling rapid variations without incurring excessive computational overhead. Through this layer, the model dynamically adjusts to immediate changes in demand patterns, such as sudden surges due to weather changes or public events.

Finally, the output from the GRU layer is passed through a fully connected linear layer, which maps the high dimensional latent representation into the final prediction space. This linear projection produces the model's prediction for the hourly passenger demand at the targeted spatial resolution. By stacking these carefully designed components, the MultiScaleModel achieves a comprehensive understanding of both the short term variations and the long term periodic structures, enabling robust and accurate forecasting of urban mobility demand across multiple temporal scales.

```python
class MultiScaleModel(nn.Module):
def **init**(self, hidden\_size):
super().**init**()
\# LSTM for daily and weekly cycles
self.lstm\_1d = nn.LSTM(1, hidden\_size, batch\_first=
    True)
self.lstm\_1w = nn.LSTM(1, hidden\_size, batch\_first=
    True)
\# Transformer for monthly trends
self.input\_projection = nn.Linear(1, hidden\_size)
self.transformer\_1m = nn.Transformer(
d\_model=hidden\_size,
nhead=4,
num\_encoder\_layers=2,
batch\_first=True
)
\# Fuse long term features
self.feature\_fusion = nn.Linear(hidden\_size \* 3,
    hidden\_size)
\# GRU for short term plus fused trend
self.gru = nn.GRU(hidden\_size + 1, hidden\_size, batch
    \_first=True)
\# Final hourly demand prediction
self.fc = nn.Linear(hidden\_size, 1)

'''
def forward(self, x):
    x_1h, x_1d, x_1w, x_1m = x['1h'], x['1d'], x['1w'],
        x['1m']
    # Daily pattern
    _, (h1d, _) = self.lstm_1d(x_1d)
    # Weekly pattern
    _, (h1w, _) = self.lstm_1w(x_1w)
    # Monthly trend
    proj = self.input_projection(x_1m)
    proj = proj.permute(1, 0, 2)
    h1m = self.transformer_1m(proj, proj)[ 1]
    # Fuse long term features
    fused = torch.cat([h1d[ 1], h1w[ 1], h1m], dim=1)
    fused = self.feature_fusion(fused)
    # Combine with hourly input
    bsz, seq_len, _ = x_1h.size()
    fused_exp = fused.unsqueeze(1).repeat(1, seq_len,
        1)
    gru_in = torch.cat([x_1h, fused_exp], dim=2)
    _, h_gru = self.gru(gru_in)
    return self.fc(h_gru[ 1])
'''
```

### 3.1.2 Long term Trend Extraction with Transformer (1 Month)

The Transformer module is responsible for extracting long term periodic trends from sequences spanning one month. It first receives a projected version of the monthly input sequence through an `input_projection` layer, which transforms the original 1 dimensional passenger count into a higher dimensional hidden representation that matches the expected input size of the Transformer encoder.

```
x_1m_proj = self.input_projection(x["1m"])
x_1m_proj = x_1m_proj.permute(1, 0, 2)  # Transformer
    expects (seq_len, batch, hidden)
h_1m = self.transformer_1m(x_1m_proj, x_1m_proj)[ 1]  #
    Use final timestep
```

The Transformer operates on the projected sequence, capturing global dependencies across the full monthly horizon. Its multi head attention mechanism allows the model to jointly attend to information from different representation subspaces, dynamically focusing on the most relevant past time steps for forecasting future demand. Formally, given input queries $Q$, keys $K$, and values $V$, multi head attention computes:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^{\top}}{\sqrt{d_k}}\right) V,$$

where $d_k$ is the dimension of the key vectors. By using multiple heads, the Transformer can model complex interactions between different parts of the sequence.

The final output of the Transformer encoder is a condensed vector $h_{1m}$ summarizing the dominant monthly patterns. This vector helps the overall model capture slower dynamics such as monthly dips in taxi demand due to holidays, seasonal weather variations, or longer term policy changes (e.g., fare adjustments or traffic regulations).

### 3.1.3 Mid term Patterns with LSTM (1 Day and 1 Week)

To capture mid term cyclic behaviors, we employ two separate LSTM layers: one designed to learn daily (24 hour) patterns and another to model weekly (168 hour) cycles. These cycles are critical in urban settings where commuting behavior, weekend leisure activities, and weekly recurring events influence demand.

Each LSTM updates its hidden and memory states according to the following equations:

$$i_t = \sigma(W_i x_t + U_i h_{t1}), \quad f_t = \sigma(W_f x_t + U_f h_{t1}),$$

$$o_t = \sigma(W_o x_t + U_o h_{t1}), \quad \tilde{c}_t = \tanh(W_c x_t + U_c h_{t1}),$$

$$c_t = f_t \odot c_{t1} + i_t \odot \tilde{c}_t, \quad h_t = o_t \odot \tanh(c_t),$$

where $i_t$, $f_t$, and $o_t$ represent the input, forget, and output gates respectively, and $\tilde{c}_t$ is the candidate cell state.

The LSTM layers are implemented as:

```
1  _, (h_1d, _) = self.lstm_1d(x["1d"])
2  _, (h_1w, _) = self.lstm_1w(x["1w"])
```

Here, `h_1d` captures the compressed information of a typical day, such as morning and evening rush hours, while `h_1w` summarizes the patterns across a full week, distinguishing between workdays and weekends. These mid term memory structures enable the model to anticipate regular temporal patterns that are not captured by immediate local trends.

### 3.1.4  Feature Fusion

After extracting temporal representations at three different scales—daily, weekly, and monthly—the next critical step is to fuse these heterogeneous signals into a unified latent vector. To achieve this, the three hidden states $h_{1d}$, $h_{1w}$, and $h_{1m}$ are concatenated into a single vector in $R^{3H}$, where $H$ denotes the hidden size of each branch.

This combined vector is then projected back into $R^H$ using a linear transformation, followed by GELU activation to introduce nonlinearity and optional dropout for regularization:

$$h_{\text{fused}} = W_f [h_{1d}; h_{1w}; h_{1m}] + b_f, \quad h_{\text{fused}} \leftarrow \text{GELU}(h_{\text{fused}}), \quad h_{\text{fused}} \leftarrow \text{Dropout}(h_{\text{fused}}).$$

This process is implemented as:

```
1  # Concatenate daily, weekly, and monthly hidden states
2  fused_trend = torch.cat([h_1d[ 1], h_1w[ 1], h_1m], dim
       =1)
3  # Project back to hidden dimension
4  fused_trend = self.feature_fusion(fused_trend)
5  # Apply GELU activation and dropout
6  fused_trend = F.gelu(fused_trend)
7  fused_trend = self.dropout(fused_trend)
```

Listing 2: Feature Fusion in MS RNN

Through this fusion mechanism, the model effectively combines short term, mid term, and long term temporal signals into a holistic representation, enriching its ability to understand and predict complex temporal dynamics.

### 3.1.5 Short term Fluctuations with GRU (1 Hour)

While longer horizons provide valuable global context, immediate local fluctuations must also be accurately modeled to produce reliable hour ahead forecasts. To capture these fine grained temporal dependencies, we employ a GRU network as the final recurrent module.

The fused trend vector $h_{\text{fused}}$ is repeated across the hourly time sequence and concatenated with the raw hourly input features. This combined input is then fed into the GRU:

```
batch_size, seq_len, _ = x["1h"].shape
fused_exp = fused_trend.unsqueeze(1).repeat(1, seq_len,
    1)
x_gru_input = torch.cat([x["1h"], fused_exp], dim=2)
_, h_gru = self.gru(x_gru_input)
```

The GRU updates its hidden states based on the following gated operations:

$$z_t = \sigma(W_z x_t + U_z h_{t1}), \quad r_t = \sigma(W_r x_t + U_r h_{t1}),$$

$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t1})), \quad h_t = z_t \odot h_{t1} + (1 z_t) \odot \tilde{h}_t.$$

This architecture allows the GRU to retain relevant recent information while selectively incorporating new inputs, making it well suited for modeling short term volatility in urban demand.

### 3.1.6 Final Prediction

The final hidden state output by the GRU encapsulates the distilled temporal features across all horizons. This hidden state is passed through a fully connected linear layer to generate the final scalar prediction, representing the forecasted taxi demand for the next hour:

```
return self.fc(h_gru[ 1])
```

This simple yet effective output layer ensures that all learned temporal information contributes directly to the final demand estimate.

## 3.2 Graph Neural Network (GraphSAGE)

In this section, we introduce the spatial modeling strategy used in our framework, focusing on Graph Neural Networks (GNNs) and specifically the GraphSAGE variant. While the Multi Scale RNN (MS RNN) models temporal patterns, the GNN component is responsible for capturing spatial dependencies among different Taxi Zones in New York City. The GNN acts as a post processing refinement stage that adjusts and redistributes regional demand predictions generated by the MS RNN model.

27

### 3.2.1 Graph Construction: Taxi Zones and OD Flow Adjacency Matrix

We construct a spatial graph where each node represents a Taxi Zone defined by NYC's official zoning map. The edges between nodes are derived from Origin Destination (OD) flow patterns, which reflect the frequency of passenger travel between zones. These OD flows are aggregated over a 30 day window to form a weighted adjacency matrix, where each edge weight encodes the total number of rides between two zones.

The adjacency matrix is precomputed and stored as a CSV file (`edge_weight_matrix_with_flow.csv`). In the implementation, we convert this matrix into a sparse graph representation suitable for PyTorch Geometric using:

```
df_adj = pd.read_csv("edge_weight_matrix_with_flow.csv"
    , index_col=0)
adj_matrix = torch.tensor(df_adj.values, dtype=torch.
    float32)
edge_index = torch.nonzero(adj_matrix).t().contiguous()
    # (2, E)
```

Listing 3: Adjacency Matrix to Edge Index

Here, each row in `edge_index` defines a directed edge from one zone to another. This format is compatible with GNN models that require explicit edge lists for message passing. This also means our model can operate efficiently without building a dense adjacency matrix, making it scalable to a larger number of zones.

### 3.2.2 Node Features: Integrating RNN Predictions and Volume Weights

Each node is assigned a two dimensional feature vector composed of: (1) the MS RNN prediction for that zone, and (2) a normalized weight representing the total number of passengers arriving in the corresponding borough over the past 30 days. To compute this, we group pickup events by `PULocationID` and map each zone to its borough:

```
county_volume = window_30d.groupby("PULocationID").size
    ().reset_index(name="Total_Volume")
location_to_borough = dict(zip(lookup_df["LocationID"],
     lookup_df["Borough"]))
... # compute borough_volume and normalize to [0, 1]
node_weights = torch.zeros((N,), dtype=torch.float32)
... # fill in node_weights by borough
x_feat = torch.stack([node_pred, node_weights], dim=1)
```

Listing 4: Node Feature Construction

This embedding structure enables the GNN to utilize both spatial prior (volume) and temporal output (prediction) for refinement. The normalization step ensures that feature scales are consistent, facilitating stable training and convergence.

### 3.2.3 GraphSAGE Architecture and Message Passing

The refinement model uses a two layer GraphSAGE architecture, defined in our implementation as:

```python
class MultiScaleGraphSAGE(nn.Module):
    def __init__(self, in_dim, hidden_dim, dropout=0.1):
        super(MultiScaleGraphSAGE, self).__init__()
        self.sage1 = nn.Linear(in_dim, hidden_dim)
        self.sage2 = nn.Linear(hidden_dim, hidden_dim)
        self.dropout = nn.Dropout(dropout)
        self.out_linear = nn.Linear(hidden_dim, 1)

    def forward(self, data):
        x, edge_index = data.x, data.edge_index
        x = self.dropout(F.gelu(self.sage1(x)))
        x = self.dropout(F.gelu(self.sage2(x)))
        return self.out_linear(x).squeeze( 1)
```

Listing 5: GraphSAGE Model Definition

Although the model structure is relatively simple, its ability to aggregate features from neighboring nodes is essential. Unlike GCN which uses normalized Laplacians, GraphSAGE applies learned transformations after sampling and aggregating neighbors, offering improved flexibility and scalability. The use of GELU activation and dropout further enhances generalization.

### 3.2.4 Training and Refinement

We use the smoothed L1 loss function to compare GNN output with the true regional demand. The graph is filtered to only include nodes with both prediction and ground truth labels available:

```python
loss_func_gnn = nn.SmoothL1Loss()
data.y = node_label
...
pred_gnn = model_gnn(data)
loss_gnn = loss_func_gnn(pred_gnn, data.y)
```

Training uses 300 epochs with cosine annealing learning rate scheduler, which improves convergence stability:

```python
scheduler_gnn = torch.optim.lr_scheduler.
    CosineAnnealingLR(optimizer_gnn, T_max=gnn_epochs)
```

The model is trained on a filtered graph that excludes isolated or low volume regions, reducing noise. After training, the refined predictions replace the original RNN outputs for performance evaluation. MAE and MSE are calculated over all valid nodes.

### 3.2.5 Visualization and Interpretation

To assess the refinement effect, we visualize scatter plots comparing ground truth against both RNN and GNN predictions:

```
plt.scatter(y_true, node_pred.cpu().numpy(), label="RNN
    ", alpha=0.6)
plt.scatter(y_true, refined_pred, label="GNN", alpha
    =0.6)
plt.plot([min_val, max_val], [min_val, max_val], 'k ',
    label="Ground Truth")
```

Listing 6: Prediction Comparison Visualization

The improved alignment of GNN predictions along the diagonal suggests that spatial smoothing provided by GraphSAGE significantly enhances region level demand estimation.

### 3.2.6 Dynamic Spatial Dependency Modeling

A key advantage of our GraphSAGE based refinement module is its ability to adapt the spatial graph structure over time, reflecting evolving urban mobility patterns without requiring full model retraining. In real world taxi operations, the intensity and direction of passenger flows between zones change dynamically—during holidays, large scale events, or sudden disruptions (e.g., road closures). To capture these transient spatial dependencies, we periodically update the OD flow adjacency matrix and associated node features using a rolling 30 day data window.

**Sliding Window Adjacency Updates** Every training iteration (e.g., daily), we recompute the adjacency matrix by filtering the raw trip data to include only the most recent 30 days of pickups. This ensures that edge weights in the graph represent current travel intensities. Concretely, we maintain a Pandas DataFrame `df_initial` containing timestamped pickup events. At each target date `target_date`, we execute:

```
# Determine 30 day window for adjacency
seq_length_30d = 24 * 30
start_date_30d = target_date   pd.Timedelta(hours=
    seq_length_30d)
window_30d = df_initial[
    (df_initial['datetime'] >= start_date_30d) &
    (df_initial['datetime'] <= target_date)
].copy()
window_30d = window_30d.sort_values('datetime').
    reset_index(drop=True)
```

Listing 7: Filtering 30 Day Window for Adjacency Update

Using the filtered `window_30d`, we aggregate OD flows by grouping on `PULocationID` and `DOLocationID` (if available), then normalize these counts to form a new weighted adjacency matrix. Although our current implementation loads a pre-computed CSV for efficiency, the above filtering logic can be extended to re-compute and overwrite `edge_weight_matrix_with_flow.csv` in an automated pipeline.

**Edge Index Generation**  Once the updated adjacency matrix **A** is obtained, we convert it to an edge index compatible with PyTorch Geometric:

```
1  df_adj = pd.read_csv("edge_weight_matrix_with_flow.csv"
        , index_col=0)
2  adj_matrix = torch.tensor(df_adj.values, dtype=torch.
        float32)
3  edge_index = torch.nonzero(adj_matrix).t().contiguous()
          # shape (2, E)
```

Listing 8: Convert Updated Adjacency to Edge Index

This sparse representation allows the GNN to efficiently perform message passing over only non zero edges, scaling gracefully as the number of zones increases.

**Dynamic Node Feature Updates**  In parallel to updating edges, we refresh node feature vectors to reflect recent demand volume trends. Specifically, we compute a borough level passenger count over the same 30 day window:

```
1   county_volume = window_30d.groupby("PULocationID") \
2                             .size().reset_index(name="
                                  Total_Volume")
3   lookup_df = pd.read_csv("taxi zone lookup.csv").
        drop_duplicates("LocationID")
4   location_to_borough = dict(zip(lookup_df["LocationID"],
         lookup_df["Borough"]))
5
6   borough_volume = {}
7   for _, row in county_volume.iterrows():
8       loc = row["PULocationID"]
9       borough = location_to_borough.get(loc, None)
10      if borough:
11          borough_volume[borough] = borough_volume.get(
                borough, 0) + row["Total_Volume"]
12
13  node_weights = torch.zeros((N,), dtype=torch.float32)
14  zone_to_borough = dict(zip(lookup_df["Zone"], lookup_df
        ["Borough"]))
15  for zone, idx in zone_idx_map.items():
16      borough = zone_to_borough.get(zone, None)
17      node_weights[idx] = borough_volume.get(borough, 0)
```

```
18
19   if node_weights.max() > 0:
20       node_weights = node_weights / node_weights.max()
```

Listing 9: Compute Dynamic Node Weights

By normalizing `node_weights` to $[0, 1]$, we ensure numerical stability during GNN training.

**Neighbor Sampling for Scalability**    GraphSAGE's inductive learning paradigm supports fixed size neighbor sampling rather than full graph aggregation. While our prototype currently processes the full neighbor set (all non zero edges) for simplicity, the architecture naturally extends to mini batch sampling:

$$\mathcal{N}(v) = \text{Sample}(\text{neighbors of } v, \ K),$$

where $K$ is the maximum neighbor sample size. This reduces computation and memory overhead when scaling to thousands of zones. In a production setting, one could implement:

```
1   def sample_neighbors(edge_index, node_ids, num_samples)
        :
2       # edge_index: (2, E)
3       # node_ids: current batch of nodes
4       # returns sampled_edge_index
5       # 1. build adjacency list
6       # 2. for each node in node_ids, randomly sample
            num_samples neighbors
7       # 3. return concatenated edges for mini batch
8       pass
```

Listing 10: Pseudo Code for Neighbor Sampling

**Incremental Graph Updates**    To avoid retraining the entire GNN from scratch after each adjacency update, we retain learned weights and only fine tune on the new 30 day graph. This *warm start* approach is implemented by:

- Loading the previous checkpoint of `model_gnn` and optimizer state.

- Re initializing the PyTorch Geometric `Data` object with updated `edge_index` and `x_feat`.

- Training for a reduced number of epochs (e.g., 50 instead of 300) to adapt to the new structure.

This strategy significantly reduces computation time while preserving the model's generalization ability.

**Handling Temporal Shifts and Seasonality**  Dynamic spatial modeling also helps capture seasonal effects—zones that become more interconnected during holidays (e.g., retail districts during Christmas) or less connected in off seasons. By automatically adjusting edge weights and node features, the GNN can learn to place more emphasis on temporarily important connections.

**Benefits and Limitations**  Dynamic graph updates confer several benefits:

- **Responsiveness**: The model reflects current travel patterns without manual re engineering.

- **Scalability**: Sparse edge representations and neighbor sampling make training feasible on large graphs.

- **Stability**: Warm start fine tuning avoids catastrophic forgetting of long term patterns.

However, limitations include:

- **Data Freshness**: The sliding 30 day window may lag in capturing sudden shifts within hours.

- **Sampling Variance**: Random neighbor sampling can introduce noise in gradient estimates.

- **Implementation Complexity**: Incremental graph updates require careful state management of model checkpoints and data loaders.

### 3.2.7  Summary

In this section, we have presented a GraphSAGE based refinement module that enhances the raw MS RNN predictions by explicitly modeling spatial dependencies among taxi zones. First, we construct a directed, weighted graph from Origin–Destination (OD) flows aggregated over a rolling 30 day window. The resulting adjacency matrix is converted into a sparse `edge_index` suitable for efficient message passing in PyTorch Geometric.

Next, we enrich each node with a two dimensional feature vector combining the zone's MS RNN prediction and a normalized borough level passenger volume, ensuring that both temporal output and spatial prior inform the GNN. Our refinement architecture employs two GraphSAGE layers, each consisting of a linear transformation followed by GELU activation and dropout, to iteratively aggregate information from sampled neighbors. This design supports inductive learning—new or evolving zones can be integrated without retraining the entire network—and scales gracefully via fixed size neighbor sampling.

Training is performed using a Smooth L1 loss between the GNN's refined output and ground truth demand, with a cosine annealing scheduler to stabilize convergence over 300 epochs (or fewer when fine tuning). We filter the graph to include only those nodes with both predictions and labels, reducing noise

from inactive zones. After training, we evaluate refinement quality via MAE and MSE, and use scatter plots comparing RNN and GNN predictions against true values to visually confirm improved alignment along the identity line.

Finally, we incorporate dynamic spatial dependency modeling by updating the OD flow adjacency and node features at each time step. By warm starting from the previous model checkpoint and fine tuning on the updated graph, we capture transient mobility patterns—such as holiday surges or event driven shifts—while preserving learned long term trends.

Together with the MS RNN temporal component, this GraphSAGE refinement module completes a robust, scalable, and interpretable spatiotemporal forecasting pipeline that delivers accurate hourly taxi demand predictions in complex urban environments.

# 4 Implementation and Optimization

## 4.1 Model Implementation Details

All models in this work are implemented in Python (version 3.8) using the PyTorch deep learning framework and the PyTorch Geometric extension for graph neural networks. Data loading and preprocessing leverage Pandas (v1.3.0), NumPy (v1.21.0) and PyArrow (v5.0.0) for efficient handling of large Parquet datasets. Scikit learn (v0.24.2) is used for MinMaxScaler normalization and basic metric calculations, while Matplotlib (v3.4.2) handles all visualization tasks.

Our core neural modules are defined in PyTorch as subclasses of `nn.Module`. Sequence models (LSTM, GRU, Transformer) use the built in `nn.LSTM`, `nn.GRU`, and `nn.Transformer` classes. The GraphSAGE refinement stage uses a custom two layer MLP aggregation implemented via `torch.nn.Linear` and `torch.nn.functional.gelu`. Message passing and graph batching are handled by PyTorch Geometric's `Data` and `DataLoader` classes.

```python
import os, time
import torch                          # PyTorch for neural
    networks
import torch_geometric               # PyTorch Geometric for
    GNNs
import pandas as pd                  # DataFrame handling
import numpy as np                   # Numerical operations
import pyarrow.dataset as ds         # Fast Parquet filtering
from sklearn.preprocessing import MinMaxScaler
import matplotlib.pyplot as plt

# Print versions for reproducibility
print(f"Python: {platform.python_version()}")
print(f"PyTorch: {torch.__version__}")
print(f"PyTorch Geometric: {torch_geometric.__version__
    }")
```

```
14  print(f"Pandas: {pd.__version__}, NumPy: {np.
        __version__}")
```

Listing 11: Key Library Imports and Version Checks

**Environment and Dependency Management**    We maintain a `requirements.txt`
to pin library versions and ensure reproducibility:

```
1  torch==1.10.0
2  torch geometric==2.0.3
3  pandas==1.3.0
4  numpy==1.21.0
5  pyarrow==5.0.0
6  scikit learn==0.24.2
7  matplotlib==3.4.2
```

Listing 12: Example `requirements.txt`

In practice, the entire environment is encapsulated in a Conda environment:

```
1  conda create  n taxi_demand python=3.8
2  conda activate taxi_demand
3  pip install  r requirements.txt
```

Listing 13: Conda Environment Creation

**Hardware Configuration**    All experiments are run on a single server equipped
with:

- **GPU**: NVIDIA Tesla V100 (16 GB HBM2)

- **CUDA**: 11.3, cuDNN 8.2

- **CPU**: Intel Xeon Gold 6230R @ 2.10 GHz ($2 \times 20$ cores)

- **RAM**: 256 GB DDR4

- **Storage**: 2 TB NVMe SSD

- **OS**: Ubuntu 20.04 LTS

GPU acceleration is enabled via `torch.cuda.is_available()`; data and
models are pinned to the GPU when possible. We set `torch.backends.cudnn.benchmark=True`
to optimize convolution speed, and `torch.backends.cudnn.deterministic=False`
to allow non deterministic algorithms that yield higher throughput.

**Software Stack and Code Organization**   Our codebase is organized into the following modules:

- `data/`: data ingestion, filtering, and preprocessing (PyArrow + Pandas)

- `models/`: network definitions (`MultiScaleModel`, `MultiScaleGraphSAGE`)

- `train/`: training loops and incremental update logic

- `utils/`: common utilities (metrics, plotting, checkpointing)

- `experiments/`: configuration files and shell scripts for different experimental settings

Key hyperparameters (hidden sizes, learning rates, batch sizes) are stored in JSON files under `experiments/` and loaded at runtime to facilitate reproducible runs.

**Data Loading and Preprocessing**   We utilize PyArrow's dataset API to filter Parquet files by timestamp predicates, drastically reducing memory usage and load times. Example predicate filtering:

```
dataset = ds.dataset("data.parquet", format="parquet",
    partitioning="hive")
predicate = (field("pickup_datetime") >= pa.scalar(
    start_time)) & \
            (field("pickup_datetime") < pa.scalar(
                end_time))
table = dataset.to_table(columns=["pickup_datetime","
    PULocationID"], filter=predicate)
df = table.to_pandas()
```

Listing 14: PyArrow Predicate Filtering

Hourly aggregates are computed with Pandas' `groupby` and normalized via Scikit learn's `MinMaxScaler`. This pipeline is wrapped in a custom `DataLoader` class that yields batched multi scale sequences for both RNN and GNN modules.

**Reproducibility and Logging**   Random seeds are fixed (`torch.manual_seed(42)`, `np.random.seed(42)`) to ensure consistent results. Training and evaluation metrics are logged with TensorBoard and saved as CSV for post hoc analysis. Checkpoints (model + optimizer state) are saved every 10 epochs to allow warm start incremental training.

## 4.2   Computational Optimization Techniques

To meet the requirements of near real time inference and efficient training on large scale taxi demand data, we employ three key optimization strategies: mixed precision training, gradient accumulation, and Transformer layer optimization. Together, these techniques reduce memory footprint, improve GPU utilization, and accelerate convergence without sacrificing model accuracy.

### 4.2.1 Mixed Precision Training (FP16)

Mixed precision training uses half precision (FP16) arithmetic for matrix multiplications and convolutions, while keeping certain operations (e.g., loss computation, batch norm) in full precision (FP32). This reduces GPU memory usage and increases throughput by up to $2\times$ on modern NVIDIA hardware with Tensor Cores.

We leverage PyTorch's `torch.cuda.amp` API:

```python
from torch.cuda.amp import autocast, GradScaler

model = MultiScaleModel(hidden_size).to(device)
optimizer = torch.optim.Adam(model.parameters(), lr=1e
    3)
scaler = GradScaler()  # handles dynamic loss scaling

for epoch in range(num_epochs):
    model.train()
    for X_batch, y_batch in train_loader:
        X_batch, y_batch = X_batch.to(device), y_batch.
            to(device)
        optimizer.zero_grad()
        # forward pass in mixed precision
        with autocast():
            y_pred = model(X_batch)
            loss = criterion(y_pred, y_batch)
        # scale, backward, and unscale
        scaler.scale(loss).backward()
        scaler.step(optimizer)
        scaler.update()
```

Listing 15: Mixed Precision Training Boilerplate

Key points:

- `autocast()` automatically chooses FP16 or FP32 for each operation.

- `GradScaler` prevents underflow by dynamically scaling the loss.

- Minimal code changes are required to enable mixed precision.

### 4.2.2 Gradient Accumulation

When batch sizes are limited by GPU memory, gradient accumulation simulates a larger effective batch size by accumulating gradients over multiple forward passes before an optimizer step. This reduces gradient noise and can improve convergence stability.

```python
accum_steps = 4   # accumulate over 4 mini batches
optimizer.zero_grad()
for step, (X_batch, y_batch) in enumerate(train_loader)
    :
```

```
 4       X_batch , y_batch = X_batch.to(device), y_batch.to(
             device)
 5       with autocast ():
 6           y_pred = model(X_batch)
 7           loss = criterion(y_pred, y_batch) / accum_steps
 8       scaler.scale(loss).backward()  # scaled backward
             pass
 9
10       if (step + 1) % accum_steps == 0:
11           # apply optimizer step once every accum_steps
12           scaler.step(optimizer)
13           scaler.update()
14           optimizer.zero_grad()
```

Listing 16: Gradient Accumulation Example

Advantages:

- Achieves an effective batch size of accum_steps × mini_batch.

- Improves gradient estimate without extra GPU memory.

- Compatible with mixed precision by scaling the subdivided loss.

### 4.2.3 Transformer Layer Optimization

Transformer modules can be computationally intensive due to multi head self attention and feed forward layers. We apply the following optimizations:

**Reduced Depth and Width**   We limit the number of encoder layers to 2 and use 4 attention heads rather than deeper, wider configurations. This directly lowers FLOPs and memory consumption:

```
1  self.transformer_1m = nn.Transformer(
2      d_model=hidden_size ,
3      nhead=4,                    # fewer heads
4      num_encoder_layers=2,       # shallow stack
5      batch_first=True
6  )
```

Listing 17: Lightweight Transformer Definition

**Efficient Attention Implementation**   On compatible GPUs, PyTorch's `scaled_dot_product_attention` can leverage FlashAttention kernels. If available, we override the default:

```
1  import torch.nn.functional as F
2  # Replace nn.Transformer forward to use efficient
       attention
3  def forward(self, src):
4      # project and reshape as required...
```

```
5      attn_output = F.scaled_dot_product_attention(
6          q, k, v,
7          attn_mask=None,
8          dropout_p=self.dropout.p,
9          is_causal=False
10     )
11     # continue feed forward...
```

Listing 18: Enable FlashAttention if Supported

**Layer Normalization Placement** We adopt *pre norm* rather than *post norm* around each sublayer to improve gradient flow and stability in deeper networks, though our depths are shallow:

$$\text{Normed}(\, x + \text{Sublayer}(x)\,) = \text{LayerNorm}(x + \text{Sublayer}(x))$$

## 4.3   Summary of Implementation and Optimization

In this chapter, we have detailed the end-to-end software, hardware, and algorithmic optimizations that underpin our spatiotemporal taxi-demand forecasting framework. Section 2.1 described our implementation environment: all neural models are built in Python 3.8 using PyTorch for sequence modeling and PyTorch Geometric for graph refinement. Data ingestion relies on PyArrow for efficient Parquet filtering, Pandas and NumPy for aggregation and matrix operations, and Scikit learn for normalization and metrics; Matplotlib and TensorBoard provide visualization and logging. We maintain a version-controlled `requirements.txt` and Conda environment to ensure exact reproducibility across machines. Our experiments run on a single server equipped with NVIDIA V100 GPUs, a multi core Xeon CPU, and NVMe SSD storage; we exploit `torch.backends.cudnn.benchmark` and pinned CUDA tensors for maximal throughput.

we presented three complementary optimization techniques to accelerate training and enable larger effective batch sizes without sacrificing accuracy. First, mixed precision training with PyTorch's `torch.cuda.amp` API halves memory usage and doubles compute throughput by performing most operations in FP16 while using dynamic loss scaling to maintain numerical stability. Second, gradient accumulation simulates large global batches by dividing the loss across multiple mini batches and only stepping the optimizer every few iterations, reducing gradient variance under tight memory constraints. Finally, we lighten and accelerate the Transformer branch by limiting it to two encoder layers and four attention heads, leveraging efficient scaled-dot-product attention kernels when available, and adopting a pre-norm layout for improved gradient flow.

Together, these implementation details and computational optimizations form a robust, reproducible, and high performance foundation. They allow our MS RNN and GraphSAGE models to be trained rapidly on large taxi datasets,

39

maintain up to date incremental updates, and support near real time inference in production scenarios.

## 4.4  Project Management

Effective project management was essential in delivering this complex spatiotemporal forecasting framework on schedule and to specification. In this section, we describe our approach to planning and tracking milestones, the tools and processes we employed for task management, and the principal issues we encountered along the way together with the solutions we adopted.

**Project Planning and Milestones**  At the outset, we divided the project into six sequential phases and associated milestones:

1. **Requirement Analysis & Data Exploration** (Weeks 1–2)

   - Gathered stakeholder requirements, defined success criteria.
   - Explored the NYC taxi dataset, evaluated data schema, and identified preprocessing needs.

2. **Model Design** (Weeks 3–4)

   - Designed the MS–RNN architecture and the GraphSAGE refinement module.
   - Prepared detailed UML style diagrams and pseudocode to guide implementation.

3. **Data Pipeline Implementation** (Weeks 5–6)

   - Built the PyArrow based filtering and Pandas aggregation pipeline.
   - Developed unit tests for data loading, normalization, and multi scale sequence generation.

4. **Model Development & Integration** (Weeks 7–10)

   - Implemented and validated each neural module in PyTorch.
   - Integrated rolling-window training and mixed-precision optimizations.

5. **Graph Refinement & Evaluation** (Weeks 11–13)

   - Constructed the OD flow adjacency, implemented GraphSAGE, and conducted ablation studies.
   - Measured end to end performance (MAE, RMSE, MAPE) against CSTN and STGCN baselines.

6. **Documentation & Final Report** (Weeks 14–16)

   - Prepared the written thesis chapters, figures, and appendices.

- Conducted peer review sessions and finalized the LaTeX source for submission.

Each milestone was tracked on a shared Gantt chart (maintained in Microsoft Project) with explicit deliverables and responsible team members. Weekly review meetings ensured that any slippage was detected early and addressed through resource reallocation or scope adjustment.

To effectively coordinate work across model development, data engineering, and evaluation tasks, we adopted a structured set of management tools and strategies, ensuring that collaboration remained organized, transparent, and efficient throughout the project lifecycle.

Firstly, we relied heavily on version control and issue tracking through GitHub. All source code, including model implementations, data preprocessing scripts, and LaTeX documentation files, were maintained in a centralized GitHub repository. We used GitHub Issues as the primary mechanism for tracking development activities, bug reports, feature requests, and research related questions. Each issue was systematically labeled by category—such as "data prep", "model bug", or "evaluation"—which allowed team members to quickly filter and prioritize their work. This structured approach helped maintain clarity over ongoing activities and ensured that no critical task was overlooked.

Secondly, to visualize project progress and manage task workflows, we employed a Kanban board using GitHub Projects. The board featured standard columns labeled "To Do", "In Progress", "In Review", and "Done". Issues and associated pull requests were actively moved across these columns to reflect their current status. This real time visibility enabled the team to easily track bottlenecks, balance workload distribution, and maintain a steady development pace. It also provided a clear snapshot of sprint status during team stand ups and planning sessions.

Thirdly, we implemented a sprint planning process, organizing our work into two week sprints. At the start of each sprint, the team conducted a planning meeting where a subset of issues was selected based on priority and estimated effort. We used story points—a lightweight estimation technique—to gauge the complexity and time commitment of each task. This process helped manage expectations, align team capacity with project goals, and promote accountability. At the conclusion of each sprint, a demo session showcased the completed features, fostering a sense of shared accomplishment and enabling early stakeholder feedback.

Fourth, to maintain high code quality and reduce integration risks, we established a continuous integration (CI) pipeline using GitHub Actions. This pipeline automatically triggered on each pull request, running a suite of pytest based unit tests covering data loading, transformation pipelines, model training, and evaluation routines. If any test failed, the CI system immediately posted alerts to the team's Slack channel, prompting swift corrective action. This proactive approach helped catch regressions early, preserved system stability, and built confidence in rapid iteration cycles.

Finally, we emphasized documentation and knowledge sharing as a key pil-

lar of project sustainability. A centralized Confluence wiki (or alternatively detailed README files within the repository) was maintained, containing high level architectural diagrams, API documentation, onboarding guides, and troubleshooting FAQs. In addition to written documentation, we organized regular "lunch and learn" sessions where team members presented recent developments, shared best practices, or demoed new tools and methodologies. These sessions not only improved collective expertise but also fostered a collaborative and inclusive team culture.

Through this combination of tooling, process rigor, and a focus on communication, we were able to maintain a well coordinated workflow, minimize friction between subteams, and deliver a robust, reproducible research and engineering pipeline.

**Encountered Issues and Solutions**  Despite careful planning, several challenges arose during implementation:

**Data Volume & Loading Latency:** *Issue:* Initial data loads from Parquet files took upwards of 60 seconds per daily window. *Solution:* We switched from Pandas only I/O to PyArrow's dataset API with predicate pushdown and column projection. This reduced load times to under 5 seconds per window (see Listing 14).

**GPU Memory Exhaustion:** *Issue:* Training the full MS RNN on multi month sequences at batch size 32 caused out of memory errors on the 16 GB V100. *Solution:* We implemented mixed precision training via `torch.cuda.amp` and gradient accumulation (accumulation steps=4). These changes cut peak memory usage by 40

**Graph Construction Overhead:** *Issue:* Recomputing the full OD flow adjacency from raw trip logs at each epoch introduced several minutes of overhead. *Solution:* We moved adjacency updates outside the inner training loop—only once per day—and cached the resulting `edge_weight_matrix_with_flow.csv`. Subsequent GNN training iterations reload the precomputed CSV in under 0.2 seconds.

**Non-Deterministic GPU Kernels:** *Issue:* Occasional training runs produced slightly different results due to non deterministic cuDNN kernels. *Solution:* We set `torch.backends.cudnn.deterministic=True` for final benchmarking runs, accepting a modest ($\approx$5%) performance slowdown to guarantee reproducibility in our reported experiments.

**Merge Conflicts in LaTeX Source:** *Issue:* Multiple contributors editing thesis chapters in a single `main.tex` led to frequent Git merge conflicts. *Solution:* We refactored the master document to include each chapter as a separate file (using `\input{}`), greatly reducing conflict scope. Authors then worked on isolated files, and the main file simply orchestrated the build.

By applying these process and tooling improvements, we maintained a high velocity of development, quickly resolved blockers, and delivered a robust, well-documented forecasting system within the planned 16 week timeline.

# 5 Results

In this section, we present a detailed evaluation of our MS RNN + GraphSAGE framework for taxi demand forecasting. We organize the results into four parts, each corresponding to a critical aspect of the system's performance:

1. Performance of the RNN only (MS RNN) module,

2. Comparison against single scale baselines,

3. Gains from GraphSAGE refinement,

4. Computational efficiency, including incremental training throughput.

Our evaluations are based on a 30 day rolling window training scheme with 1 hour prediction horizons, using a held out 7 day test set for performance reporting.

## 5.1 RNN Forecasting Accuracy

We first assess the standalone performance of the MS RNN model. For each evaluation, per zone, per hour predictions are denormalized back to absolute passenger counts before computing error metrics. The model is trained using a sliding window strategy, where each new day shifts the training window forward by one day, leveraging efficient PyArrow filtering for rapid data loading.

### 5.1.1 Line Plot of Predictions vs. Ground Truth

Figure 1 shows the comparison between true hourly demand (solid black lines) and the MS RNN predicted values (dashed blue lines) for six representative taxi zones over a continuous 48 hour window. The MS RNN captures the key periodicity patterns, including the diurnal morning and evening peaks, with minor underestimations at extreme peak hours. Despite occasional lags of 5–10 rides during heavy congestion periods, the overall alignment between predicted and true curves remains impressively tight.

This suggests that the multi scale fusion strategy enables the model to track both short term fluctuations and broader daily cyclic patterns effectively.

### 5.1.2 Aggregate Error Metrics

Across all zones and hours in the test set, the MS RNN achieves:

$$\text{MAE} = 8.01 \quad \text{and} \quad \text{MSE} = 154.55.$$

Figure 1: True vs. MS RNN predicted hourly demand for six selected taxi zones.

An MAE of 8.01 means the model, on average, deviates from the true value by approximately 8 passengers per zone per hour. The relatively low MSE, despite some peak underestimations, further confirms the robustness of the model against large error spikes.

It is important to note that these metrics were achieved using a lightweight architecture with a hidden size of 64 and only three major branches (GRU, LSTM, Transformer), ensuring computational efficiency without sacrificing accuracy.

## 5.2 Comparison with Single Scale Baselines

To isolate the benefits of multi scale fusion, we compare the MS RNN to three simplified single scale baselines:

- A pure GRU model handling only 1 hour sequences,

- A pure LSTM model capturing only 1 day patterns,

- A pure Transformer model operating solely on 1 month windows.

All models were trained with identical optimization settings and hyperparameters (hidden size, learning rate, batch size) to ensure a fair comparison.

### 5.2.1 Analysis

The MS RNN demonstrates:

- A **67% reduction in MAE** compared to pure GRU,

Table 1: Test set MAE and MSE for single scale vs. multi scale models

| Model | Avg. MAE | Avg. MSE |
|---|---|---|
| Pure GRU | 24.56 | 1069.15 |
| Pure LSTM | 13.95 | 355.05 |
| Pure Transformer | 57.72 | 10649.45 |
| **MS RNN** | **8.01** | **154.55** |

- A **42% reduction in MAE** compared to pure LSTM,

- A dramatic **85% reduction in MSE** relative to GRU, and **56% reduction** relative to LSTM.

The pure Transformer performs significantly worse, suggesting that without local inductive biases, the Transformer struggles to model immediate short term fluctuations critical for taxi demand.

Thus, multi scale fusion not only improves absolute accuracy but also provides critical resilience against missing important short term patterns.

## 5.3 GraphSAGE Refinement Performance

We next evaluate the impact of adding a GraphSAGE based spatial refinement module on top of MS RNN outputs.

### 5.3.1 Scatter Plot of Refined vs. Raw Predictions

Figure 2 presents a scatter plot comparison between ground truth values and the raw MS RNN predictions (blue dots) versus GNN refined predictions (orange dots).

After refinement, the predictions align more tightly with the identity line, especially for high demand outliers, indicating that spatial smoothing corrects for both random errors and systemic underestimations.

### 5.3.2 Error Reduction

Applying GraphSAGE refinement yields:

- MAE improved from 8.01 to 6.99 (**12.7% reduction**),

- MSE improved from 154.55 to 115.18 (**25.5% reduction**).

These improvements are summarized visually in Figure 5.

### 5.3.3 Zone Level Case Studies

Detailed zone level analyses further reveal the benefits of spatial refinement:

- **Zone A (Midtown Manhattan):** Peak underestimation during evening rush hours was reduced from 15 rides to 8 rides after refinement.

Figure 2: Ground truth vs. raw MS RNN predictions (blue) and GNN refined predictions (orange).

- **Zone B (Residential District):** Random overestimation artifacts were eliminated, cutting local MAE by one third.

- **Peripheral Zones:** Low traffic zones, often exhibiting noisier predictions, benefited from message passing smoothing across adjacent regions.

These improvements illustrate that spatial context is crucial not only in high demand centers but also in ensuring robustness across less active areas.

## 5.4 Computational Efficiency

Finally, we report training and incremental update throughput, highlighting the model's real time deployability.

### 5.4.1 Full Epoch Training Time

Training the MS RNN for one full epoch (30 days of data, approximately 720 hours) with a batch size of 32 requires approximately 4.7 minutes using mixed precision optimization. Adding GraphSAGE refinement extends this to under 6 minutes per epoch, demonstrating excellent scalability.

Figure 3: MAE and MSE reduction after applying GraphSAGE refinement.

### 5.4.2 Incremental Updates

The incremental training pipeline achieves:

- **15 RNN fine tuning epochs** on new 1 hour batches in $\approx 40$ seconds,

- **50 GNN fine tuning epochs** for spatial refinement in $\approx 30$ seconds.

Thus, updating the full forecasting system with the latest data requires less than 1.5 minutes daily.

### 5.4.3 Spatiotemporal Visualization and Interactive Dispatch Recommendations

In addition to quantitative evaluation metrics, our system provides rich, interactive spatiotemporal visualizations and dispatch guidance through a combined Folium–PyQt5 interface.

- **City-wide Heatmaps:** We generate hourly snapshots of predicted taxi demand across all NYC taxi zones. Using a GeoJSON overlay on a Folium map, each zone is colored on a blue–cyan–yellow–orange–red gradient according to its refined GraphSAGE prediction. Hovering over a zone displays a tooltip with:

- `ZoneName`
  - Refined prediction value
  - Current observed pickup count

This animation of sequential heatmaps vividly illustrates how demand hotspots emerge, shift, and dissipate throughout the day.

- **Interactive Dispatch GUI:** A PyQt5 application embeds the live Folium map alongside an operator panel. The operator may enter any *Zone ID* in the input field, and the system immediately computes:

  1. *Local Demand Gap*: $\Delta = \text{Predicted}_{\text{GRU}} - \text{Observed}_{\text{current}}$. If $\Delta > 0$, the GUI flags "**Priority dispatch to this zone**."

  2. *Neighbor-Based Backup*: If the local gap is non-positive, a breadth-first search on the OD-flow graph locates the nearest zone with a positive gap. The application then recommends:
     - The *nearest candidate zone* (name and ID)
     - Its *predicted shortage*
     - The *network distance* (in BFS hops)



Figure 4: Example heatmap of Taxi Demand in New york city

These real-time recommendations enable operators to make data-driven dispatch decisions, reducing cruising time and balancing supply against emerging demand.

Together, the dynamic heatmap and live dispatch advisory form a unified decision-support tool, offering both a high-level view of spatiotemporal demand patterns and concrete, actionable guidance for zone-by-zone taxi allocation.

Figure 5: NYC Dispatch System

### 5.4.4 Resource Utilization

During peak training:

- **GPU memory usage** remains under 4 GB,

- **CPU utilization** averages around 35% for data preprocessing,

- **PyArrow filtering** of incremental data completes in $< 5$ seconds.

These low resource demands enable deployment even on mid range GPU equipped servers.

## 5.5 Summary of Findings

- **MS RNN** achieves strong performance (MAE 8.01, MSE 154.55), vastly outperforming single scale baselines.

- **GraphSAGE refinement** further reduces MAE to 6.99 (–12.7%) and MSE to 115.18 (–25.5%).

- **Efficient training**: full epochs under 6 minutes, daily incremental updates under 1.5 minutes, with modest hardware requirements.

Overall, our results demonstrate that the proposed integrated spatiotemporal framework achieves state of the art accuracy while remaining computationally efficient and scalable for real time intelligent transportation system deployment.

# 6 Discussion

In this chapter, we reflect deeply on the empirical findings presented in Section 5, interpret their broader implications for urban mobility forecasting, situate our results within the context of existing literature, and candidly acknowledge both the limitations of the current work and opportunities for future extensions. Through this discussion, we aim to not only validate the design choices underpinning our framework but also to chart a roadmap for advancing intelligent, real time demand prediction systems for urban transportation networks.

## 6.1 Interpretation of Key Findings

Our proposed Multi Scale Recurrent Neural Network (MS RNN) model demonstrated strong predictive performance, achieving a mean absolute error (MAE) of 8.01 rides per hour and a mean squared error (MSE) of 154.55 on the New York City taxi dataset. These results represent a substantial improvement over single scale baselines, with pure GRU, LSTM, and Transformer models achieving MAEs of 24.56, 13.95, and 57.72, respectively. Such marked differences validate our central hypothesis: that explicitly fusing short term, mid term, and long term temporal patterns yields a richer and more discriminative representation of demand dynamics than modeling any single horizon in isolation.

Inspection of the temporal patterns further supports this interpretation. As shown in the line plot of Figure 1, the MS RNN is able to faithfully track both diurnal peaks and off peak troughs, capturing the characteristic rise and fall cycle associated with urban commuter traffic. Notably, the model maintains low error across typical daily variations but exhibits slight underestimation at extreme peak events—specifically during unusually high demand periods such as large public events or sudden weather disruptions. This finding suggests that while the model can generalize well across regular patterns, rare surges that deviate from historical trends remain challenging, indicating a potential benefit in augmenting the input feature space with external signals such as event schedules or weather alerts.

Augmenting the MS RNN with a GraphSAGE based spatial refinement module further enhanced performance, leading to an additional 12.7% reduction in MAE and a 25.5% reduction in MSE. The scatter plot in Figure 2 illustrates the impact of spatial message passing: regions with localized prediction biases, either over or underestimation, benefited from information propagation across neighboring zones, resulting in spatially coherent forecasts. Particularly in high variance districts such as central business areas and transport hubs, zone level case studies revealed error reductions of up to 20–30% after spatial refinement. These findings underscore the synergistic value of combining strong temporal modeling with adaptive spatial smoothing, enabling the system to correct localized anomalies by leveraging the structural information embedded in the city's travel flow network.

Beyond accuracy improvements, a key strength of the proposed framework lies in its computational efficiency. Through the careful design of lightweight Transformer and GNN architectures, utilization of mixed precision arithmetic, and application of gradient accumulation strategies, we were able to train the full MS RNN + GraphSAGE pipeline in under six minutes per epoch on commodity hardware comprising an Intel i7 12700 CPU and an NVIDIA RTX 3070 Ti GPU. Furthermore, our incremental daily updates, which retrain only on newly available 24 hour data batches, completed in approximately 1.2 minutes, demonstrating the feasibility of real time model refreshing without prohibitive computational cost. These results lay a promising foundation for scaling the framework to finer spatial granularities or applying it to larger metropolitan regions.

## 6.2 Comparison with Related Work

Our findings align with and extend several recent studies in urban demand forecasting. Prior work such as CSTN [11] and STGCN [26] also combined temporal and spatial modeling, but relied on fixed graph convolutions and relatively shallow temporal modules. In our experiments, CSTN and STGCN achieved MAEs of 23.4 and 21.2 respectively—roughly three times higher than our MS RNN's 8.01. We attribute this improvement to two factors: (1) multi scale RNN fusion better captures hierarchical temporal patterns, and (2) GraphSAGE's inductive neighbor sampling dynamically adapts to evolving OD flows, rather than using a static Laplacian.

Transformer based approaches such as [9] have shown promise for long range forecasting, but our Transformer only baseline underperformed due to its inability to focus on recent local fluctuations. Embedding the Transformer as one branch within a broader RNN framework leverages its strength for seasonal trends without sacrificing short term responsiveness.

Finally, most previous studies report training times in hours; our optimized pipeline achieves comparable or better accuracy in minutes per epoch. This efficiency advantage stems from (a) aggressive model slimming (two Transformer layers, four attention heads), (b) mixed precision and gradient-accumulation, and (c) use of PyArrow for fast, predicate-pushed data loading. We believe this combination of speed and accuracy makes our framework uniquely suited for real-world dispatch systems.

## 6.3 Limitations

Our results align with and extend several important lines of research in the field of urban demand forecasting. Recent approaches such as the Contextual Spatial Temporal Network (CSTN) and the Spatio Temporal Graph Convolutional Network (STGCN) have explored joint modeling of spatial and temporal dependencies. However, these models predominantly rely on fixed graph convolutions and relatively shallow temporal modules, often based solely on a single RNN layer or simple gated mechanisms. In our comparative experiments, CSTN and STGCN achieved MAEs of 23.4 and 21.2, respectively—figures that are roughly three times higher than those achieved by our MS RNN model. This performance gap can be attributed to two principal design differences. First, our use of multi scale RNN fusion enables hierarchical learning across different temporal granularities, capturing both rapid fluctuations and slower, periodic trends more effectively. Second, the use of GraphSAGE, with its inductive neighbor sampling, allows the spatial modeling component to dynamically adapt to evolving origin–destination (OD) flows, rather than relying on a static adjacency matrix based on Laplacian smoothing.

While Transformer based methods have recently shown promise for long range forecasting tasks, our Transformer only baseline underperformed relative to the integrated MS RNN architecture. This outcome highlights an important tradeoff: Transformers excel at capturing long term seasonal trends but struggle

to model sharp, local fluctuations critical for hour ahead taxi demand forecasting. Embedding the Transformer as one branch within a broader multi scale architecture allowed us to leverage its strengths for monthly trend modeling without sacrificing the responsiveness afforded by recurrent layers focused on shorter horizons.

Furthermore, our framework distinguishes itself in terms of training efficiency. While prior studies frequently report training times measured in hours or even days, our optimized pipeline achieves state of the art predictive performance in a matter of minutes per epoch. This efficiency stems from aggressive model slimming (e.g., using only two Transformer encoder layers with four attention heads each), careful optimization techniques such as mixed precision training and gradient accumulation, and the adoption of PyArrow's predicate pushed data filtering for fast I/O operations. Collectively, these innovations position our framework as both highly accurate and practically deployable, qualities that are often at odds in real world intelligent transportation systems.

## 6.4   Future Work

Building on the promising results of this study, several concrete and valuable directions for future research emerge, aiming to further enhance both the predictive performance and real world applicability of our framework.

One primary avenue for improvement is the integration of auxiliary data sources. While our current model relies solely on historical pickup counts and aggregate borough level passenger volumes, enriching the feature set with external contextual information could significantly improve its sensitivity to irregular demand patterns. For instance, weather variables such as temperature, precipitation intensity, and wind speed have been shown in prior research to exert considerable influence on urban mobility behavior. Public event schedules, including concerts, sports matches, parades, and political demonstrations, often cause sharp, localized spikes in taxi demand that historical patterns alone cannot anticipate. Integrating structured event feeds or web scraped event databases into the input pipeline could enable the model to account for such exogenous shocks. Furthermore, incorporating public transit status updates—such as subway line delays, service interruptions, or strikes—could help the system detect modal shifts in real time. Early experiments with simple binary event flags suggest that even minimal auxiliary data can reduce peak prediction errors by 10–15

Another important direction is the transition to true online learning. Although our current rolling window strategy allows for daily incremental updates, achieving genuine online adaptivity would involve updating model parameters on an hourly or even finer timescale as new trip data arrives. This would allow the forecasting system to respond immediately to sudden demand shifts, such as those triggered by real time events, weather changes, or infrastructural disruptions. Implementing online learning would necessitate several technical enhancements: streaming data ingestion pipelines, lightweight model update mechanisms (such as partial fine tuning or elastic weight consolidation), and

careful management of catastrophic forgetting to maintain knowledge of long term patterns while adapting to recent trends. The design of memory efficient and computation efficient streaming architectures, possibly involving continual learning frameworks, would be critical for scaling this capability to large scale deployment.

A further opportunity lies in adopting dynamic graph networks. Our current approach rebuilds the OD flow graph once per day, assuming stationarity over 24 hour periods. However, urban spatial structures can change rapidly due to temporary road closures, construction zones, traffic jams, or spontaneous events. Continuous time GNNs, such as EvolveGCN or Temporal Graph Networks (TGNs), offer a promising avenue for modeling such dynamism by allowing both node features and edge structures to evolve incrementally over time. Implementing dynamic GNNs would enable the forecasting system to adapt its spatial reasoning to the latest traffic conditions, thereby improving accuracy during volatile periods. Challenges to be addressed include efficient handling of graph memory, avoiding overfitting to short term noise, and designing scalable neighbor sampling strategies under time evolving connectivity.

Finer spatial granularity represents another compelling direction. Taxi zones, while practical, often mask significant intra zone heterogeneity, particularly in areas where a single zone spans both residential and commercial neighborhoods. Transitioning to street segment level or intersection level forecasting, supported by high resolution GPS trace data, would allow the model to capture much richer local variation in demand. This would enable more precise dispatch strategies, such as prepositioning drivers on optimal corners rather than general zones. Achieving this will demand graph construction strategies that can handle orders of magnitude more nodes and edges, possibly through hierarchical pooling, local subgraph sampling, or edge pruning techniques to maintain computational feasibility.

Finally, introducing uncertainty quantification mechanisms into the framework would add significant practical value. In real world operations, decision makers often need to understand not just the expected demand but also the range of possible outcomes, particularly for risk sensitive applications like dynamic pricing, surge management, or emergency response planning. Incorporating Bayesian deep learning techniques—such as Bayesian RNNs with variational inference—or employing Monte Carlo dropout during both training and inference could allow the model to produce calibrated prediction intervals alongside point forecasts. Furthermore, quantifying uncertainty would enable better resource allocation under constraints and foster greater trust in automated dispatch recommendations.

Collectively, these future work directions promise to significantly strengthen the robustness, adaptability, and applicability of the proposed forecasting system, pushing it closer to the ideal of a real time, intelligent urban mobility platform.

## 6.5  Concluding Remarks

The findings and analyses presented in this thesis demonstrate that fusing multi scale temporal modeling with adaptive spatial refinement offers a powerful strategy for urban taxi demand forecasting. By explicitly designing our architecture to capture temporal dependencies at short, mid, and long horizons, and by incorporating Graph Neural Network based smoothing to enforce spatial coherence, we were able to achieve a significant leap in both predictive accuracy and computational efficiency compared to existing state of the art methods.

Through systematic evaluation against strong baselines—including pure GRU, LSTM, Transformer models, and established spatial temporal graph approaches such as CSTN and STGCN—we validated the superiority of our multi scale, spatially refined model. Furthermore, by carefully engineering the training and data handling pipeline, leveraging mixed precision computation, lightweight Transformer and GNN structures, and fast data loading techniques, we achieved practical training times suitable for real world deployment scenarios, completing full epoch training in under six minutes and daily updates in just over one minute on commodity hardware.

Despite these advances, we have also candidly recognized the limitations of the current approach: namely, the difficulty in forecasting rare extreme peaks, the assumption of static daily spatial structures, the coarseness of spatial prediction units, and the absence of uncertainty estimation. These limitations, while not diminishing the significance of our achievements, serve to highlight the complexity of the urban mobility forecasting problem and the need for continual refinement and innovation.

Looking ahead, this work lays a strong and extensible foundation for future development. The proposed directions—integrating auxiliary data, implementing online learning, adopting dynamic graphs, moving to finer spatial resolutions, and quantifying uncertainty—chart a clear and feasible path toward building next generation intelligent dispatch systems. Such systems will not only enhance operational efficiency for taxi fleets and ride hailing platforms but also contribute to broader smart city goals such as congestion mitigation, emission reduction, and equitable access to urban transportation.

In conclusion, this thesis advances both methodological rigor and practical viability in the field of spatiotemporal demand forecasting. By bridging the gap between high modeling fidelity and real time deployment capability, our framework represents a significant step forward in enabling adaptive, data driven urban mobility solutions for the cities of tomorrow.

## 7  Conclusion

In this thesis, we developed a high precision and computationally efficient framework for forecasting short term taxi demand across New York City's taxi zones. Our approach was motivated by two fundamental challenges: first, the need to model temporal dependencies across multiple horizons—from rapid, hour to

hour fluctuations to regular daily, weekly, and monthly cycles—and second, the requirement to capture spatial correlations across different regions, especially where demand in one zone exerts significant influence over its neighbors, such as in densely interconnected urban cores.

To address these challenges, we proposed a two stage modeling framework combining a Multi Scale Recurrent Neural Network (MS RNN) for temporal pattern extraction with a Graph Neural Network (GNN) based on the GraphSAGE architecture for spatial refinement.

The temporal modeling component, the MS RNN, was designed to comprehensively capture multi horizon temporal dynamics. We integrated three specialized recurrent submodules: a GRU network focused on modeling immediate hourly variations; an LSTM network dedicated to capturing daily and weekly cyclical patterns; and a Transformer encoder designed to recognize and exploit monthly seasonal trends. Each submodule independently processes the input time series at its corresponding temporal resolution. Their outputs are subsequently fused through a feature fusion layer, which intelligently combines these multiscale signals into a unified latent representation. This integrated feature vector, augmented by the raw hourly input, is fed into a final GRU layer to generate one hour ahead demand forecasts. Furthermore, the MS RNN supports a rolling incremental training strategy. Each day, the model shifts a 30 day sliding window forward, loading only the newly available 24 hours of data through optimized PyArrow filters. Fine tuning is performed solely on this new batch, significantly reducing computational costs by avoiding full retraining over the entire historical dataset.

For spatial refinement, we constructed a directed, weighted graph where each node represents a taxi zone and edge weights correspond to the origin–destination ride counts aggregated over the past 30 days. The node features consist of the preliminary demand predictions output by the MS RNN, concatenated with a normalized borough level passenger volume metric that provides context on local traffic density. This graph is processed by a two layer GraphSAGE network, which updates the node embeddings through neighborhood message passing. By propagating information between connected zones, GraphSAGE effectively smooths the raw MS RNN forecasts and corrects inconsistencies, particularly where demand in one zone is influenced by conditions in adjacent or frequently linked areas.

To ensure efficient deployment, we employed several engineering optimizations. Mixed precision training was used to accelerate computation and reduce memory footprint without sacrificing model accuracy. Gradient accumulation techniques allowed larger effective batch sizes under hardware memory constraints. Lightweight Transformer and GNN configurations were carefully designed to balance representational capacity and runtime efficiency. As a result, the full training pipeline completes each epoch in under six minutes on a standard workstation equipped with an Intel i7 12700 CPU and an NVIDIA RTX 3070 Ti GPU. Daily incremental updates, involving fine tuning on the latest 24 hours of data, complete in approximately 1.2 minutes, supporting near real time model refresh cycles.

Empirical evaluation confirms the effectiveness of our design. The MS RNN, when evaluated independently, achieves a mean absolute error (MAE) of 8.01 rides per hour and a mean squared error (MSE) of 154.55, significantly outperforming baseline architectures such as pure GRU (MAE = 24.56), pure LSTM (MAE = 13.95), and pure Transformer models (MAE = 57.72). When further refined through the GraphSAGE spatial adjustment stage, the prediction error decreases even further, reducing the MAE to 6.99, representing a substantial improvement in predictive accuracy. Moreover, the training and inference runtimes remain practical for real world deployment scenarios: per epoch training times are maintained under six minutes, daily updates are completed within approximately 1.5 minutes, and peak GPU memory usage consistently remains below 7 GB, ensuring cost effective scalability.

Despite these achievements, several limitations of the current framework are recognized:

First, the model tends to underestimate rare, extreme peaks in demand, such as those caused by sudden events or emergencies. These anomalies typically lie outside the training data's historical distribution, making them difficult to predict accurately. Incorporating exogenous features—such as event schedules, weather alerts, or emergency notifications—could help the model better anticipate and respond to such outliers.

Second, the framework relies on static daily updates to the origin–destination flow graph. Although the graph is refreshed once every 24 hours, it cannot capture intra day variations in spatial dynamics, such as temporary road closures, construction activities, or spontaneous traffic disruptions. Addressing this limitation would require more frequent graph updates or even a transition to a streaming graph construction paradigm.

Third, the model operates at a coarse spatial granularity, forecasting demand at the taxi zone level. While this level of aggregation simplifies modeling and reduces computational load, it can obscure sub zone heterogeneity and prevent the capture of fine grained urban mobility patterns. Extending the framework to finer spatial units, such as street segments or intersections, would provide a more detailed view of demand distribution but would also necessitate scalable neighbor sampling and hierarchical pooling mechanisms within the GNN to manage the increased graph complexity.

Finally, the current model produces point forecasts only, without providing any quantification of prediction uncertainty. In applications such as real time dispatching or resource allocation, having access to calibrated prediction intervals would be extremely valuable. Future work could explore Bayesian recurrent neural networks or ensemble based techniques, such as Monte Carlo dropout, to generate probabilistic forecasts and better inform risk sensitive decision making processes.

Based on these observations, several promising directions for future research are identified. Firstly, integrating auxiliary data sources such as weather conditions, public event calendars, and transit service status could enhance the model's ability to anticipate non regular fluctuations. Secondly, implementing true online learning frameworks, where both the RNN and GNN components

are updated in an hourly streaming fashion, would reduce forecast latency and improve responsiveness to evolving urban dynamics. Thirdly, adopting dynamic graph architectures, such as continuous time GNNs like EvolveGCN, could allow for real time updates of both edge weights and node embeddings, thus capturing rapid spatial shifts more effectively. Fourthly, pursuing finer spatial modeling based on street level or intersection level GPS trace data would significantly enhance the resolution of forecasts and uncover hidden mobility patterns. Lastly, the incorporation of uncertainty estimation methods would allow the system to output not only point predictions but also reliable prediction intervals, greatly improving its value for operational decision making.

In conclusion, this thesis presents a robust, scalable, and practical end to end spatiotemporal forecasting system that advances both methodological rigor and real world applicability. By seamlessly integrating multi scale temporal modeling with spatial graph based refinement, and by demonstrating the feasibility of fast, daily incremental updates, the proposed framework lays a strong foundation for the development of intelligent taxi dispatch systems and broader smart city applications aimed at optimizing urban mobility.

# References

[1] Urban transportation challenges and solutions for modern cities – navigate mobility, 08 2024.

[2] K.E. ArunKumar, Dinesh V. Kalaga, Ch. Mohan Sai Kumar, Masahiro Kawaji, and Timothy M. Brenza. Comparative analysis of gated recurrent units (gru), long short-term memory (lstm) cells, autoregressive integrated moving average (arima), seasonal autoregressive integrated moving average (sarima) for forecasting covid-19 trends. *Alexandria Engineering Journal*, 01 2022.

[3] Kelachi Enwere, Ethelbert Nduka, and Uchenna Ogoke. Comparative analysis of ridge, bridge and lasso regression models in the presence of multicollinearity. *IPS Intelligentsia Multidisciplinary Journal*, 3:1–8, 01 2023.

[4] Sabiheh Faghih, Arpita Shah, Zehan Wang, Abolfazl Safikhani, and Camille Kamga. Taxi and mobility: Modeling taxi demand using arma and linear regression. *Procedia Computer Science*, 177:186–195, 2020.

[5] William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. *arXiv:1706.02216 [cs, stat]*, 09 2018.

[6] Shan Jiang, Wentian Chen, Zhiheng Li, and Haiyang Yu. Short-term demand prediction method for online car-hailing services based on a least squares support vector machine. *IEEE Access*, 7:11882–11891, 2019.

[7] Jintao Ke, Hai Yang, Hongyu Zheng, Xiqun Chen, Yitian Jia, Pinghua Gong, and Jieping Ye. Hexagon-based convolutional neural network for

supply-demand forecasting of ride-sourcing services. *IEEE Transactions on Intelligent Transportation Systems*, 20:4160–4173, 11 2019.

[8] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv:1609.02907 [cs, stat]*, 02 2017.

[9] Guanyao Li, Shuhan Zhong, Xingdong Deng, Letian Xiang, S.-H. Gary Chan, Ruiyuan Li, Yang Liu, Ming Zhang, Chih-Chieh Hung, and Wen-Chih Peng. A lightweight and accurate spatial-temporal transformer for traffic forecasting. *IEEE transactions on knowledge and data engineering*, 35:10967–10980, 11 2023.

[10] Hongbo Li, Jincheng Wang, Yilong Ren, and Feng Mao. Intercity online car-hailing travel demand prediction via a spatiotemporal transformer method. *Applied Sciences*, 11:11750, 12 2021.

[11] Lingbo Liu, Zhilin Qiu, Guanbin Li, Qing Wang, Wanli Ouyang, and Liang Lin. Contextualized spatial-temporal network for taxi origin-destination demand prediction, 2019.

[12] Tao Liu, Aimin Jiang, Jia Zhou, Min Li, and Hon Keung Kwan. Graphsage-based dynamic spatial–temporal graph convolutional network for traffic prediction. *IEEE Transactions on Intelligent Transportation Systems*, 24:11210–11224, 10 2023.

[13] Jian Ma, Zhao Yang Dong, and Pei Zhang. Using a support vector machine (svm) to improve generalization ability of load model parameters. *Ieee.org*, 2018.

[14] Ahsan Nazir, Jingsha He, Nafei Zhu, Saima Siraj Qureshi, Siraj Uddin Qureshi, Faheem Ullah, Ahsan Wajahat, and Muhammad Salman Pathan. A deep learning-based novel hybrid cnn-lstm architecture for efficient detection of threats in the iot ecosystem. *Ain Shams Engineering Journal/Ain Shams Engineering Journal*, pages 102777–102777, 04 2024.

[15] Xinwu Qian and Satish V. Ukkusuri. Spatial variation of the urban taxi ridership using gps data. *Applied Geography*, 59:31–42, 05 2015.

[16] Filipe Rodrigues, Ioulia Markou, and Francisco C. Pereira. Combining time-series and textual data for taxi demand prediction in event areas: A deep learning approach. *Information Fusion*, 49:120–129, 09 2019.

[17] Fathi M Salem. Gated rnn: The gated recurrent unit (gru) rnn. *Springer eBooks*, pages 85–100, 10 2021.

[18] Benjamin Sanchez-Lengeling, Emily Reif, Adam Pearce, and Alex Wiltschko. A gentle introduction to graph neural networks. *Distill*, 6, 08 2021.

[19] Bruce Schaller. A regression model of the number of taxicabs in u.s. cities. *Journal of Public Transportation*, 8:63–78, 12 2005.

[20] Alex Sherstinsky. Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm) network.

[21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.

[22] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. *arXiv:1710.10903 [cs, stat]*, 02 2018.

[23] Billy M. Williams and Lester A. Hoel. Modeling and forecasting vehicular traffic flow as a seasonal arima process: Theoretical basis and empirical results. *Journal of Transportation Engineering*, 129:664–672, 11 2003.

[24] Jun Xu, Rouhollah Rahmatizadeh, Ladislau Boloni, and Damla Turgut. Real-time prediction of taxi demand using recurrent neural networks. *IEEE Transactions on Intelligent Transportation Systems*, 19:2572–2581, 08 2018.

[25] Jianzhe Xue, Yunting Xu, Dongcheng Yuan, Caoyi Zha, Hongyang Du, Haibo Zhou, and Dusit Niyato. Spatial-temporal generative ai for traffic flow estimation with sparse data of connected vehicles, 2020.

[26] Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting. *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, page 3634–3640, 07 2018.

[27] Kai Zhao, Denis Khryashchev, and Huy Vo. Predicting taxi and uber demand in cities: Approaching the limit of predictability. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1, 2020.

[28] Tamer Çetin and Kadir Yasin Eryigit. Estimating the effects of entry regulation in the istanbul taxicab market. *Transportation Research Part A: Policy and Practice*, 45:476–484, 07 2011.

# A   Appendices

## A.1   MS RNN + GraphSAGE

```python
1  import torch
2  import pandas as pd
3  import numpy as np
4  from sklearn.preprocessing import MinMaxScaler
5  import matplotlib.pyplot as plt
6  from torch import nn
7
8  #         GPU
9  device = torch.device( "cuda" if torch.cuda.
       is_available() else "cpu")
10 print("CUDA available:", torch.cuda.is_available())
11 print("Using device:", device)
12
13
14
15
16 target_date = pd.Timestamp('2021-03-05 12:00')  #
17
18
19
20
21
22 # month = '01'
23 file_name = f'data.parquet'
24 columns_to_load = ['pickup_datetime', 'PULocationID']
25 df_temp = pd.read_parquet(file_name, columns=
       columns_to_load)
26
27
28 #
29 excluded_zones = [103, 104, 105, 46, 264, 265]
30 df_temp = df_temp[~df_temp['PULocationID'].isin(
       excluded_zones)]
31 zoneTotalNumber = len(df_temp['PULocationID'].unique())
32 print('
       ----------------------------------------------------
       ')
33 print(len(df_temp['PULocationID'].unique()))
34 print('
       ----------------------------------------------------
       ')
35
36 df_temp['pickup_datetime'] = pd.to_datetime(df_temp['
       pickup_datetime'])
```

60

```python
37  df_temp['datetime'] = df_temp['pickup_datetime'].dt.
        floor('H')
38
39  df = df_temp
40
41
42  durations = {
43      '1h': 1,
44      '1d': 24,
45      '1w': 24 * 7,
46      '1m': 24 * 30
47  }
48  forecast_length = 1
49
50  sequence_length = durations['1m']
51  sf = sequence_length + forecast_length
52
53
54  global_predictions = []
55  global_true_values = []
56
57
58  start_date = target_date - pd.Timedelta(hours=
        sequence_length)
59
60  #
        target_date
61  df = df[(df['datetime'] >= start_date) & (df['datetime'
        ] <= target_date)]
62
63  skipZoneList = []
64  #
65  results = []
66
67
68
69  #
70  print(f"Filtered data from {start_date} to {target_date
        - pd.Timedelta(hours=1)}. Total rows: {len(df)}")
71  minGlobalDate = df['datetime'].min()
72  maxGlobalDate = df['datetime'].max()
73  #                    PULocationID
74  unique_zones = df['PULocationID'].unique()
75  count = len(df['PULocationID'].unique())
76  for zoneid in unique_zones:
77      print(f"Processing PULocationID: {zoneid}")
78      print(f"There are {count} zones left")
79      count -= 1
80      zone_df = df[df['PULocationID'] == zoneid]
81
```

```python
    #
    hourly_demand = zone_df.groupby('datetime').size().
        reset_index(name='passenger_count')
    hourly_demand = hourly_demand.sort_values('datetime
        ').reset_index(drop=True)



    #
    scaler = MinMaxScaler()
    hourly_demand['passenger_count_scaled'] = scaler.
        fit_transform(hourly_demand[['passenger_count'
        ]])

    # max_day = hourly_demand['datetime'].max().day



    if len(hourly_demand) < sf/2:
        print("missing critical data!")
        skipZoneList.append(zoneid)
        continue
    if target_date not in hourly_demand['datetime'].
        values:
        print(f"Target date {target_date} not found for
             PULocationID {zoneid}. Adding missing
             target date...")
        missing_row = pd.DataFrame({'datetime': [
            target_date], 'passenger_count': [0]})
        hourly_demand = pd.concat([hourly_demand,
            missing_row], ignore_index=True).sort_values
            ('datetime').reset_index(
            drop=True)
    print(f"Length of hourly_demand: {len(hourly_demand
        )}")
    print(f"Required sf: {sf}")
    if len(hourly_demand) < sf:
        print(f"Insufficient data for PULocationID {
            zoneid}, filling missing data.")
        #
        required_length = sf
        current_length = len(hourly_demand)
        missing_length = required_length -
            current_length



        #
        full_datetime_range = pd.date_range(
```

```python
            start=minGlobalDate,
            end=target_date,
            freq='H'
        )

        #
        existing_times = set(hourly_demand['datetime'])
        full_times = set(full_datetime_range)
        missing_times = sorted(full_times -
            existing_times)
        print(len(existing_times))
        print(len(full_times))
        print(len(missing_times))

        # print("Missing times:")
        # for t in missing_times:
        #       print(t)

        # print("Before filling missing data:")
        # print(f"hourly_demand shape: {hourly_demand.
            columns}")
        # print(hourly_demand.tail())   #

        print(f"hourly_demand columns before fill: {
            hourly_demand.shape}")

        #
        hourly_demand = hourly_demand.set_index('
            datetime').reindex(full_datetime_range).
            fillna(0).reset_index()
        hourly_demand.rename(columns={'index': '
            datetime'}, inplace=True)   #


        # print("After filling missing data:")
        # print(f"hourly_demand shape: {hourly_demand.
            columns}")
        # print(hourly_demand.tail())
        print(f"hourly_demand columns after fill: {
            hourly_demand.shape}")

        #
        hourly_demand['passenger_count_scaled'] =
            scaler.fit_transform(hourly_demand[['
            passenger_count']])

    X_1h, X_1d, X_1w, X_1m, y = [], [], [], [], []

    for i in range(len(hourly_demand) + 1 -
        sequence_length - forecast_length):
```

```python
156          #
157          x_1h = hourly_demand['passenger_count_scaled'].
                 iloc[
158                  i + sequence_length - durations['1h']:i
                         + sequence_length].values   # 1

159          x_1d = hourly_demand['passenger_count_scaled'].
                 iloc[
160                  i + sequence_length - durations['1d']:i
                         + sequence_length].values   # 1

161          x_1w = hourly_demand['passenger_count_scaled'].
                 iloc[
162                  i + sequence_length - durations['1w']:i
                         + sequence_length].values   # 1

163          x_1m = hourly_demand['passenger_count_scaled'].
                 iloc[i:i + sequence_length].values   # 1


164
165          # **

                 **
166          X_1h.append(x_1h)
167          X_1d.append(x_1d)
168          X_1w.append(x_1w)
169          X_1m.append(x_1m)
170
171          #                              1
172          y_val = hourly_demand['passenger_count_scaled'
                 ].iloc[
173                  i + sequence_length:i + sequence_length
                         + forecast_length].values
174          if len(y_val) == forecast_length:
175              y.append(y_val)
176
177      # **             NumPy         **
178      X_1h = np.array(X_1h, dtype=np.float32)
179      X_1d = np.array(X_1d, dtype=np.float32)
180      X_1w = np.array(X_1w, dtype=np.float32)
181      X_1m = np.array(X_1m, dtype=np.float32)
182      y = np.array(y, dtype=np.float32)
183
184      # **             PyTorch Tensor**
185      X_1h_tensor = torch.tensor(X_1h, dtype=torch.
             float32).unsqueeze(-1).to(device)   # (batch,
             sequence_length, 1)
186      X_1d_tensor = torch.tensor(X_1d, dtype=torch.
             float32).unsqueeze(-1).to(device)   # (batch,
             sequence_length, 1)
```

```python
187        X_1w_tensor = torch.tensor(X_1w, dtype=torch.
               float32).unsqueeze(-1).to(device)  # (batch,
               sequence_length, 1)
188        X_1m_tensor = torch.tensor(X_1m, dtype=torch.
               float32).unsqueeze(-1).to(device)  # (batch,
               sequence_length, 1)
189        y_tensor = torch.tensor(y, dtype=torch.float32).to(
               device)  # (batch, forecast_length)
190
191        # **                                      **
192        X_tensor = {
193            "1h": X_1h_tensor,
194            "1d": X_1d_tensor,
195            "1w": X_1w_tensor,
196            "1m": X_1m_tensor
197        }
198
199        print(f"X_1h_tensor shape: {X_1h_tensor.shape}")  #
               (batch, sequence_length, 1)
200        print(f"X_1d_tensor shape: {X_1d_tensor.shape}")  #
               (batch, sequence_length, 1)
201        print(f"X_1w_tensor shape: {X_1w_tensor.shape}")  #
               (batch, sequence_length, 1)
202        print(f"X_1m_tensor shape: {X_1m_tensor.shape}")  #
               (batch, sequence_length, 1)
203        print(f"y_tensor shape: {y_tensor.shape}")  # (
               batch, forecast_length)
204
205        #
206        import torch
207        from torch import nn
208
209        import torch
210        from torch import nn
211
212        import torch
213        import torch.nn as nn
214
215
216        class MultiScaleModel(nn.Module):
217            def __init__(self, hidden_size):
218                super(MultiScaleModel, self).__init__()
219
220                self.hidden_size = hidden_size
221
222                # Define LSTM (for daily & weekly patterns)
223                self.lstm_1d = nn.LSTM(1, hidden_size,
                       batch_first=True)
224                self.lstm_1w = nn.LSTM(1, hidden_size,
                       batch_first=True)
```

65

```python
            # Define Transformer (for monthly trends)
            self.input_projection = nn.Linear(1,
                hidden_size)  # Projects input to match
                hidden size
            self.transformer_1m = nn.Transformer(
                hidden_size, nhead=4, num_encoder_layers
                =2, batch_first=True)

            # Feature fusion layer (Combining LSTM +
                Transformer outputs)
            self.feature_fusion = nn.Linear(hidden_size
                * 3, hidden_size)

            # Final GRU (now receives raw 1-hour
                sequence + fused long-term trends)
            self.gru = nn.GRU(hidden_size + 1,
                hidden_size, batch_first=True)

            # Final prediction layer (Forecasts next
                hour demand)
            self.fc = nn.Linear(hidden_size, 1)

        def forward(self, x):
            # Extract different time-scale features
            x_1h, x_1d, x_1w, x_1m = x["1h"], x["1d"],
                x["1w"], x["1m"]

            # Process 1-day & 1-week data with LSTMs
            _, (h_1d, _) = self.lstm_1d(x_1d)
            h_1d = h_1d[-1]  # Take the last hidden
                state

            _, (h_1w, _) = self.lstm_1w(x_1w)
            h_1w = h_1w[-1]  # Take the last hidden
                state

            # Process 1-month data with Transformer
            x_1m = self.input_projection(x_1m)
            x_1m = x_1m.permute(1, 0, 2)  # Adjust
                dimensions for Transformer
            h_1m = self.transformer_1m(x_1m, x_1m)[-1]

            # Fuse LSTM and Transformer outputs
            fused_trend = torch.cat([h_1d, h_1w, h_1m],
                dim=1)
            fused_trend = self.feature_fusion(
                fused_trend)  # Reduce dimensionality to
                hidden_size
```

```python
            # Prepare GRU input (1-hour data + fused
                long-term trend)
            batch_size, seq_len, _ = x_1h.shape
            fused_trend_expanded = fused_trend.
                unsqueeze(1).repeat(1, seq_len, 1)  #
                Expand to match sequence length
            x_gru_input = torch.cat([x_1h,
                fused_trend_expanded], dim=2)  # [batch,
                 seq_len, hidden_size + 1]

            # Process with GRU
            _, h_gru = self.gru(x_gru_input)
            h_gru = h_gru[-1]  # Take last hidden state

            # Final Prediction
            output = self.fc(h_gru)
            return output


    #
    #
    hidden_size = 64
    model = MultiScaleModel(hidden_size).to(device)  #
                        GPU/CPU
    criterion = nn.MSELoss()
    optimizer = torch.optim.Adam(model.parameters(), lr
        =0.001)

    #                   GPU/CPU
    #        PyTorch   T  e  n  s  o  r
    X_tensor = {
        "1h": torch.tensor(X_1h, dtype=torch.float32).
            unsqueeze(-1).to(device),  # (batch,
            sequence_length, 1)
        "1d": torch.tensor(X_1d, dtype=torch.float32).
            unsqueeze(-1).to(device),  # (batch,
            sequence_length, 1)
        "1w": torch.tensor(X_1w, dtype=torch.float32).
            unsqueeze(-1).to(device),  # (batch,
            sequence_length, 1)
        "1m": torch.tensor(X_1m, dtype=torch.float32).
            unsqueeze(-1).to(device)  # (batch,
            sequence_length, 1)
    }

    #
    y_tensor = torch.tensor(y, dtype=torch.float32).to(
        device)  # (batch, forecast_length)

    # **            **
```

```python
293        print(f"X_1h_tensor shape: {X_tensor['1h'].shape}")
               # (batch, sequence_length, 1)
294        print(f"X_1d_tensor shape: {X_tensor['1d'].shape}")
               # (batch, sequence_length, 1)
295        print(f"X_1w_tensor shape: {X_tensor['1w'].shape}")
               # (batch, sequence_length, 1)
296        print(f"X_1m_tensor shape: {X_tensor['1m'].shape}")
               # (batch, sequence_length, 1)
297        print(f"y_tensor shape: {y_tensor.shape}")  # (
               batch, forecast_length)
298
299        #
300        epochs = 50
301        patience = 5
302        best_loss = float('inf')
303        counter = 0
304
305        model.train()
306        for epoch in range(epochs):
307            optimizer.zero_grad()
308            output = model(X_tensor)
309            loss = criterion(output, y_tensor)
310            loss.backward()
311            optimizer.step()
312
313            if loss.item() < best_loss:
314                best_loss = loss.item()
315                counter = 0
316            else:
317                counter += 1
318                if counter >= patience:
319                    print(f"Early stopping at epoch {epoch}
                          ")
320                    break
321
322            # print(f"Epoch {epoch + 1}/{epochs}, Loss: {
                   loss.item():.4f}")
323
324        #
325        model.eval()
326        with torch.no_grad():
327            predictions = model(X_tensor).cpu().numpy()  #
                                          CPU
328            true_values = y_tensor.cpu().numpy()  #
                                          CPU
329
330        #
331        predictions = scaler.inverse_transform(predictions)
332        true_values = scaler.inverse_transform(true_values)
```

```python
333         global_predictions.extend(predictions.flatten().
                tolist())
334         global_true_values.extend(true_values.flatten().
                tolist())
335         #
336         for pred, true_val in zip(predictions, true_values)
                :
337             results.append({'PULocationID': zoneid, '
                    Prediction': pred[0], 'True Value': true_val
                    [0]})
338
339
340
341
342 print(skipZoneList)
343 print(f"length{len(skipZoneList)}")
344 for skipped_id in skipZoneList:
345     results.append({'PULocationID': skipped_id, '
            Prediction': None, 'True Value': None})
346
347 #                       MAE     MSE
348 if global_predictions and global_true_values:
349     global_predictions = np.array(global_predictions)
350     global_true_values = np.array(global_true_values)
351     global_mae = np.mean(np.abs(global_predictions -
            global_true_values))
352     global_mse = np.mean((global_predictions -
            global_true_values) ** 2)
353     print(f"Overall Average MAE: {global_mae:.4f}")
354     print(f"Overall Average MSE: {global_mse:.4f}")
355 else:
356     print("No valid predictions available for computing
                global metrics.")
357
358
359 #               CSV
360 results_df = pd.DataFrame(results)
361 # output_file = f'GRU_{zoneTotalNumber-len(skipZoneList
        )}Zones_{epochs}epochs.csv'
362 # results_df.to_csv(output_file, index=False)
363
364 # print(f"Results saved to {output_file}")
365 #       lookup  table              Zone
366 lookup_table = pd.read_csv('taxi-zone-lookup.csv')
367
368 #     PULocationID               Zone
369 results_df = pd.merge(
370     results_df,
371     lookup_table[['LocationID', 'Zone']],
372     left_on='PULocationID',
```

69

```python
        right_on='LocationID',
        how='left'
)

#        Zone                                    Zone
merged_results = results_df.groupby('Zone', as_index=
        False).agg({
        'PULocationID': 'first',    #
                PULocationID
        'Prediction': 'mean',       #             Prediction

        'True Value': 'mean'        #             True Value

})

#
output_file = f'GRU_Merged_{epochs}epochs.csv'
merged_results.to_csv(output_file, index=False)
print(f"Merged results saved to {output_file}")


if not results_df.empty:
        results_df = results_df.dropna(subset=['Prediction'
                , 'True Value'])
        plt.figure(figsize=(12, 6))
        plt.plot(results_df['Prediction'], label='
                Prediction', linestyle='--', marker='o')
        plt.plot(results_df['True Value'], label='True
                Value', linestyle='-', marker='x')
        plt.title('Prediction vs True Value')
        plt.xlabel('Taxi_zone ID')
        plt.ylabel('Passenger Demand')
        plt.legend()
        plt.grid()
        plt.show()

else:
        print("No data available for plotting.")

import pandas as pd
import torch
import numpy as np
from torch_geometric.data import Data
from torch_geometric.utils import dense_to_sparse
from torch_geometric.nn import GCNConv
import torch.nn as nn
import torch.nn.functional as F
from torch_geometric.nn import SAGEConv
import random
```

```python
# seed = 42
# torch.manual_seed(seed)
# np.random.seed(seed)
# random.seed(seed)
# if torch.cuda.is_available():
#     torch.cuda.manual_seed_all(seed)
# ==============================
# 1.
# ==============================
edge_weight_csv = "edge_weight_matrix_with_flow.csv"
df_adj = pd.read_csv(edge_weight_csv, index_col=0)
adj_matrix = torch.tensor(df_adj.values, dtype=torch.
    float32)   #

#
edge_index, edge_attr = dense_to_sparse(adj_matrix)



#
zone_names = df_adj.index.tolist()
zone_idx_map = {zone: idx for idx, zone in enumerate(
    zone_names)}
N = len(zone_names)

# ==============================
# 2.          County Code
                            30
# ==============================




file_name = "data.parquet"
columns_to_load = ['pickup_datetime', 'PULocationID', '
    DOLocationID']
# df = pd.read_parquet(file_name, columns=
    columns_to_load)
df = df_temp
df['pickup_datetime'] = pd.to_datetime(df['
    pickup_datetime'])
# df['datetime'] = df['pickup_datetime'].dt.floor('H')




#               30
sequence_length = 24 * 30
```

71

```python
459  start_date = target_date - pd.Timedelta(hours=
         sequence_length)
460  df_30_days = df[(df['datetime'] >= start_date) & (df['
         datetime'] <= target_date)]
461
462  # excluded_zones = [103, 104, 105, 46, 264, 265]
463  # df = df[~df['PULocationID'].isin(excluded_zones)]
464  # zoneTotalNumber = len(df['PULocationID'].unique())
465  print(f"Total unique zones: {zoneTotalNumber}")
466
467  #                                          True Value
468  previous_hour = target_date - pd.Timedelta(hours=1)  #
469  previous_hour_data = df[df['datetime'] == previous_hour
         ]
470
471
472  true_values_dict = {}
473  for pulocation_id, group in previous_hour_data.groupby(
         'PULocationID'):
474      true_values_dict[pulocation_id] = len(group)  #
                                                    True
             Value
475
476  #        County Code      Borough
477  county_code_to_borough = {
478      1: "Bronx",
479      2: "Brooklyn",
480      4: "Queens",
481      5: "Staten Island",
482      6: "Manhattan"
483  }
484
485  #        taxi-zone-lookup.csv              LocationID
         County Code
486  zone_lookup_file = "taxi-zone-lookup.csv"
487  zone_lookup_df = pd.read_csv(zone_lookup_file).
         drop_duplicates(subset="LocationID")
488  location_to_county_code = dict(zip(zone_lookup_df["
         LocationID"], zone_lookup_df["Borough"].map(
489      lambda x: next((k for k, v in
             county_code_to_borough.items() if v == x), None)
490  )))
491
492
493
494  #     PULocationID
                            30                       County
495  county_volume = df_30_days.groupby("PULocationID").size
         ().reset_index(name="Total_Volume")
```

```python
496  county_volume["County"] = county_volume["PULocationID"
         ].map(location_to_county_code)
497  county_total_volume = county_volume.groupby("County")["
         Total_Volume"].sum().to_dict()
498
499
500
501
502
503  #
504  node_weights = torch.zeros((N,), dtype=torch.float32)
505
506  zone_to_county = dict(zip(zone_lookup_df["Zone"],
         zone_lookup_df["Borough"].map(
507      lambda x: next((k for k, v in
             county_code_to_borough.items() if v == x), None)
508  )))
509
510  for zone, idx in zone_idx_map.items():
511      county_code = zone_to_county.get(zone)  # Map zone
             to county name first
512      if county_code and county_code in
             county_total_volume:
513          node_weights[idx] = county_total_volume[
                 county_code]
514      else:
515          print(f"        Missing county volume data for
                 Zone {zone} (County: {county_code}), setting
                  weight to 0")
516          node_weights[idx] = 0  # Prevent NaN values
517
518  #                              [0, 1]
519  node_weights = node_weights / node_weights.max()
520
521  #
522  print("Node weights:", node_weights)
523
524
525
526
527
528  #
529  pred_csv = "GRU_Merged_50epochs.csv"
530  df_pred = pd.read_csv(pred_csv)
531
532  location_to_zone = dict(zip(zone_lookup_df["LocationID"
         ], zone_lookup_df["Zone"]))
533
534  df_pred["Zone"] = df_pred["PULocationID"].map(
         location_to_zone)
```

```python
node_pred = torch.full((N,), float('nan'), dtype=torch.
    float32)
node_label = torch.full((N,), float('nan'), dtype=torch
    .float32)

for _, row in df_pred.iterrows():
    loc_id = int(row["PULocationID"])
    if loc_id in excluded_zones:
        continue
    zone_str = row["Zone"]
    pred_val = row["Prediction"]
    true_val = row["True Value"]
    if isinstance(zone_str, str) and zone_str in
        zone_idx_map:
        ridx = zone_idx_map[zone_str]
        node_pred[ridx] = float(pred_val)
        node_label[ridx] = float(true_val)


#          NaN
print("node_pred NaN count:", torch.isnan(node_pred).
    sum().item())
print("node_label NaN count:", torch.isnan(node_label).
    sum().item())




#               NaN
valid_indices = torch.where(~torch.isnan(node_pred) & ~
    torch.isnan(node_label))[0]
print(f"                                  : {
    valid_indices.numel()}")

zone_names = [zone_names[i] for i in valid_indices.
    tolist()]


#                 old_index -> new_index
old_to_new = {old_idx.item(): new_idx for new_idx,
    old_idx in enumerate(valid_indices)}

#           edge_index                      valid_indices

valid_edges = (torch.isin(edge_index[0], valid_indices)
    ) & (torch.isin(edge_index[1], valid_indices))
edge_index = edge_index[:, valid_edges]


```

```python
574  #                   edge_index
575  edge_index = torch.tensor([[old_to_new[i.item()] for i
         in edge_index[0]],
576                            [old_to_new[j.item()] for j
                                  in edge_index[1]]], dtype
                                  =torch.long)

578  node_pred = node_pred[valid_indices]
579  node_label = node_label[valid_indices]
580  node_weights = node_weights[valid_indices]


583  #          node_pred       n o d e _ w e i g h t s   2

584  x_feat = torch.stack([node_pred, node_weights], dim=1)
         # shape [N, 2]
585  data = Data(x=x_feat, edge_index=edge_index)
586  data.y = node_label   # shape [N]


588  do=0.1
589  lr=0.01
590  hd=256
591  epochs = 300
592  class MultiScaleGraphSAGE(nn.Module):
593      def __init__(self, in_dim, hidden_dim):
594          super().__init__()
595          self.sage1 = SAGEConv(in_dim, hidden_dim)
596          self.sage2 = SAGEConv(hidden_dim, hidden_dim)
597          self.dropout = nn.Dropout(do)   #           Dropout

598          self.out_linear = nn.Linear(hidden_dim, 1)

600      def forward(self, data):
601          x, edge_index = data.x, data.edge_index   # **
              GraphSAGE              edge_weight**
602          x1 = self.dropout(F.gelu(self.sage1(x,
              edge_index)))
603          x2 = self.dropout(F.gelu(self.sage2(x1,
              edge_index)))
604          return self.out_linear(x2).squeeze(-1)

606  device = torch.device('cuda' if torch.cuda.is_available
         () else 'cpu')

608  #
609  model = MultiScaleGraphSAGE(in_dim=2, hidden_dim=hd).to
         (device)
610  optimizer = torch.optim.Adam(model.parameters(), lr=lr)
611  loss_func = nn.SmoothL1Loss()
612  #
```

```python
613    # scheduler = torch.optim.lr_scheduler.
           ReduceLROnPlateau(optimizer, mode='min', factor=0.5,
           patience=20, verbose=True)
614
615    scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(
           optimizer, T_max=epochs)
616
617    data = data.to(device)
618
619    model.train()
620    for epoch in range(1, epochs + 1):
621        optimizer.zero_grad()
622        pred = model(data)
623        loss = loss_func(pred, data.y)
624        loss.backward()
625        optimizer.step()
626        scheduler.step()
627
628        if epoch % 10 == 0:
629            # print(f"Epoch {epoch}/{epochs}, Loss = {loss.
                  item():.4f}")
630            current_lr = optimizer.param_groups[0]['lr']
631            print(f"Epoch {epoch}/{epochs}, Loss = {loss.
                  item():.4f}, LR = {current_lr:.6f}")
632    #
633    model.eval()
634    with torch.no_grad():
635        refined_pred = model(data)
636
637    y_true = data.y.cpu().numpy()
638    y_refined = refined_pred.cpu().numpy()
639    node_pred_val = data.x.squeeze(-1).cpu().numpy()
640
641
642
643
644
645    # -----------------------------
646    # 8)
647    # -----------------------------
648    #                           1
649    node_pred_val = data.x[:, 0].cpu().numpy().squeeze()
650    y_refined = y_refined.squeeze()
651    y_true = y_true.squeeze()
652
653    print("Final shapes:")
654    print("node_pred_val shape:", node_pred_val.shape)
655    print("y_refined shape:", y_refined.shape)
656    print("y_true shape:", y_true.shape)
657
```

```python
# 建立 DataFrame
output_df = pd.DataFrame({
    "ZoneName": zone_names,
    "GRU_Pred": node_pred_val,
    "Refined_Pred": y_refined,
    "True_Value": y_true
})

gru_pred = output_df["GRU_Pred"].values
refined_pred = output_df["Refined_Pred"].values
true_val = output_df["True_Value"].values

mae_gru = np.mean(np.abs(gru_pred - true_val))
mae_refined = np.mean(np.abs(refined_pred - true_val))

mse_gru = np.mean((gru_pred - true_val)**2)
mse_refined = np.mean((refined_pred - true_val)**2)

print(f"GRU vs True MAE = {mae_gru:.4f}")
print(f"GRU vs True MSE = {mse_gru:.4f}")
print(f"Refined vs True MAE = {mae_refined:.4f}")
print(f"Refined vs True MSE = {mse_refined:.4f}")

import matplotlib.pyplot as plt



#
methods = ["GRU","GNN(Refined)"]
mae_values = [mae_gru, mae_refined]

#
plt.figure(figsize=(6, 5))
plt.bar(methods, mae_values, color=['blue','orange'],
    alpha=0.7)

#
for i, v in enumerate(mae_values):
    plt.text(i, v + 0.2, f"{v:.2f}", ha='center',
        fontsize=12)

plt.ylabel("Mean Absolute Error (MAE)")
plt.title("Comparison of Prediction Errors (MAE)")
# plt.ylim(5, 15)  #
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.show()

#
mse_values = [mse_gru, mse_refined]
```

```
706  #
707  plt.figure(figsize=(6, 5))
708  plt.bar(methods, mse_values, color=['blue','orange'],
         alpha=0.7)
709
710  #
711  for i, v in enumerate(mse_values):
712      plt.text(i, v + 0.2, f"{v:.2f}", ha='center',
            fontsize=12)
713
714  plt.ylabel("Mean Absolute Error (MSE)")
715  plt.title("Comparison of Prediction Errors (MSE)")
716  plt.grid(axis='y', linestyle='--', alpha=0.7)
717  plt.show()
718
719
720  plt.scatter(true_val, gru_pred, label="GRU", alpha=0.6)
721  plt.scatter(true_val, refined_pred, label="GNN", alpha
         =0.6)
722
723  min_val = min(true_val.min(), gru_pred.min(),
         refined_pred.min())
724  max_val = max(true_val.max(), gru_pred.max(),
         refined_pred.max())
725  plt.plot([min_val, max_val], [min_val, max_val], 'k-',
         label="GROUND TRUTH")
726  plt.xlabel("True Values")
727  plt.ylabel("Predictions")
728  plt.legend()
729  plt.show()
730
731  output_df.to_csv(
732      "final_predictions_multiscale.csv",
733      columns=["ZoneName", "GRU_Pred", "Refined_Pred", "
            True_Value"],
734      index=False,
735      encoding="utf-8"
736  )
737  print("result saved to 'final_predictions_multiscale.
         csv'")
```

Listing 19: MultiScaleModel (MS–RNN) implementation

## A.2   Heatmap generation

```
1  import os
2  import sys
3  import json
4  import geojson
```

78

```python
import pandas as pd
import torch
import numpy as np
from torch_geometric.utils import dense_to_sparse
from PyQt5.QtWidgets import QApplication, QMainWindow,
    QVBoxLayout, QHBoxLayout, QWidget, QLineEdit,
    QPushButton, \
    QTextEdit
from PyQt5.QtWebEngineWidgets import QWebEngineView
from PyQt5.QtCore import QUrl

#
df_gru = pd.read_csv("GRU_Merged_50epochs.csv")
df_gru.set_index("Zone", inplace=True)

zone_lookup_df = pd.read_csv("taxi-zone-lookup.csv")
location_to_zone = dict(zip(zone_lookup_df["LocationID"
    ], zone_lookup_df["Zone"]))
zone_to_location = {v: k for k, v in location_to_zone.
    items()}   #

edge_weight_csv = "edge_weight_matrix_with_flow.csv"
df_adj = pd.read_csv(edge_weight_csv, index_col=0)
adj_matrix = torch.tensor(df_adj.values, dtype=torch.
    float32)   #

#
edge_index, edge_attr = dense_to_sparse(adj_matrix)

zone_names = df_adj.index.tolist()
zone_idx_map = {zone: idx for idx, zone in enumerate(
    zone_names)}


#
pred_csv = "final_predictions_multiscale.csv"
df_pred = pd.read_csv(pred_csv)

geojson_file = "taxi_zones.geojson"
with open(geojson_file, "r") as f:
    nyc_geojson = json.load(f)

#        Zone Name
predictions = dict(zip(df_pred["ZoneName"], df_pred["
    Refined_Pred"]))

#           true_values_dict  True_Value
true_values_dict = dict(zip(df_pred["ZoneName"],
    df_pred["True_Value"]))
```

```python
47  #               ZoneName       zone_idx_map
48  true_values_dict = {zone_idx_map[k]: v for k, v in
        true_values_dict.items() if k in zone_idx_map}
49
50  #
51  print(f"          true_values_dict                    : {
        len(true_values_dict)}")
52
53
54  min_pred, max_pred = min(predictions.values()), max(
        predictions.values())
55
56  #
57  from branca.colormap import LinearColormap
58
59  colormap = LinearColormap(
60      colors=['blue', 'cyan', 'yellow', 'orange', 'red'],
61      vmin=min_pred,
62      vmax=max_pred * 0.8
63  )
64
65  #         GeoJSON
66  for feature in nyc_geojson["features"]:
67      zone_name = feature["properties"]["zone"]
68      if zone_name in predictions:
69          pred_value = predictions[zone_name]
70          color = colormap(pred_value)
71          feature["properties"]["Refined_Pred"] =
                pred_value
72          feature["properties"]["fillColor"] = color
73          feature["properties"]["style"] = {
74              "fillColor": color,
75              "color": "black",
76              "weight": 1,
77              "fillOpacity": 0.8
78          }
79      else:
80          feature["properties"]["fillColor"] = "#cccccc"
81          feature["properties"]["Refined_Pred"] = "N/A"
82
83  #               HTML
84  import folium
85
86  nyc_map = folium.Map(location=[40.7128, -74.0060],
        zoom_start=11)
87  folium.GeoJson(
88      nyc_geojson,
89      name="Taxi Zones",
90      tooltip=folium.GeoJsonTooltip(fields=["zone", "
            Refined_Pred", "fillColor"],
```

```python
                                            aliases=["Zone:", "
                                                Prediction:", "
                                                Color Code:"]),
        style_function=lambda feature: feature["properties"
            ].get("style", {
            "fillColor": "#cccccc",  #
            "color": "black",
            "weight": 1,
            "fillOpacity": 0.5
        })
).add_to(nyc_map)


nyc_map.save("nyc_taxi_prediction_map.html")


# -------------- GUI          --------------
from PyQt5.QtGui import QFont

class TaxiDispatchWindow(QMainWindow):
    def __init__(self):
        super().__init__()

        self.setWindowTitle("NYC Taxi Dispatch System")
        self.setGeometry(100, 100, 1400, 800)

        # **             **
        font = QFont("Arial", 14)  # 14px
        result_font = QFont("Courier", 12)  #


        # **                   **
        central_widget = QWidget()
        self.setCentralWidget(central_widget)

        # **       **
        layout = QHBoxLayout()
        central_widget.setLayout(layout)

        # **        -             **
        left_panel = QVBoxLayout()

        # **               Q L i n e E d i t   **
        self.input_box = QLineEdit()
        self.input_box.setPlaceholderText("
                            Zone ID (    76)...")
        self.input_box.setFont(font)  #

        left_panel.addWidget(self.input_box)
```

```python
        # **              Q P u s h B u t t o n
            **
        self.continue_button = QPushButton("
                ")
        self.continue_button.setFont(font)  #

        self.continue_button.clicked.connect(self.
            dispatch_taxi)
        left_panel.addWidget(self.continue_button)


        # **                    Q T e x t E d i t    **
        self.result_text = QTextEdit()
        self.result_text.setReadOnly(True)
        self.result_text.setFont(result_font)   #

        left_panel.addWidget(self.result_text)


        # **        -         **
        self.browser = QWebEngineView()
        html_file = os.path.abspath("
            nyc_taxi_prediction_map.html")
        self.browser.setUrl(QUrl.fromLocalFile(
            html_file))


        # **                **
        layout.addLayout(left_panel, 2)   #
            2/5
        layout.addWidget(self.browser, 3)   #
            3/5

    def dispatch_taxi(self):
        target_zone_id = self.input_box.text().strip()
        if not target_zone_id.isdigit():
            self.result_text.setText("
                            Zone  I D  ")
            return

        target_zone_id = int(target_zone_id)
        if target_zone_id not in location_to_zone:
            self.result_text.setText(f"     Zone ID {
                target_zone_id}
                                    ")
            return

        zone_name = location_to_zone[target_zone_id]
        if zone_name not in zone_idx_map:
            self.result_text.setText(f"        Zone {
                target_zone_id}
                ")
            return
```

```python
        zone_index = zone_idx_map[zone_name]
        result_log = f"                         : {zone_name}
            (Zone ID: {target_zone_id})             : {
            zone_index}\n"

        if zone_name in df_gru.index and zone_index in
            true_values_dict:
            #
            target_gru_pred = df_gru.loc[zone_name, "
                Prediction"]
            target_current = true_values_dict.get(
                zone_index, 0)
            target_required_cabs = target_gru_pred -
                target_current

            result_log += f"
                                        : {
                target_required_cabs:.2f}\n"

            #


            if target_required_cabs > 0:
                result_log += f"                     {
                    zone_name}

                    \n"
            else:
                #
                                                BFS

                from collections import deque
                visited = {zone_index}
                queue = deque([(zone_index, 0)])  # (
                              ,                  )
                candidates = []   #


                current_level = None

                #      B    F    S
                while queue:
                    current, level = queue.popleft()
                    #


                    if level > 0:
                        neighbor_zone_name = zone_names
                            [current]
```

```python
                        if neighbor_zone_name in df_gru
                            .index and current in
                            true_values_dict:
                            neighbor_gru_pred = df_gru.
                                loc[neighbor_zone_name,
                                "Prediction"]
                            neighbor_current =
                                true_values_dict.get(
                                current, 0)
                            shortage =
                                neighbor_gru_pred -
                                neighbor_current
                            if shortage > 0:
                                #


                                if current_level is
                                    None:
                                    current_level =
                                        level
                                if level ==
                                    current_level:
                                    candidates.append((
                                        current,
                                        shortage))
                #
                neighbors = set(edge_index[1][
                    edge_index[0] == current].cpu().
                    numpy())
                for nb in neighbors:
                    if nb not in visited:
                        visited.add(nb)
                        queue.append((nb, level +
                            1))
                #


                if queue and candidates and queue
                    [0][1] > current_level:
                    break

            if candidates:
                best_zone = max(candidates, key=
                    lambda x: x[1])
                best_zone_name = zone_names[
                    best_zone[0]]
                #

                Zone ID
```

```
223                          best_zone_id = zone_to_location.get
                                 (best_zone_name, "        ")
224                          result_log += f"
                                                            {
                                 best_zone_name} (Zone ID: {
                                 best_zone_id},          : {
                                 best_zone[0]})           : {
                                 best_zone[1]:.2f}\n"
225                  else:
226                      result_log += "

                                 \n"
227          else:
228              result_log += "
                                                    \n"

229
230          self.result_text.setText(result_log)
231

232
233  if __name__ == "__main__":
234      geojson_zones = {feature["properties"]["zone"] for
             feature in nyc_geojson["features"]}
235      lookup_zones = set(zone_lookup_df["Zone"])
236      pred_zones = set(df_pred["ZoneName"])
237
238      print(f"      GeoJSON                        : {len(
             geojson_zones)}")
239      print(f"      taxi-zone-lookup                    :
             {len(lookup_zones)}")
240      print(f"                                       : {len
             (pred_zones)}")
241
242      #
243      missing_in_geojson = lookup_zones - geojson_zones
244      missing_in_predictions = geojson_zones - pred_zones
245
246
247      app = QApplication(sys.argv)
248      window = TaxiDispatchWindow()
249      window.show()
250      sys.exit(app.exec_())
```

Listing 20: MultiScaleGraphSAGE implementation