

The Institution of  
Engineering and Technology

## ORIGINAL RESEARCH OPEN ACCESS

# Graph Neural Networks Empowered Origin-Destination Learning for Urban Traffic Prediction

Chuanting Zhang<sup>1,2</sup> | Guoqing Ma<sup>3</sup> | Liang Zhang<sup>3</sup> | Basem Shihada<sup>3</sup>

<sup>1</sup>School of Software, Shandong University, Jinan, China | <sup>2</sup>Shandong Key Laboratory of Intelligent Communication and Sensing-Computing Integration, Jinan, China | <sup>3</sup>Computer, Electrical and Mathematical Science and Engineering Division, King Abdullah University of Science and Technology (KAUST), Thuwal, Saudi Arabia

**Correspondence:** Basem Shihada ([basem.shihada@kaust.edu.sa](mailto:basem.shihada@kaust.edu.sa)) | Chuanting Zhang ([chuanting.zhang@sdu.edu.cn](mailto:chuanting.zhang@sdu.edu.cn))

**Received:** 10 January 2023 | **Revised:** 14 June 2024 | **Accepted:** 4 July 2024

**Handling Editor:** Wenguan Wang

**Funding:** This research was supported by the National Natural Science Foundation of China, Grant/Award Number: 62401338, by the Shandong Province Excellent Youth Science Fund Project (Overseas), Grant/Award Number: 2024HWYQ-028, and by the Fundamental Research Funds of Shandong University.

**Keywords:** data mining | deep neural networks | intelligent transportation systems | learning (artificial intelligence) | traffic engineering computing

## ABSTRACT

Urban traffic prediction with high precision is always the unrelenting pursuit of intelligent transportation systems and is instrumental in bringing smart cities into reality. The fundamental challenges for traffic prediction lie in the accurate modelling of spatial and temporal traffic dynamics. Existing approaches mainly focus on modelling the traffic data itself, but do not explore the traffic correlations implicit in origin-destination (OD) data. In this paper, we propose STOD-Net, a dynamic spatial-temporal OD feature-enhanced deep network, to simultaneously predict the in-traffic and out-traffic for each and every region of a city. We model the OD data as dynamic graphs and adopt graph neural networks in STOD-Net to learn a low-dimensional representation for each region. As per the region feature, we design a gating mechanism and operate it on the traffic feature learning to explicitly capture spatial correlations. To further capture the complicated spatial and temporal dependencies among different regions, we propose a novel joint feature learning block in STOD-Net and transfer the hybrid OD features to each block to make the learning process spatiotemporal-aware. We evaluate the effectiveness of STOD-Net on two benchmark datasets, and experimental results demonstrate that it outperforms the state-of-the-art by approximately 5% in terms of prediction accuracy and considerably improves prediction stability up to 80% in terms of standard deviation.

## 1 | Introduction

Smart transportation systems shape the blood vessels of a city and promote the rapid development of social society. Continuing this trend, highly accurate urban traffic prediction plays an essential role in intelligent transportation systems and facilitates the realisation of smart cities [1, 2]. Based on the knowledge of traffic prediction, intelligent traffic control can be achieved to enhance travel efficiency, reduce traffic congestion and improve the quality of life of citizens. Traffic prediction can

also improve public safety by predicting traffic volume for each and every region of a city [3]. In addition, a significant number of vehicular applications, such as route planning, navigation and travel time estimation, are heavily based on traffic condition evaluation [4–8]. Consequently, research on traffic prediction problems [9–12] has been active for decades and has received massive attention from both academia and industry.

Traffic prediction refers to the problem of predicting future traffic statuses, such as volume, speed and congestion, by

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2025 The Author(s). CAAI Transactions on Intelligence Technology published by John Wiley & Sons Ltd on behalf of The Institution of Engineering and Technology and Chongqing University of Technology.

modelling historical traffic data. The prediction horizon can be either short term, for example, 15 or 30 min, or long term, for example, several hours later. Depending on the application scenarios and data structures, traffic prediction can be applied to different road sensors/segments [13] or to any region of a city in arbitrary granularities [14]. In this study, we focus on short-term traffic volume prediction for all regions of a city.

Although it has huge importance in intelligent transportation systems, traffic prediction with high precision remains an extremely challenging task, as traffic status evolves in a highly complicated and nonlinear way, especially when the status transits between free flow, unstable, congestion and recovery [15]. Thus, the traffic volume series is highly nonlinear, and there exist both spatial and temporal dependencies among traffic series of different locations.

Many studies have been proposed to solve the traffic prediction problem. Some initial works focus on using statistical models, such as Kalman filtering [16] and autoregressive integrated moving average (ARIMA) [17]. These models have clear mathematical definitions, leading to high interpretability. However, it is essential to recognise that they have few parameters and come with low predictability limitations. In particular, they become inefficient when traffic series are nonlinear and have spatial dependencies.

The advancement of machine learning techniques [18] and the increasing availability of big traffic data have made data-driven models strong competitors to statistical models for traffic prediction. Support vector machines, random forests and neural networks have been explored for traffic prediction [19]. In particular, deep neural networks have drawn the most significant attention as their strong representation ability in an end-to-end manner [20, 21]. Long short-term memory networks (LSTM) and convolutional neural networks (CNN) are widely investigated in current literature. LSTM can model temporal dependence, as it has feedback connections between different neurons. Together with the inputs of other traffic series, the spatial dependence can also be captured by LSTM. Similarly, CNN is able to capture both spatial and temporal dependencies when the input has multiple sequences, as the convolution operation is a natural way to fuse information [22]. In particular, a deep CNN framework was proposed to solve the spatial-temporal modelling of traffic prediction [23]. Deep residual learning was also introduced into traffic prediction for collectively forecasting the in-traffic and out-traffic of a region simultaneously [24]. Apart from CNN-oriented frameworks, the growing prominence of graph neural networks (GNN) [25–29] has spurred interest in designing graph-based learning methods with the ability to naturally support spatial modelling of traffic prediction. Along this line, spatial-temporal graph convolutional networks were designed to extend traffic prediction from grid structures to general domains [30–32].

Most of the above work models spatial-temporal dependencies of different city regions using historical traffic data only. However, relying solely on historical traffic data may not capture the heterogeneity of spatial dependencies well. This is because almost all traffic sequences are time series with autocorrelation

and periodicity characteristics, resulting in non-negligible spatial correlations for any two different traffic series. The prediction model faces considerable challenges when learning from this kind of data and is unable to distinguish which traffic series is relevant and which is not, but may be, noise, especially for urban traffic prediction on a citywide prediction [33]. Besides, both spatial and temporal traffic dependencies are not static. Instead, they are highly time-varying, depending on the current traffic situation and road network. Take, for example, the places that are penetrated by highways: there exist more areas with higher spatial correlation because the drivers drive fast and cover longer distances. On the contrary, in the city centre, the speed of cars is relatively low, and the areas with higher spatial correlations are fewer.

To cope with the above challenges and enhance spatial-temporal dependency modelling, we resort to introducing origin-destination (OD) data into traffic prediction for two reasons. First, the OD data records the traffic interactions among regions of a city and provides detailed statistics that generate the traffic data of each region. Thus, the OD data can reflect the spatial correlation more directly. Second, OD data, together with static road network data, can be regarded as auxiliary information for regional dependence modelling from the perspective of data fusion. They provide different dimensions for understanding spatial-temporal dependencies of regions. Thus, based on the historical traffic data, the OD data and the static road network information, we propose a novel deep learning framework called STOD-Net, that is, spatial-temporal OD feature-enhanced deep network, for urban traffic prediction. To the best of our knowledge, we are the first to predict grid-based traffic volume by incorporating OD data and real-world road network data. We leverage CNN to model the historical traffic data, as they are sequences of matrices. Moreover, because the OD data are inherently in a graph structure, leading to the ineffectiveness of traditional CNN, we use GNNs to model them to obtain hidden representations. GNNs are novel neural networks designed especially for non-Euclidean data and are gaining increasing popularity in the deep learning field.

STOD-Net consists of four components. The first component is traffic feature learning, from which we obtain the hidden representations of different regions. In the second component, we perform OD feature learning, followed by a gating mechanism that functions as an enhanced spatial dependence modelling. Then, in the third component, we perform joint feature learning for the historical traffic data and the OD data using convolutional layers with introduced dense connectivity. To enhance the modelling of spatial-temporal dependencies, we further transfer the learnt OD features to the subsequent learning blocks, and this yields the fourth component of STOD-Net. In summary, our main contributions are listed as follows:

- We propose STOD-Net, a spatial-temporal deep learning framework for urban traffic prediction. STOD-Net accepts not only historical traffic data but also the OD data and static road network and can model these two kinds of data simultaneously.
- We introduce GNNs to model the OD data, as they retain the traffic interactions among regions. Based on the OD

representations, we further propose a gating mechanism that operates on the traffic representations and acts as a guide to enhance spatial-temporal dependency modelling. We present a joint modelling scheme for both dynamic and static OD data, making STOD-Net more robust in capturing spatial dependence.

- We validate the effectiveness of our STOD-Net on two real-world datasets, and the results demonstrate that STOD-Net significantly improves the prediction performance.

The rest of this paper is organised as follows: Section 2 is devoted to the related work on traffic prediction. Section 3 gives the problem formulation and preliminaries on graph neural networks. The details of our proposed STOD-Net are introduced in Section 4. We report the experimental results in Section 5 and conclude our paper in Section 6.

## 2 | Related Works

Current studies on traffic prediction can be classified into statistical model-based methods, traditional machine learning-based methods, deep learning-based methods and OD-related traffic prediction. This section is devoted to a brief review of related works in these three categories.

### 2.1 | Statistical Model-Based Traffic Prediction

Methods that fall into this category mainly centre on ARIMA and its variants. One of the earliest works that used the ARIMA model for traffic prediction was completed in Ref. [34] by Ahmed and Cook, in which they discovered that the ARIMA model can be more accurate than moving average and exponential smoothing in representing traffic volume data. Afterwards, researchers explored other versions of ARIMA for traffic prediction, including ARIMAX [35] and subset ARIMA [36]. Notwithstanding the popularity of ARIMA-based models in traffic prediction, they face significant limitations on predictive ability. The reason is that they are simply linear models, which assume the traffic is stationary. Therefore, they usually fail when dealing with highly complicated nonlinear traffic data [37]. Despite a few nonlinear models being proposed, such as heteroskedasticity-based models, the analytical study of nonlinear models is still in its infancy compared to linear models.

### 2.2 | Traditional Machine Learning-Based Traffic Prediction

Because of the inefficacy of statistical models in solving the traffic prediction problem, researchers resort to machine learning models, which are flexible, data-driven and can adjust the parameters automatically. Wu et al. applied support vector regression to traffic prediction [38] and achieved better prediction performances than statistical models because of its strong generalisation ability and guaranteed global minimum.  $k$ -nearest neighbour models were also explored for traffic prediction [39, 40] and found that they can achieve better prediction

performance than support vector regression in terms of mean squared error. Besides, random forests and artificial neural networks were also adopted for traffic prediction [41]. Though machine learning models have stronger predictive abilities than statistical models, they were more challenging to train 10 years ago, particularly for neural networks. Besides, traditional machine learning models have few parameters, leading to limited performance gains for traffic prediction.

### 2.3 | Deep Learning-Based Traffic Prediction

Deep learning models have been the mainstream for traffic prediction since the rise of deep neural networks. Different types of neural networks can be adopted for traffic prediction based on the data on hand. According to a recent survey [2], point-based data mainly adopt LSTM and GNN and trajectory-based data use CNN more.

For point-based data, LSTM [42], bidirectional LSTM [43] and LSTM with feature enhancement [44] are investigated in current literature, respectively. These models are good at capturing the temporal dependency of traffic volume. To further model the spatial dependency, GNN is introduced into traffic prediction [45, 46], as it can model the spatial relationships of different locations. Afterwards, many works on GNN-based traffic prediction methods were proposed, such as spatial-temporal GNN [47], GNN with attention scheme [13] and multiattention scheme [20].

For trajectory-based data, most prediction methods are built on top of CNN, such as DeepST [23] and ST-ResNet [24]. The key difference between these two models is that ST-ResNet introduces a residual learning strategy into traffic prediction, leading to a significant performance improvement. CNN is also the basic learning block of several other frameworks based on meta-learning [14], context-aware learning [48] and multiview learning [9, 31]. These frameworks can capture both spatial and temporal dependencies of different regions citywise but from different perspectives.

### 2.4 | OD-Related Traffic Prediction

Recently, some studies on adopting OD data for traffic prediction of a location or for predicting the OD matrix were proposed. For example, OD traffic was predicted by a hybrid spatial-temporal network [49], and a dynamic graph convolutional recurrent network was proposed [50] to model the dynamic characteristics of correlations among locations. Besides, the ideas of multitask learning and generative adversarial networks were also explored in the current literature. In Ref. [51], a multitask adversarial spatial-temporal network model was proposed to simultaneously predict the traffic flow and OD flow. An encoder-decoder structure was designed to capture the spatial-temporal dependencies of different regions, and a discriminative loss on task classification and an adversarial loss on shared feature extraction were incorporated to reduce information redundancy. The discriminative loss was also adopted in Ref. [52] for traffic flow prediction, yet under the framework

of generative adversarial nets (GAN) and the proposed model is TrafficGAN. One advantage of TrafficGAN is that a deformable convolution kernel for CNN is adopted to better handle the input road network data. More recently, Wang et al. proposed MC-STGCN [32], a multivariate correlation-aware spatiotemporal graph convolutional network for multiscale traffic prediction. In MC-STGCN, auxiliary information such as traffic speed and traffic occupancy rate was introduced into traffic prediction for correlation measurement and accuracy improvements.

The fundamental difference between our work and the above-mentioned studies lies in the fact that we use the OD data and static road network information to enhance spatial dependency modelling, instead of relying on the static distance between locations. We further propose a gating mechanism that can effectively fuse the historical traffic data and the OD data.

### 3 | Problem Formulation and Preliminaries

This section is devoted to the problem formulation of urban traffic prediction and the introduction of graph convolutional networks.

#### 3.1 | Problem Formulation

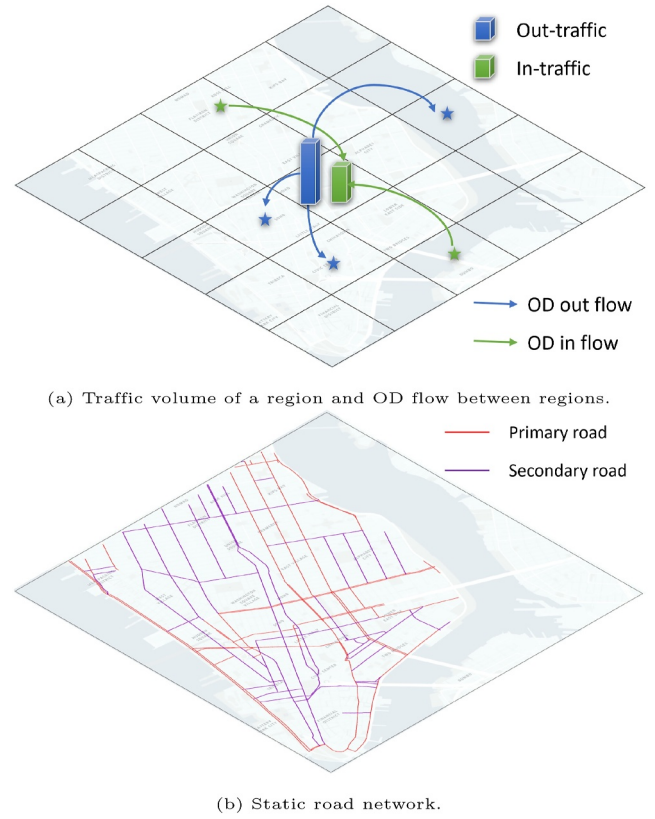
**Definition 1.** (Region). Many feasible ways exist to define a region of a city, and in this paper, we consider the following scenario: The city is evenly partitioned into an  $I \times J$  grid map based on the longitude and latitude, and a grid  $(i, j)$  denotes a region.

**Definition 2.** (In-traffic/out-traffic). For each region  $n$ , we consider two types of traffic volume, that is, in-traffic and out-traffic. The in-traffic,  $x^{(in,ij)}$ , indicates how many vehicles enter region  $n$  from neighbouring regions. Similarly, the out-traffic,  $x^{(out,ij)}$ , indicates how many vehicles drive out of this region. We denote the traffic volume as  $\mathbf{X}_t \in \mathbb{R}^{2 \times I \times J}$  for the grid map at time slot  $t$ ,<sup>1</sup> where  $t \in [1, 2, 3, \dots, T]$ .

**Definition 3.** (OD data). At  $t$ -th time slot, some vehicles are driving from region  $(i, j)$  to region  $(i', j')$  and others may drive from  $(i', j')$  to  $(i, j)$ . These records are called origin-destination flow data, and we denote the traffic from  $(i, j)$  to  $(i', j')$  as  $d^{(i,j),(i',j')}$ . Thus, we use  $\mathbf{D}_t \in \mathbb{R}^{N \times N}$  to denote the OD flow data at the  $t$ -th time slot, where  $N = I \times J$  represents the total number of regions of the city.

**Definition 4.** (Static road network). All the regions are connected by physical roads, based on which we obtain an adjacency matrix  $\mathbf{S} \in \mathbb{R}^{N \times N}$ , denoting the static road network.

To make the understanding of these definitions easy, a toy example of the defined region, in- and out-traffic, OD flow and the corresponding static road network is displayed in Figure 1.



**FIGURE 1** | Example of the in- and out-traffic, OD in- and out-flow and the static road network of New York City.

**Definition 5.** (Definition). Given a series of historical traffic data  $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_t\}$ , the corresponding OD flow data  $\{\mathbf{D}_1, \mathbf{D}_2, \dots, \mathbf{D}_t\}$  and static road network  $\mathbf{S}$ , predict  $\mathbf{X}_{t+1}$ .

Mathematically, our problem can be formally defined as follows:

$$w^* = \underset{w}{\operatorname{argmin}} \mathcal{L}(f(\mathcal{X}_t, \mathcal{D}_t, \mathbf{S}|w), \mathbf{X}_{t+1}), \quad (1)$$

where  $w$  denotes all the parameters of our model  $f(\cdot)$ .  $\mathcal{X}_t$  and  $\mathcal{D}_t$  are the input features extracted from traffic volume data and OD data, respectively, and their construction will be detailed in Section 4. Moreover,  $\mathcal{L}(\cdot, \cdot)$  is a loss function used to measure the goodness of our prediction model.

#### 3.2 | Graph Neural Networks

Graph neural networks (GNNs) refer to frameworks for representation learning on graphs, and their goal is to automatically learn a low-dimensional representation for the nodes, edges or the whole graph in an end-to-end fashion. Most GNN frameworks have a unified architecture, which can be described as follows:

$$\mathbf{H}^{(l+1)} = g(\mathbf{H}^{(l)}, \mathcal{A}), \quad (2)$$

where  $l$  indexes the layers and  $\mathbf{H}^{(l)}$  denotes the hidden representation in layer  $l$ . Note that  $\mathbf{H}^{(0)} = \mathbf{X}$  represents the input feature matrix.  $\mathcal{A}$  is the adjacency matrix corresponding to the



graph.  $g$  is a nonlinear function, and different GNN frameworks have different choices of  $g$ .

One of the most popular GNN frameworks is graph convolutional networks (GCNs) and the update rule is denoted as follows:

$$g(\mathbf{H}^{(l)}, \mathcal{A}) = \sigma(\bar{\mathcal{A}}\mathbf{H}^{(l)}\mathbf{W}^{(l)}), \quad (3)$$

where  $\bar{\mathcal{A}}$  is the symmetrically normalised graph Laplacian of  $\mathcal{A}$  with self-loops,  $\mathbf{W}^{(l)}$  is a parameter matrix of layer  $l$ , and  $\sigma$  is the rectified linear unit (ReLU) function. Equation (3) indicates that when updating the representation for a node in the graph, GCN performs a weighted average of all neighbours' representations of that node.

The attention scheme can also be considered in  $g$ , and this yields to graph attention networks [25] or GAT for short. The update rule in GAT can be described as follows:

$$g(\mathbf{H}^{(l)}, \mathcal{A}) = \sigma(\mathbf{A}^{(l)}\mathbf{H}^{(l)}\mathbf{W}^{(l)}). \quad (4)$$

Equation (4) is similar to Equation (3) except that  $\mathbf{A}^{(l)}$  is the attention weight matrix for layer  $l$  instead of the graph Laplacian of  $\mathcal{A}$ .

In this study, we use GATs to learn representations for each region, based on which we propose a gating mechanism to enhance characterisation of the hidden spatial dependencies.

## 4 | STOD-Net Architecture

In this section, we first give a general introduction to our proposed traffic prediction framework. Then, we detail each of its components to demonstrate how we model the spatial and temporal dependencies of urban traffic.

### 4.1 | Overall Architecture

To solve Equation (1), we propose a novel traffic prediction framework, that is, STOD-Net, which can effectively fuse the information from both the traffic data and the OD data, thereby capturing the spatial and temporal correlations among different regions simultaneously. Figure 2 illustrates the details of STOD-

Net, in which its fundamental component, that is, spatial-temporal origin-destination learning (STOD), is also included.

As current literature has shown, traffic volume's temporal correlations can be effectively characterised by closeness dependence, periodicity dependence and trend dependence [24]. Thus, in this paper, we embrace this point of view and resort to three separate subnetworks to learn the hidden representations for these temporal dependencies to capture their different impacts on future traffic volume prediction. To reduce complexity, these three subnetworks are designed with a shared architecture. As shown in Figure 2, there are three types of inputs, that is, the traffic volume data  $\mathcal{X}_t = \{\mathbf{X}_t^c, \mathbf{X}_t^p, \mathbf{X}_t^r\}$ , the OD data  $\mathcal{D}_t = \{\mathbf{D}_t^c, \mathbf{D}_t^p, \mathbf{D}_t^r\}$  and the static road network data denoted by the adjacency matrix  $\mathbf{S}$ . Specifically,  $\{\mathbf{X}_t^c, \mathbf{D}_t^c, \mathbf{S}\}$ ,  $\{\mathbf{X}_t^p, \mathbf{D}_t^p, \mathbf{S}\}$  and  $\{\mathbf{X}_t^r, \mathbf{D}_t^r, \mathbf{S}\}$  are used to model the closeness dependence, periodicity dependence and trend dependence, respectively. The constructions of  $\mathbf{X}_t^c$  and  $\mathbf{D}_t^c$  are as follows:

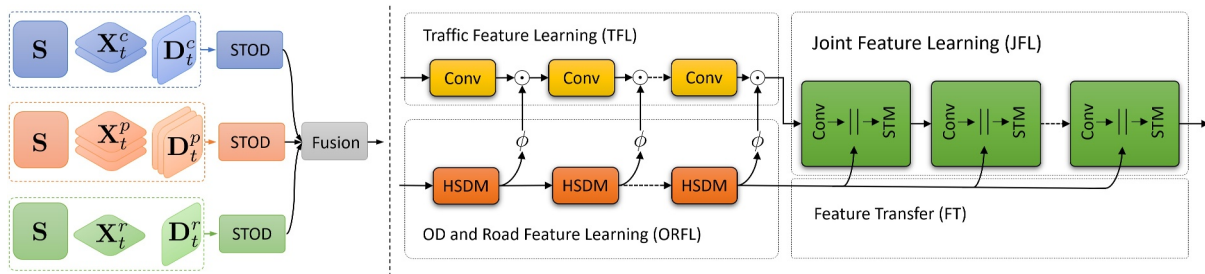
$$\begin{aligned} \mathbf{X}_t^c &= \mathbf{X}_{t-1} \parallel \mathbf{X}_{t-2} \parallel \cdots \parallel \mathbf{X}_{t-L_c}, \\ \mathbf{D}_t^c &= \mathbf{D}_{t-1} \parallel \mathbf{D}_{t-2} \parallel \cdots \parallel \mathbf{D}_{t-L_c}, \end{aligned} \quad (5)$$

where  $L_c$  is the length of the closeness sequence and  $\parallel$  denotes the concatenation operator. Similarly,  $\mathbf{X}_t^p$ ,  $\mathbf{D}_t^p$ ,  $\mathbf{X}_t^r$  and  $\mathbf{D}_t^r$  can be obtained as follows:

$$\begin{aligned} \mathbf{X}_t^p &= \mathbf{X}_{t-P} \parallel \mathbf{X}_{t-2P} \parallel \cdots \parallel \mathbf{X}_{t-L_pP}, \\ \mathbf{D}_t^p &= \mathbf{D}_{t-P} \parallel \mathbf{D}_{t-2P} \parallel \cdots \parallel \mathbf{D}_{t-L_pP}, \\ \mathbf{X}_t^r &= \mathbf{X}_{t-R} \parallel \mathbf{X}_{t-2R} \parallel \cdots \parallel \mathbf{X}_{t-L_rR}, \\ \mathbf{D}_t^r &= \mathbf{D}_{t-R} \parallel \mathbf{D}_{t-2R} \parallel \cdots \parallel \mathbf{D}_{t-L_rR}, \end{aligned} \quad (6)$$

where  $L_p$  and  $L_r$  are the lengths of the periodicity sequence and trend sequence, and  $P$  and  $R$  are the period and trend span corresponding to 1 day and 1 week, respectively. The adjacency matrix  $\mathbf{S}$  is constructed from physical roads, and the details will be given in the following subsection.

For STOD, as shown on the right of Figure 2, it consists of four components, that is, traffic feature learning (TFL), OD and road feature learning (ORFL), feature transfer (FT) and joint feature learning (JFL). The role of TFL is to learn initial representations for the traffic data. Like TFL, ORFL learns representations for the OD data and the static road network and then fuses them with the traffic representations. FT



**FIGURE 2** | STOD-Net framework and the corresponding key component STOD: convolutional (Conv), hybrid spatial dependence modelling (HSDM) and spatial-temporal modelling (STM).  $\odot$  and  $\parallel$  denote the Hadamard product and concatenation operation.  $\phi$  is a mapping function that implements the gating mechanism.

means we transfer the learnt OD and road features to subsequent layers for spatial and temporal modelling. As for JFL, its function is to perform joint feature learning for the fused representations, including the traffic information and the OD information. Once we obtain the final representations for the closeness, periodicity and trend data, we fuse them together, as shown on the left of Figure 2, followed by a nonlinear activation, and we get the predicted traffic volume  $\hat{\mathbf{X}}_{t+1}$  for time slot  $t + 1$ .

In the next section, we first introduce how we construct the adjacency matrix for the static road network, and then we give technical details for each of the components of our proposed STOD. For brevity, we use  $\{\mathbf{X}_t, \mathbf{D}_t, \mathbf{S}\}$  to denote any of the data sequences  $\{\mathbf{X}_t^c, \mathbf{D}_t^c, \mathbf{S}\}$ ,  $\{\mathbf{X}_t^p, \mathbf{D}_t^p, \mathbf{S}\}$  and  $\{\mathbf{X}_t^r, \mathbf{D}_t^r, \mathbf{S}\}$  and omit the subscript. This is because the three subnetworks share the same architecture, which indicates the inputs receive exactly the same operations; only the outputs are different.

#### 4.2 | Static Road Network Construction

We adopt an adjacency matrix to represent the static road network, in which nodes denote regions and edges denote connection weights between regions. The weights are constructed from physical roads that are crawled from OpenStreetMap. The details of weight construction are as follows:

- We extract the information on all kinds of road types and classify them into four categories, that is, primary, secondary, tertiary and others;
- For each road of these four categories, we extract the region set using a depth-first search strategy. Then, for any given region pair  $u$  and  $v$ , the connection weight is defined as follows:

$$s_{u,v}^z = \frac{1}{d_{u,v}^z}, \quad (7)$$

where  $z \in \{\text{primary, secondary, tertiary, others}\}$  denotes the road type and  $d_{u,v}^z$  the shortest path between  $u$  and  $v$  for the case of road type  $z$ .

- We assign a weight  $\rho_z$  for road type  $z$  because different roads have different capacities. For example, primary roads normally have more lanes than secondary roads; thus, it has a greater impact on traffic. For region pair  $u$  and  $v$ , we summarise the weights of different road types and yield the final weight.

$$s_{u,v} = \sum_i \rho_z s_{u,v}^z. \quad (8)$$

Figure 3 shows a toy example for the construction of the adjacency matrix of the static road network. In this example, two roads belong to different types, connecting nine regions. The weights of different road types are set to 1 for simplicity.

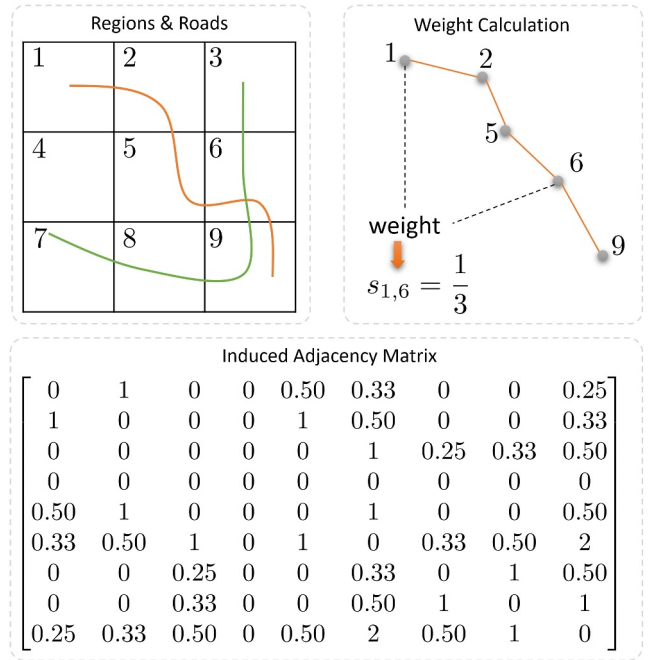


FIGURE 3 | A toy example for constructing the adjacency matrix, that is,  $\mathbf{S}$ , representing the static road network.

#### 4.3 | Traffic Feature Learning

For citywise urban traffic prediction, convolutional networks have shown superior performance in capturing the spatial correlation of traffic volume of different regions. Thus, we adopt convolutional networks to learn the hidden representations of traffic volume data. In TFL, the obtained representations at layer  $l$  are defined as follows:

$$\mathbf{X}_t^{(l)} = \sigma(\mathbf{W}_{\text{TFL}}^{(l)} * \mathbf{X}_t^{(l-1)}), \quad (9)$$

where  $*$  is the convolution operation and  $\sigma$  is the ReLU function. In TFL, there are in total  $L_{\text{TFL}}$  layers for traffic feature learning, and the final representation is denoted as  $\mathbf{X}_t^{(L_{\text{TFL}})}$ .

#### 4.4 | OD and Road Feature Learning

The OD data records the traffic volume interactions from one region to another and explicitly reflects the connectivity of different regions. In addition, the traffic dynamics are largely influenced and constrained by the static road network. Thus, we use the dynamic OD data and the static road network to enhance the spatial dependence modelling of traffic volume. The challenge herein is how to design an approach that can effectively model these two types of data and fuse them with the traffic volume representation.

To solve this challenge, we propose a hybrid spatial dependence modelling (HSDM) component by simultaneously learning the dynamic OD data and the static road network. For the dynamic graph, we use  $\mathbf{D}_t$  to denote its adjacency matrix and  $\mathbf{H}_t$  its input feature. In this study, we use the in-traffic and the out-traffic of

each node as the features. Therefore,  $\mathbf{H}_t \in \mathbb{R}^{N \times 2}$  is obtained from  $\mathbf{X}_t \in \mathbb{R}^{2 \times I \times J}$  and they have the same elements, excluding the dimensions. For the static graph, the adjacency matrix is  $\mathbf{S}$  and the input feature is  $\mathbf{H} = \sum_t \mathbf{H}_t$ .

As shown in Figure 2, we use  $L_{\text{ORFL}}$  HSDM components to perform representation learning for the dynamic graph and the static graph. Specifically, each HSDM component includes two GAT models, in which the multihead attention scheme is introduced to stabilise the learning process. Figure 4 illustrates the implementation details of our HSDM component and the multihead attention scheme. For the  $l$ -th component, the output can be expressed as follows:

$$\mathbf{H}_t^{(l)} = \parallel_{m=1}^M \sigma(\mathbf{A}_{d,m}^{(l)} \mathbf{H}_t^{(l-1)} \mathbf{W}_{d,m}^{(l)}), \quad (10)$$

where  $\mathbf{A}_{d,m}^{(l)}$  and  $\mathbf{W}_{d,m}^{(l)}$  denote the masked self-attention weights and learnable parameters of the  $m$ -th head in the  $l$ -th HSDM, respectively. For each node pair  $u$  and  $v$  in  $\mathbf{A}_{d,m}^{(l)}$ , the attention weight  $a_{u,v}^{(l)}$  is calculated by the following equation:

$$a_{u,v}^{(l)} = \frac{\exp(e_{u,v}^{(l)})}{\sum_{k \in \mathcal{N}_{(u)}} \exp(e_{u,k}^{(l)})}, \quad (11)$$

$$e_{u,v}^{(l)} = \sigma\left\{\left(\mathbf{h}_u^{(l)} \mathbf{W}_{d,m}^{(l)} \parallel \mathbf{h}_v^{(l)} \mathbf{W}_{d,m}^{(l)}\right) \mathbf{a}^{(l)}\right\},$$

where  $\mathbf{a}^{(l)}$  denotes a learnable parameter vector and  $\parallel$  is the concatenation operator. The core idea here is to compute the hidden representations of each node in the graph by attending to its neighbours, following a self-attention strategy.

The representation learning for the static road network can be obtained similarly, and we describe it as follows:

$$\mathbf{H}^{(l)} = \parallel_{m=1}^M \sigma(\mathbf{A}_{s,m}^{(l)} \mathbf{H}^{(l-1)} \mathbf{W}_{s,m}^{(l)}). \quad (12)$$

where  $\mathbf{A}_{s,m}^{(l)}$  and  $\mathbf{W}_{s,m}^{(l)}$  are the masked self-attention weights and learnable parameters of the  $m$ -th head in the  $l$ -th component of HSDM, respectively.

After we obtain the representations of the dynamic graph and the static graph, we design a gating mechanism by using a mapping function  $\phi: \mathbb{R}^{N \times F} \rightarrow \mathbb{R}^{I \times J \times F}$  and apply it to the traffic volume

representations, and then these two types of representations are fused, resulting in the following feature representation:

$$\mathbf{G}_t^{(l)} = \beta \phi(\mathbf{H}_t^{(l)}) + (1 - \beta) \phi(\mathbf{H}^{(l)}), \quad (13)$$

where  $\beta$  is a predefined parameter that balances the different effects of the dynamic and static spatial dependences on the final representations.  $\mathbf{G}_t^{(l)}$  denotes the learnt representations for each region and reflects the spatial dependencies of traffic volume among different regions. Then,  $\mathbf{G}_t^{(l)}$  is modulated with the traffic feature, and the learnt representation of Equation (9) can be updated as follows:

$$\mathbf{X}_t^{(l)} = \sigma(\mathbf{W}_{\text{TFL}}^{(l)} * \mathbf{X}_t^{(l-1)}) \odot \mathbf{G}_t^{(l)}, \quad (14)$$

where  $\odot$  is the Hadamard product.

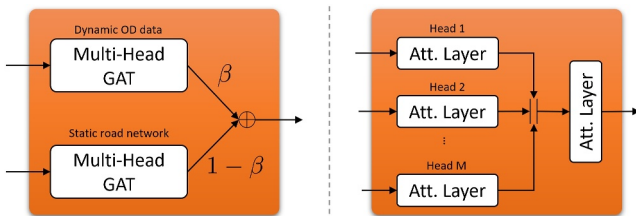
#### 4.5 | Joint Feature Learning

The spatial and temporal correlations of traffic volume among different regions are quite challenging to capture, as they are affected by many factors such as geographic distance and urban layout. Some regions that are far away from each other may have high spatial correlations when they are similar functional urban areas. Thus, a deeper network is preferred to learn the representations of traffic volume and capture the long-distance spatial correlations. We design a JFL component and perform joint feature learning for  $\mathbf{X}_t^{(L_{\text{TFL}})}$ . As shown in Figure 2, we stack learning blocks sequentially in the JFL component, and the design of each block is detailed in Figure 5. The blocks are connected by a convolutional layer, whose purpose is to reduce the dimension of feature maps and, consequently, to reduce the number of parameters in our STOD-Net.

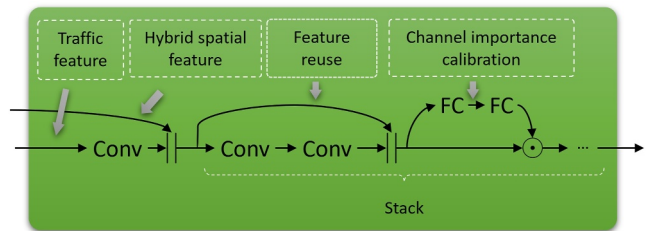
For each block  $b \in \{1, 2, \dots, B\}$ , the input to it can be written as follows:

$$\mathbf{Y}_{t,b} = (\mathbf{Y}_{t,b-1} * \mathbf{W}_{\text{JFL},b}) \parallel \mathbf{G}_t^{(L_{\text{TFL}})}, \quad (15)$$

where  $\mathbf{W}_{\text{JFL},b}$  denotes the parameters in the convolutional layer of block  $b$  in the JFL component and  $\mathbf{Y}_{t,0} = \mathbf{X}_t^{(L_{\text{TFL}})}$ . Note that in Equation (15), we transfer (FL component) the hybrid spatial representations from dynamic OD data and static road network to the JFL module by adding shortcut connections between the last HSDM block and the JFL component. This is because shortcut connections among layers of neural networks are



**FIGURE 4** | Design details of HSDM (left) and multihead GAT (right).



**FIGURE 5** | Design of joint feature learning module.

beneficial for the optimisation of deeper models [53]. They can alleviate gradient vanishing or exploding problems and strengthen feature propagation, thereby improving the prediction performance of neural networks. We use concatenation instead of summation because concatenation can increase the variation in the input of subsequent layers and improve efficiency [54].

We then forward  $\mathbf{Y}_{t,b}$  to the STM module for spatial-temporal modelling. STM consists of  $L_{JFL}$  layers and each layer performs the following equation:

$$\mathbf{Y}_{t,b}^{(l)} = \mathcal{H}_b^{(l)}(\mathbf{Y}_{t,b}^{(l-1)}; \mathbf{W}_{JFL,b}^{(l)}), \quad (16)$$

where  $\mathcal{H}_b^{(l)}$  represents the composite function of layer  $l$  in block  $b$  and it implements six consecutive operations: BN-ReLU-Conv( $1 \times 1$ )-BN-ReLU-Conv( $3 \times 3$ ). Note that  $\mathbf{Y}_{t,b}^{(0)} = \mathbf{Y}_{t,b}$  represents the initial input to block  $b$ .

To further capture channelwise importance and make the learning process flexible, we employ a gating mechanism with a sigmoid activation on  $\mathbf{Y}_{t,b}^{(l)}$ . Particularly, we parameterise the gating mechanism by forming a bottleneck with two fully connected (FC) layers:

$$\mathbf{C}_{t,b}^{(l)} = \text{sigmoid}(\mathbf{W}_{FC,2}^{(l)} \sigma(\mathbf{W}_{FC,1}^{(l)} \mathbf{Y}_{t,b}^{(l)})). \quad (17)$$

Then, the output  $\mathbf{Y}_{t,b}^{(l)}$  is updated as  $\mathbf{Y}_{t,b}^{(l)} = \mathbf{Y}_{t,b}^{(l)} \odot \mathbf{C}_{t,b}^{(l)}$ . After  $B$  blocks' representation learning, the final output of the JFL component is  $\mathbf{Y}_{t,B}$ .

#### 4.6 | Final Fusion

The above subsections summarise the detailed implementations of our STOD. As shown on the left of Figure 2, there exist three inputs, that is,  $\{\mathbf{X}_t^c, \mathbf{D}_t^c\}$ ,  $\{\mathbf{X}_t^p, \mathbf{D}_t^p\}$  and  $\{\mathbf{X}_t^r, \mathbf{D}_t^r\}$ . Thus, three outputs can be obtained:  $\mathbf{Y}_{t,B}^{(c)}$ ,  $\mathbf{Y}_{t,B}^{(p)}$  and  $\mathbf{Y}_{t,B}^{(r)}$ . They denote the learnt spatiotemporal representations for the closeness, periodicity and trend dependencies, respectively. The prediction can be obtained after performing fusion on the three kinds of outputs:

$$\hat{\mathbf{X}}_{t+1} = \text{sigmoid}(\mathbf{Y}_{t,B}^{(c)} + \mathbf{Y}_{t,B}^{(p)} + \mathbf{Y}_{t,B}^{(r)}), \quad (18)$$

where  $\text{sigmoid}(\cdot)$  stands for the sigmoid activation function. We adopt mean squared error (MSE) as our objective function; thus, the loss function in Equation (1) can be detailed as follows:

$$\mathcal{L} = \|\hat{\mathbf{X}}_{t+1} - \mathbf{X}_{t+1}\|_2^2, \quad (19)$$

where  $\|\cdot\|_2^2$  denotes the square of the Frobenius norm. The parameters of our STOD-Net, that is,  $w$ , can be obtained by optimising Equation (19) over all the training dataset.

## 5 | Experimental Results

In this section, we introduce two real-world datasets used in the experiments, evaluation metrics and several baseline methods to

which we compared. We also explain the detailed experimental settings, followed by thorough experimental results and the corresponding analysis.

### 5.1 | Datasets and Preprocessing

The two real-world datasets come from New York City (NYC), which are publicly available [55]. They are taxi and bike trip records, and we denote them as NYC-Taxi and NYC-Bike, respectively. NYC-Taxi includes approximately 22.3 million trip records from 1 January 2015 to 1 March 2015, whereas the NYC-Bike dataset includes about 2.2 million trip records from 1 July 2016 to 29 August 2016.

The NYC is split into  $10 \times 20$  regions, and each region's area is approximately  $1 \text{ km}^2$ . We set the time interval as 30 min and calculate the traffic volume in each region and the OD data from one region to another. We use the first 50 days of data to train our model and the last 10 days of data to test the model's performance. Before training our model, both the traffic volume data and the OD data are scaled into  $[0, 1]$  by using min-max normalisation and are rescaled back to their original scale during evaluation. Besides, traffic values that are less than a certain threshold are ignored when evaluating performance because low traffic is of little interest, and this assumption is also commonly adopted in industry and academia [31]. We follow the same settings in Refs. [9, 31] and set the threshold to 10.

### 5.2 | Baselines and Evaluation Metrics

We compare our model, STOD-Net, with the following several baselines:

- Historical average (HA) method. The next time slot prediction is set to the average of all historical traffic values.
- Naive method. This method treats the last observation as the prediction of the next time slot.
- Autoregressive integrated moving average (ARIMA) model. ARIMA is one of the most frequently used time series prediction approaches and can model the autocorrelations in traffic data.
- Linear regression (LR). LR denotes a linear approach to traffic prediction.
- Multilayer perceptron (MLP). MLP represents a class of feedforward artificial neural networks. MLP can be used to model the nonlinear relationships hidden in the traffic data.
- Spatiotemporal residual network (ST-ResNet) [24]. ST-ResNet is a deep learning framework for solving citywise traffic prediction and has become a seminal work since its publication. It can model the spatial and temporal dynamics of traffic volume simultaneously.
- Spatial-temporal graph convolutional networks (STGCN) [4]. Instead of applying regular convolutional units, STGCN formulates the traffic prediction on graphs and builds the



model with complete convolutional structures, which enables a much faster training speed with fewer parameters.

- Spatiotemporal dynamic network (STDN) [31]. STDN can model the different spatial dynamics among different locations and the temporal shifting in traffic data. It achieves state-of-the-art performance on both the NYC-Taxi and NYC-Bike datasets.

We adopt three metrics to evaluate the prediction performance of different approaches. The metrics are root mean square error (RMSE), mean absolute error (MAE) and mean absolute percentage error (MAPE), respectively.

$$\text{RMSE} = \sqrt{\frac{1}{M} \sum_{t=1}^{T'} \sum_{i=1}^I \sum_{j=1}^J (x_t^{(i,j)} - \hat{x}_t(i,j))^2}, \quad (20)$$

$$\text{MAE} = \frac{1}{M} \sum_{t=1}^{T'} \sum_{i=1}^I \sum_{j=1}^J |x_t^{(i,j)} - \hat{x}_t(i,j)|, \quad (21)$$

$$\text{MAPE} = \frac{1}{M} \sum_{t=1}^{T'} \sum_{i=1}^I \sum_{j=1}^J \left| \frac{x_t^{(i,j)} - \hat{x}_t(i,j)}{x_t^{(i,j)}} \right| \quad (22)$$

where  $T'$  denotes the time slots in the test dataset and  $M = T' \times I \times J$ . In addition,  $x_t^{(i,j)}$  ( $\hat{x}_t(i,j)$ ) denotes the real (predicted) in-/out-traffic value at time slot  $t$  for region  $(i,j)$ .

### 5.3 | Experimental Settings

The experimental results are obtained with the following settings. For the lengths of closeness, periodicity and trend sequences, we set them to 5, 4 and 1, and these values are obtained based on a grid search over  $L_c, L_p, L_r \in \{1, 2, 3, 4, 5\}$ . When constructing the static road network, we set  $\rho_z = \{0.4, 0.3, 0.2, 0.1\}$ . In the TFL component, we use two convolutional layers for preliminary traffic representation learning, and the hidden representations in each layer are 24. In the ORFL component, there are two HSDM blocks, and each block has two  $M$ -heads GAT for the static and dynamic spatial dependency modelling, and  $M$  equals 2. The parameter,  $\beta$ , that balances the different effects of the static and dynamic spatial dependences, is set to 0.5. In the JFL component, there are 3 blocks. Each block has  $L_{\text{JFL}} = 8$  layers, and each layer outputs  $K = 12$  hidden representations.

We use Adam to optimise our STOD-Net with a learning rate of 0.0001, and the learning rate is adjusted via the one-cycle policy. In addition, STOD-Net is trained using batch 128 for 500 epochs on the NYC-Taxi dataset and 100 epochs on the NYC-Bike dataset.

### 5.4 | Citywise Prediction Results

Table 1 summarises the citywise prediction performance comparisons of different methods on the NYC-Taxi and NYC-Bike datasets. To obtain Table 1, we carry out 10 independent experiments and report the mean and standard deviation ( $\pm$ )

results of three evaluation metrics. The best results are marked in bold for clearness.

From this table, we can clearly observe that (1) deep neural network-based models can normally achieve better prediction performance than linear models (LR and ARIMA), followed by simple forecasting methods (HA and Naive). For instance, the RMSE result on the in-traffic of the NYC-Taxi dataset improves from 36.96 (Naive) to 30.09 (MLP) to 22.32 (STDN), which validates the advantages of deep neural networks in modelling spatiotemporal dependencies of traffic volume. (2) Our proposed method, that is, STOD-Net, achieves the best prediction performance in terms of RMSE, MAE and MAPE, on both datasets. We take the MAPE metric on the NYC-Taxi dataset as an example to illustrate the effectiveness of STOD-Net. Specifically, the obtained MAPE results of STOD-Net for the in-traffic (out-traffic) are about 15.28% (15.33%), whereas the best baseline's (STDN) MAPE results are about 16.15% (16.13%). These numbers indicate that approximately a 5.4% (5.0%) performance gain can be acquired by STOD-Net. For the other two metrics, that is, RMSE and MAE, similar performance gains can be obtained as well if we do the calculation, but we omit their specific values here for brevity.

Aside from lower prediction errors, STOD-Net remains stable across different experiments, as it yields lower standard deviations than baselines on these three metrics. Consider STDN as an example: the standard deviations of MAPE on the NYC-Taxi dataset are 0.62 (in-traffic) and 0.52 (out-traffic), respectively. In comparison, STOD-Net's standard deviations of MAPE on the same dataset are approximately 0.09 (in-traffic) and 0.10 (out-traffic), respectively, which are much lower than those of STDN.

We also analyse the prediction errors in Figure 6, which shows the histograms of the absolute prediction errors (APEs) of three methods, that is, ST-ResNet, STDN and STOD-Net. Figure 6 left/right shows the results on the NYC-Taxi/NYC-Bike dataset. From this figure, we can see that for STOD-Net, a larger portion of prediction errors tend to be zero compared with ST-ResNet and STDN. The results of Figure 6 demonstrate that STOD-Net has lower prediction errors than its competitive counterparts.

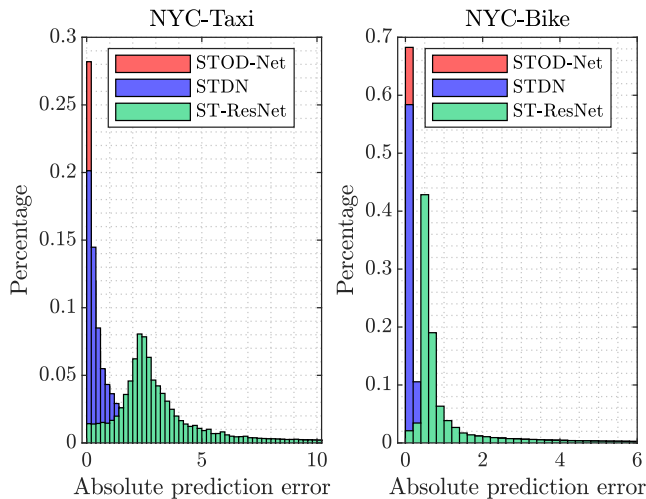
The prediction results of Table 1 and Figure 6 confirm the effectiveness of STOD-Net. Thus, we claim it yields better prediction performance than baselines. The reasons behind the success of STOD-Net can be attributed to twofold: (1) The introduced ORFL component models the hidden relationships of different regions, leading to an accurate characterisation of their spatial dependencies. (2) The complex spatial and temporal dependencies among regions are simultaneously captured and learnt by the JFL component; thus, more effective representations can be obtained to enhance STOD-Net's prediction ability.

### 5.5 | Regionwise Prediction Results

The above subsection reports the overall quantitative prediction results for all regions. In this subsection, we compare the regionwise prediction performance of different methods. To be specific, we randomly select several regions and plot the ground

**TABLE 1** | Prediction performance of different methods on the NYC-Taxi and NYC-Bike datasets.

Data	Method	In-traffic			Out-traffic		
		RMSE	MAE	MAPE	RMSE	MAE	MAPE
NYC-Taxi	HA	71.02	41.10	38.06%	59.90	32.55	36.23%
	Naive	36.96	22.72	22.94%	31.78	18.32	22.96%
	ARIMA	34.92	21.97	24.85%	29.99	18.12	25.26%
	LR	30.55	18.93	19.83%	25.66	15.12	19.66%
	MLP	30.09 $\pm$ 0.21	18.57 $\pm$ 0.13	19.97 $\pm$ 0.18%	24.69 $\pm$ 0.24	14.31 $\pm$ 0.12	19.06 $\pm$ 0.16%
	ST-ResNet	23.89 $\pm$ 0.16	15.25 $\pm$ 0.07	17.16 $\pm$ 0.07%	19.47 $\pm$ 0.09	12.14 $\pm$ 0.06	16.67 $\pm$ 0.07%
	STGCN	22.78 $\pm$ 0.20	14.29 $\pm$ 0.15	16.67 $\pm$ 0.31%	18.52 $\pm$ 0.15	11.54 $\pm$ 0.13	16.56 $\pm$ 0.31%
	STDN	22.32 $\pm$ 0.22	14.09 $\pm$ 0.18	16.15 $\pm$ 0.62%	18.08 $\pm$ 0.27	11.38 $\pm$ 0.20	16.13 $\pm$ 0.52%
	STOD-Net	<b>21.44 <math>\pm</math> 0.08</b>	<b>13.37 <math>\pm</math> 0.04</b>	<b>15.28 <math>\pm</math> 0.09%</b>	<b>17.61 <math>\pm</math> 0.08</b>	<b>10.87 <math>\pm</math> 0.05</b>	<b>15.33 <math>\pm</math> 0.10%</b>
NYC-Bike	HA	17.46	11.02	37.31%	16.72	10.69	35.54%
	Naive	14.03	9.48	31.25%	13.43	9.28	30.62%
	ARIMA	12.92	8.81	28.59%	12.38	8.60	27.84%
	LR	11.89	8.07	26.76%	11.21	7.74	25.69%
	MLP	9.41 $\pm$ 0.04	6.54 $\pm$ 0.02	23.05 $\pm$ 0.10%	8.54 $\pm$ 0.06	6.12 $\pm$ 0.03	21.71 $\pm$ 0.15%
	ST-ResNet	8.96 $\pm$ 0.03	6.46 $\pm$ 0.02	22.72 $\pm$ 0.06%	8.19 $\pm$ 0.04	6.08 $\pm$ 0.03	21.49 $\pm$ 0.09%
	STGCN	8.83 $\pm$ 0.18	6.35 $\pm$ 0.12	22.42 $\pm$ 0.45%	7.89 $\pm$ 0.14	5.86 $\pm$ 0.10	20.76 $\pm$ 0.37%
	STDN	8.61 $\pm$ 0.18	6.14 $\pm$ 0.13	21.42 $\pm$ 0.22%	7.78 $\pm$ 0.18	5.73 $\pm$ 0.12	20.15 $\pm$ 0.31%
	STOD-Net	<b>8.18 <math>\pm</math> 0.03</b>	<b>5.87 <math>\pm</math> 0.02</b>	<b>20.63 <math>\pm</math> 0.05%</b>	<b>7.39 <math>\pm</math> 0.03</b>	<b>5.48 <math>\pm</math> 0.02</b>	<b>19.39 <math>\pm</math> 0.03%</b>

**FIGURE 6** | Histograms of absolute prediction errors.

truth values and the achieved predictions by different methods for these regions. In the following, we only compare STOD-Net with ST-ResNet and STDN, as they generally obtain better performance than other baselines. Figures 7 and 8 show the results of three randomly selected regions on the NYC-Taxi and NYC-Bike datasets, respectively. In Figures 7 and 8, the left three subfigures compare the predictions versus ground truth values, the middle three subfigures show the distributions of APEs and the right three subfigures illustrate the cumulative distribution functions (CDFs) of APEs and the mean prediction errors.

From these two figures, we can observe that (1) the predictions of STOD-Net and STDN, as well as ST-ResNet, all match the ground truth values fairly well, although distinct traffic patterns exist in

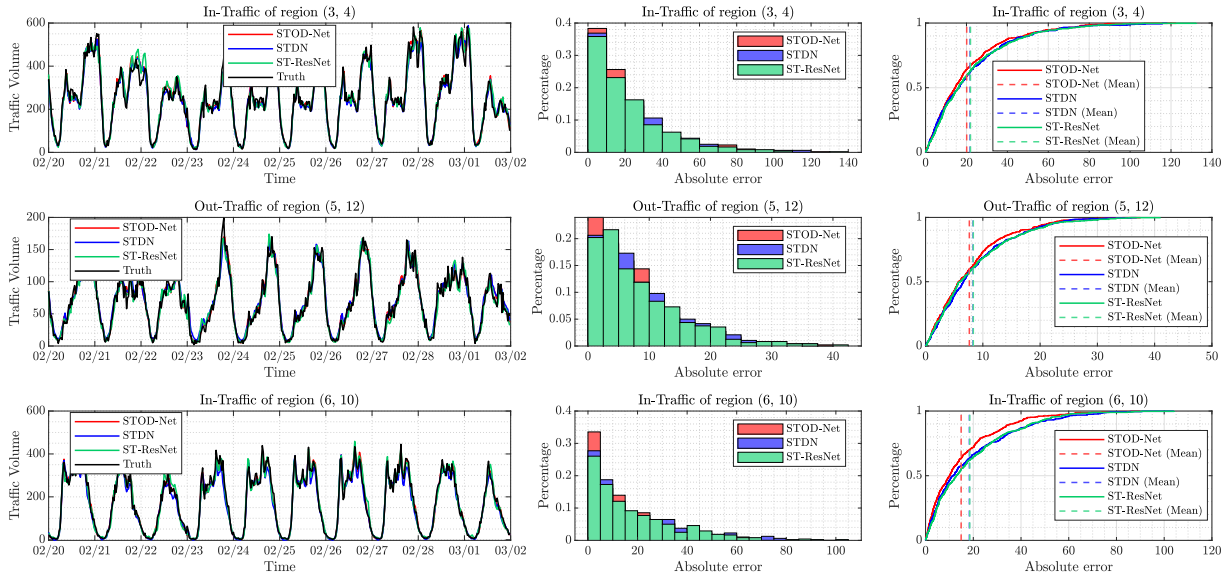
different regions or datasets. For instance, on the NYC-Taxi dataset, the traffic pattern of the region (6, 10) is different from that of the region (5, 12), whose peak traffic has a clear delay. Besides, region (6, 10) has multiple peak periods, whereas region (5, 12) tends to have a single peak period for most of the days.

Nevertheless, all three methods achieve competitive prediction results, and this validates the superiority of convolutional neural networks for urban traffic prediction. (2) The STOD-Net method achieves better regionwise prediction results than its two counterparts, which can be verified by the histograms and CDFs of APEs. Take the region (6, 10) on the NYC-Taxi dataset as an example. From the histograms of APEs, we can perceive that STOD-Net has more prediction errors tending to zero compared with STDN and ST-ResNet. This can be more quantitatively reflected by the CDFs of APEs. More specifically, about 72% of APEs are less than 20 for STOD-Net, whereas they are 66% and 64% for STDN and ST-ResNet, respectively. The mean values of APEs for STOD-Net, STDN and ST-ResNet are 14.91, 18.35 and 18.50, respectively.

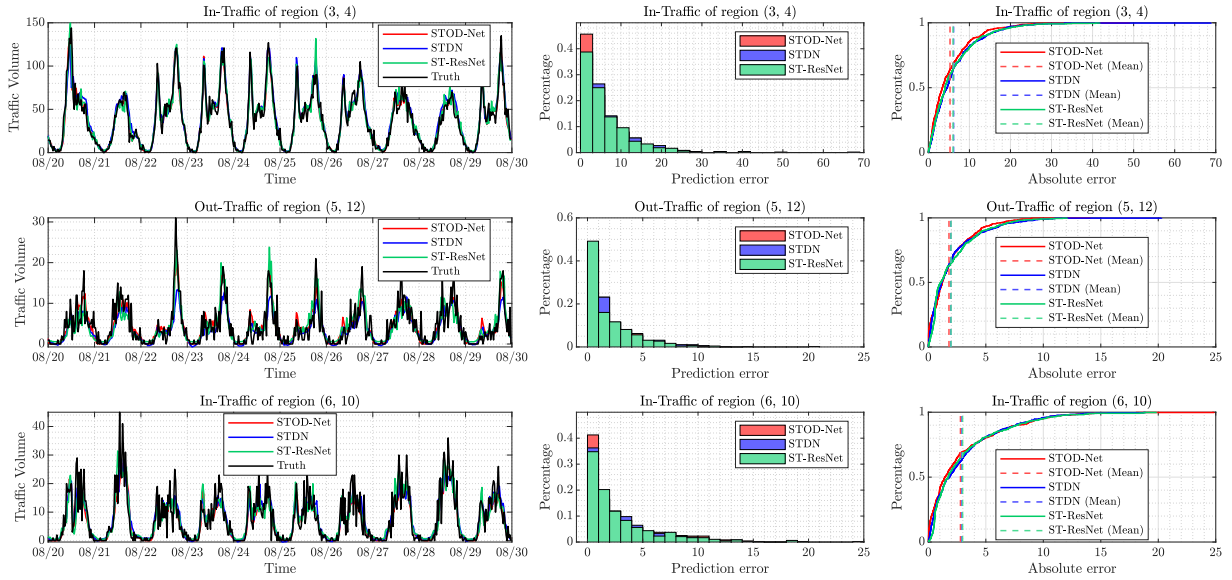
Similar results can be obtained for the other regions on the NYC-Taxi and NYC-Bike datasets. Thus, we can conclude that our proposed STOD-Net achieves better regionwise prediction results than baselines.

## 5.6 | Impacts of OD and Road Feature Learning

In this subsection, we go one step further and analyse the impacts of OD and road feature learning on the prediction performance of STOD-Net. In particular, we report the prediction results with



**FIGURE 7** | Predictions versus ground truth and the corresponding prediction error analysis for three randomly selected regions, that is, (3, 4), (5, 12) and (6, 10), on the NYC-Taxi dataset.



**FIGURE 8** | Predictions versus ground truth and the corresponding prediction error analysis for three randomly selected regions, that is, (3, 4), (5, 12) and (6, 10), on the NYC-Bike dataset.

and without OD and road feature learning to demonstrate the advantages of introducing OD feature learning for urban traffic prediction. Table 2 summarises the obtained overall results on the NYC-Taxi and NYC-Bike datasets, without distinguishing the in-traffic and out-traffic. In this table, there are three versions of STOD-Net, that is, Basic version, Basic + ORFL version and Basic + ORFL + FT version. The ‘Basic’ denotes that we only use the TFL and JFL components for traffic prediction without the OD feature. The ‘Basic + ORFL’ means the ORFL component is further introduced into traffic prediction. The ‘Basic + ORFL + FT’ denotes that both the ORFL and the FT components are introduced into traffic prediction.

From Table 2, we know that the ‘Basic’ version of STOD-Net performs the worst among all the three versions. With the

**TABLE 2** | Effectiveness of OD learning on prediction performance.

Dataset	Method	MSE	MAE	MAPE
NYC-Taxi	Basic	19.57	12.10	15.28%
	Basic + ORFL	19.33	11.92	15.09%
	Basic + ORFL + FT	19.22	11.90	15.08%
NYC-Bike	Basic	8.16	5.93	20.79%
	Basic + ORFL	8.01	5.83	20.48%
	Basic + ORFL + FT	7.88	5.73	19.97%

introduction of ORFL, the prediction errors on both datasets become lower. Besides, the ORFL component further improves prediction performance. The results in Table 2 indicate that

the OD data can indeed provide additional performance gains and serve as a new direction to consider for urban traffic prediction.

## 5.7 | Impacts of Hyperparameters

In this subsection, we explore the impacts of several related hyperparameters on the prediction performance of STOD-Net. For example, the number of HSDM blocks and the parameter  $\beta$  that balances the different effects of the dynamic and static spatial dependencies.

### 5.7.1 | The Number of HSDM Blocks

We vary the number of HSDM blocks  $L_{ORFL} \in \{0, 1, 2, 3\}$  and present the obtained RMSE and MAE results on the NYC-Taxi and NYC-Bike datasets in Figure 9. Note that when  $L_{ORFL} = 0$ , STOD-Net degenerates into the 'Basic' version, which is explained in the last subsection. We can observe from Figure 9 that with the increase of  $L_{ORFL}$ , the performances regarding RMSE and MAE improve gradually until they reach their corresponding minimums, after which the performances decrease. The results of Figure 9 are consistent with the results of Table 2 and demonstrate again the usefulness of OD feature learning. Besides, based on the results, we can conclude that the number of HSDM blocks should not be too large, as larger  $L_{ORFL}$  may break the TFL component's representation learning, thus leading to performance degradation.

### 5.7.2 | The Parameter $\beta$

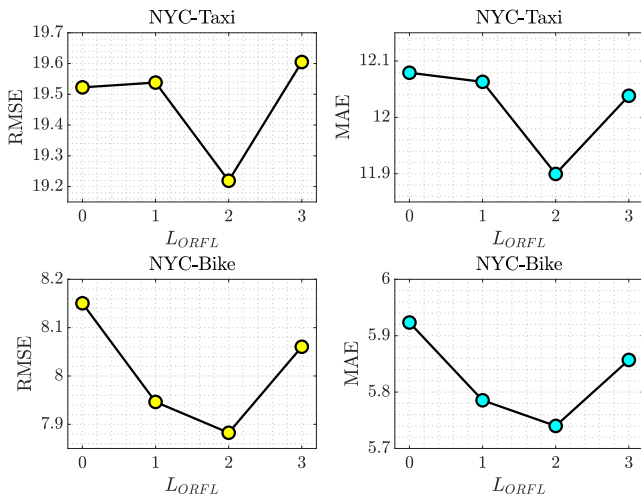
The parameter  $\beta$  in Equation (13) is a predefined value that affects the gating mechanism. It balances the contribution of the dynamic and static spatial dependencies to the final representations. We set  $\beta \in \{0, 0.5, 1\}$  and report the obtained results in Figure 10. Note that when  $\beta = 0$ , it means we only consider the

static spatial dependence. Similarly, when  $\beta = 1$ , it means we only consider the dynamic spatial dependence. When  $\beta = 0.5$ , we consider both the static and dynamic spatial dependencies, and they contribute equally to the final representations. We can see that the prediction performances vary when the value of  $\beta$  changes. Nonetheless, the gating mechanism is beneficial to predictions regardless of its values, as it achieves better performance than the 'Basic' version of STOD-Net. Additionally, we observe that considering both the static and dynamic spatial dependencies, that is,  $\beta = 0.5$ , can generally obtain lower prediction errors than only considering the static or dynamic spatial dependence. Because the static and dynamic OD data contain complementary information for capturing spatial dependencies, thus modelling them together enhances the prediction performance.

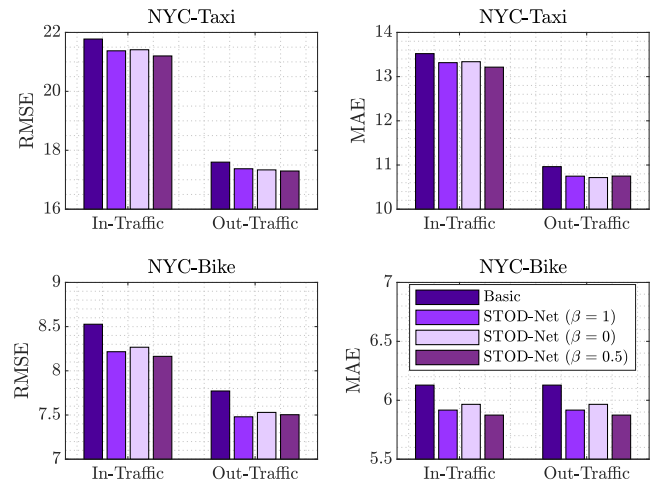
### 5.7.3 | Others

There are several other hyperparameters that influence the prediction performance of STOD-Net, and we give the experimental results here. Six hyperparameters are considered, that is, the length of closeness dependence  $L_c$ , the length of periodicity dependence  $L_p$ , the length of trend dependence  $L_r$ , the number of layers in each dense block  $L_{JFL}$ , the output representations in each convolutional layer of the dense block  $K$  and the number of heads in GAT  $M$ . Their influence on prediction results is illustrated in Figure 11. Note that we only report the RMSE results on one dataset (NYC-Bike) as other metrics have similar results to Figure 11. It can be observed from Figure 11 that with the increase of  $L_c$  and  $L_p$ , the prediction performance tends to improve because more data are used to model temporal dependence. However, for  $L_r$ , as it increases, the performance degrades gradually. This is because the increase of  $L_r$  significantly reduces the number of training samples; thus, under-fitting may occur.

For the number of layers in each dense block  $L_{JFL}$ , STOD-Net has a large capacity as it increases, leading to performance

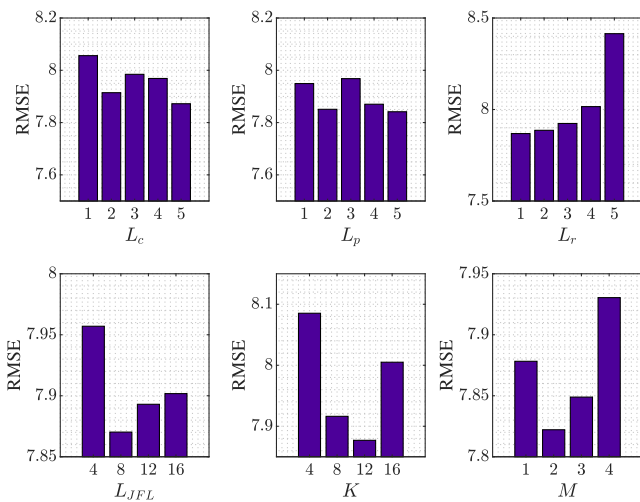


**FIGURE 9** | The number of HSDM blocks versus the prediction performance.



**FIGURE 10** | The influence of gate mechanism on prediction performance.





**FIGURE 11** | The influence of hyperparameters on prediction performance.

improvement. However, too large a  $L_{JFL}$  will make our model prone to overfit the data, thus degrading the performance. This phenomenon holds for the hyperparameters  $K$  and  $M$ . These hyperparameters control the capacity of STOD-Net. Larger values indicate large capacity and strong representation ability. However, excessively large values can easily cause overfitting and reduce the generalisation of STOD-Net. Thus, the best hyperparameters can be obtained based on a grid search strategy.

## 6 | Conclusion

In this paper, we proposed STOD-Net, a spatial-temporal origin-destination feature-enhanced deep neural network, to solve urban traffic prediction. Beyond modelling the historical traffic itself, we introduced OD data into the prediction and adopted graph neural networks to model them to capture the inter-regional spatial interaction patterns. We consider two types of OD data, namely, static and dynamic OD data, in STOD-Net and fuse them in a weighted fashion to capture different regional spatial dependences. Extensive experiments were conducted on two real-world datasets, and the results demonstrated that STOD-Net achieves superior prediction performance to state-of-the-art methods. Possible future directions include considering more fine-grained road network information and exploring the trade-off between the number of parameters and the prediction performance.

### Conflicts of Interest

The authors declare no conflicts of interest.

### Data Availability Statement

Research data are available upon request.

### Endnotes

<sup>1</sup> Time slot  $t$  denotes a time interval  $(t, t + \Delta t]$ , where  $\Delta t$  is the temporal granularity of the data.  $\Delta t$  could be 15 min or 1 h based on the data on hand.

## References

1. X. Fan, C. Xiang, C. Chen, et al., "BuildSenSys: Reusing Building Sensing Data for Traffic Prediction With Cross-Domain Learning," *IEEE Transactions on Mobile Computing* 20, no. 6 (2021): 2154–2171, <https://doi.org/10.1109/tmc.2020.2976936>.
2. D. A. Tedjopurnomo, Z. Bao, B. Zheng, F. M. Choudhury, and A. K. Qin, "A Survey on Modern Deep Neural Network for Traffic Prediction: Trends, Methods and Challenges," *IEEE Transactions on Knowledge and Data Engineering* 34, no. 4 (2022): 1544–1561, <https://doi.org/10.1109/TKDE.2020.3001195>.
3. Y. Zheng, L. Capra, O. Wolfson, and H. Yang, "Urban Computing: Concepts, Methodologies, and Applications," *ACM Transactions on Intelligent Systems and Technology* 5, no. 3 (2014): 1–55, <https://doi.org/10.1145/2629592>.
4. B. Yu, H. Yin, and Z. Zhu, "Spatio-Temporal Graph Convolutional Networks: A Deep Learning Framework for Traffic Forecasting," in *Proceedings of the 27th International Joint Conference on Artificial Intelligence* (2018), 3634–3640.
5. I. AlQerm and B. Shihada, "Energy Efficient Traffic Offloading in Multi-Tier Heterogeneous 5G Networks Using Intuitive Online Reinforcement Learning," *IEEE Transactions on Green Communications and Networking* 3, no. 3 (2019): 691–702, <https://doi.org/10.1109/TGCN.2019.2916900>.
6. I. AlQerm and B. Shihada, "A Cooperative Online Learning Scheme for Resource Allocation in 5G Systems," in *2016 IEEE International Conference on Communications (ICC)* (2016), 1–7, <https://doi.org/10.1109/ICC.2016.7511617>.
7. A. Elwhishi, P. H. Ho, K. Naik, and B. Shihada, "Self-Adaptive Contention Aware Routing Protocol for Intermittently Connected Mobile Networks," *IEEE Transactions on Parallel and Distributed Systems* 24, no. 7 (2013): 1422–1435, <https://doi.org/10.1109/TPDS.2012.23>.
8. L. Zhang, C. Zhang, and B. Shihada, "Efficient Wireless Traffic Prediction at the Edge: A Federated Meta-Learning Approach," *IEEE Communications Letters* 26, no. 7 (2022): 1573–1577, <https://doi.org/10.1109/LCOMM.2022.3167813>.
9. H. Yao, F. Wu, J. Ke, et al., "Deep Multi-View Spatial-Temporal Network for Taxi Demand Prediction," *Proceedings of the AAAI Conference on Artificial Intelligence* 32, no. 1 (2018): 2588–2595, <https://doi.org/10.1609/aaai.v32i1.11836>.
10. C. Zhang, S. Dang, B. Shihada, and M. S. Alouini, "Dual Attention-Based Federated Learning for Wireless Traffic Prediction," in *IEEE INFOCOM 2021 - IEEE Conference on Computer Communications* (2021), 1–10, <https://doi.org/10.1109/INFOCOM42981.2021.9488883>.
11. D. Zhang and X. Feng, "Dynamic Auto-Structuring Graph Neural Network: A Joint Learning Framework for Origin-Destination Demand Prediction," *IEEE Transactions on Knowledge and Data Engineering* 35, no. 4 (2023): 3699–3711, <https://doi.org/10.1109/tkde.2021.3135898>.
12. T. Qi, G. Li, L. Chen, and Y. Xue, "ADGCN: An Asynchronous Dilation Graph Convolutional Network for Traffic Flow Prediction," *IEEE Internet of Things Journal* 9, no. 5 (2022): 4001–4014, <https://doi.org/10.1109/jiot.2021.3102238>.
13. S. Guo, Y. Lin, N. Feng, C. Song, and H. Wan, "Attention Based Spatial-Temporal Graph Convolutional Networks for Traffic Flow Forecasting," *Proceedings of the AAAI Conference on Artificial Intelligence* 33, no. 1 (2019): 922–929, <https://doi.org/10.1609/aaai.v33i01.3301922>.
14. Z. Pan, Y. Liang, W. Wang, Y. Yu, Y. Zheng, and J. Zhang, "Urban Traffic Prediction From Spatio-Temporal Data Using Deep Meta Learning," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (2019), 1720–1730.

15. N. G. Polson and V. O. Sokolov, "Deep Learning for Short-Term Traffic Flow Prediction," *Transportation Research Part C: Emerging Technologies* 79 (2017): 1–17, <https://doi.org/10.1016/j.trc.2017.02.024>.
16. S. V. Kumar, "Traffic Flow Prediction Using Kalman Filtering Technique," *Procedia Engineering* 187 (2017): 582–587, <https://doi.org/10.1016/j.proeng.2017.04.417>.
17. H. Zare Moayed and M. A. Masnadi-Shirazi, "ARIMA Model for Network Traffic Prediction and Anomaly Detection," in *2008 International Symposium on Information Technology*, Vol. 4 (2008), 1–6, <https://doi.org/10.1109/ITSIM.2008.4631947>.
18. Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning," *Nature* 521, no. 7553 (2015): 436–444, <https://doi.org/10.1038/nature14539>.
19. Y. Lv, Y. Duan, W. Kang, Z. Li, and F. Wang, "Traffic Flow Prediction With Big Data: A Deep Learning Approach," *IEEE Transactions on Intelligent Transportation Systems* 16, no. 2 (2015): 865–873.
20. C. Zheng, X. Fan, C. Wang, and J. Qi, "GMAN: A Graph Multi-Attention Network for Traffic Prediction," *Proceedings of the AAAI Conference on Artificial Intelligence* 34, no. 1 (2020): 1234–1241, <https://doi.org/10.1609/aaai.v34i01.5477>.
21. J. Zhang, Y. Zheng, J. Sun, and D. Qi, "Flow Prediction in Spatio-Temporal Networks Based on Multitask Deep Learning," *IEEE Transactions on Knowledge and Data Engineering* 32, no. 3 (2020): 468–478, <https://doi.org/10.1109/TKDE.2019.2891537>.
22. C. Zhang, H. Zhang, J. Qiao, D. Yuan, and M. Zhang, "Deep Transfer Learning for Intelligent Cellular Traffic Prediction Based on Cross-Domain Big Data," *IEEE Journal on Selected Areas in Communications* 37, no. 6 (2019): 1389–1401, <https://doi.org/10.1109/jsac.2019.2904363>.
23. J. Zhang, Y. Zheng, D. Qi, R. Li, and X. Yi, "DNN-Based Prediction Model for Spatio-Temporal Data," in *Proceedings of the 24th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems* (2016), 1–4, <https://doi.org/10.1145/2996913.2997016>.
24. J. Zhang, Y. Zheng, and D. Qi, "Deep Spatio-Temporal Residual Networks for Citywide Crowd Flows Prediction," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence* (2017), 1655–1661.
25. P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph Attention Networks," preprint, arXiv, arXiv:1710.10903 (2017).
26. J. Chen, T. Ma, and C. Xiao, "FastGCN: Fast Learning With Graph Convolutional Networks via Importance Sampling," preprint, arXiv, arXiv:1801.10247 (2018).
27. S. Qi, W. Wang, B. Jia, J. Shen, and S. C. Zhu, "Learning Human-Object Interactions by Graph Parsing Neural Networks," in *Proceedings of the European Conference on Computer Vision (ECCV)* (2018).
28. W. Wang, X. Lu, J. Shen, D. J. Crandall, and L. Shao, "Zero-Shot Video Object Segmentation via Attentive Graph Neural Networks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019).
29. L. Fan, W. Wang, S. Huang, X. Tang, and S. C. Zhu, "Understanding Human Gaze Communication by Spatio-Temporal Graph Reasoning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (2019).
30. M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, "Geometric Deep Learning: Going Beyond Euclidean Data," *IEEE Signal Processing Magazine* 34, no. 4 (2017): 18–42, <https://doi.org/10.1109/msp.2017.2693418>.
31. H. Yao, X. Tang, H. Wei, G. Zheng, and Z. Li, "Revisiting Spatial-Temporal Similarity: A Deep Learning Framework for Traffic Prediction," *Proceedings of the AAAI Conference on Artificial Intelligence* 33, no. 1 (2019): 5668–5675, <https://doi.org/10.1609/aaai.v33i01.33015668>.
32. S. Wang, M. Zhang, H. Miao, Z. Peng, and P. S. Yu, "Multivariate Correlation-Aware Spatio-Temporal Graph Convolutional Networks for Multi-Scale Traffic Prediction," *ACM Transactions on Intelligent Systems and Technology* 13, no. 3 (2022): 1–22, <https://doi.org/10.1145/3469087>.
33. Z. Cui, K. Henrickson, R. Ke, and Y. Wang, "Traffic Graph Convolutional Recurrent Neural Network: A Deep Learning Framework for Network-Scale Traffic Learning and Forecasting," *IEEE Transactions on Intelligent Transportation Systems* 21, no. 11 (2020): 4883–4894, <https://doi.org/10.1109/tits.2019.2950416>.
34. M. S. Ahmed and A. R. Cook, "Analysis of Freeway Traffic Time-Series Data by Using Box-Jenkins Techniques," *Transportation Research Record* 722 (1979).
35. B. M. Williams, "Multivariate Vehicular Traffic Flow Prediction: Evaluation of ARIMAX Modeling," *Transportation Research Record* 1776, no. 1 (2001): 194–200, <https://doi.org/10.3141/1776-25>.
36. S. Lee and D. B. Fambro, "Application of Subset Autoregressive Integrated Moving Average Model for Short-Term Freeway Traffic Volume Forecasting," *Transportation Research Record* 1678, no. 1 (1999): 179–188, <https://doi.org/10.3141/1678-22>.
37. X. Dai, R. Fu, E. Zhao, et al., "DeepTrend 2.0: A Light-Weighted Multi-Scale Traffic Prediction Model Using Detrending," *Transportation Research Part C: Emerging Technologies* 103 (2019): 142–157, <https://doi.org/10.1016/j.trc.2019.03.022>.
38. C. H. Wu, J. M. Ho, and D. Lee, "Travel-Time Prediction With Support Vector Regression," *IEEE Transactions on Intelligent Transportation Systems* 5, no. 4 (2004): 276–281, <https://doi.org/10.1109/TITS.2004.837813>.
39. L. Zhang, Q. Liu, W. Yang, N. Wei, and D. Dong, "An Improved K-Nearest Neighbor Model for Short-Term Traffic Flow Prediction," *Procedia - Social and Behavioral Sciences* 96 (2013): 653–662, <https://doi.org/10.1016/j.sbspro.2013.08.076>.
40. P. Cai, Y. Wang, G. Lu, P. Chen, C. Ding, and J. Sun, "A Spatio-temporal Correlative K-Nearest Neighbor Model for Short-Term Traffic Multistep Forecasting," *Transportation Research Part C: Emerging Technologies* 62 (2016): 21–34, <https://doi.org/10.1016/j.trc.2015.11.002>.
41. G. Leshem and Y. Ritov, "Traffic Flow Prediction Using AdaBoost Algorithm With Random Forests as a Weak Learner," *Proceedings of World Academy of Science, Engineering and Technology* 19 (2007): 193–198.
42. R. Fu, Z. Zhang, and L. Li, "Using LSTM and GRU Neural Network Methods for Traffic Flow Prediction," in *2016 31st Youth Academic Annual Conference of Chinese Association of Automation (YAC)* (2016), 324–328.
43. Z. Cui, R. Ke, Z. Pu, and Y. Wang, "Deep Bidirectional and Unidirectional LSTM Recurrent Neural Network for Network-Wide Traffic Speed Prediction," preprint, arXiv, arXiv:1801.02143 (2018): 1–11.
44. B. Yang, S. Sun, J. Li, X. Lin, and Y. Tian, "Traffic Flow Prediction Using LSTM With Feature Enhancement," *Neurocomputing* 332 (2019): 320–327, <https://doi.org/10.1016/j.neucom.2018.12.016>.
45. L. Zhao, Y. Song, C. Zhang, et al., "T-GCN: A Temporal Graph Convolutional Network for Traffic Prediction," *IEEE Transactions on Intelligent Transportation Systems* 21, no. 9 (2020): 3848–3858, <https://doi.org/10.1109/TITS.2019.2935152>.
46. B. Du, X. Hu, L. Sun, J. Liu, Y. Qiao, and W. Lv, "Traffic Demand Prediction Based on Dynamic Transition Convolutional Neural Network," *IEEE Transactions on Intelligent Transportation Systems* 22, no. 2 (2021): 1237–1247, <https://doi.org/10.1109/tits.2020.2966498>.
47. X. Yao, Y. Gao, D. Zhu, E. Manley, J. Wang, and Y. Liu, "Spatial Origin-Destination Flow Imputation Using Graph Convolutional Networks," *IEEE Transactions on Intelligent Transportation Systems* 22, no. 12 (2021): 7474–7484, <https://doi.org/10.1109/tits.2020.3003310>.

48. Z. Lin, J. Feng, Z. Lu, Y. Li, and D. Jin, "DeepSTN+: Context-Aware Spatial-Temporal Neural Network for Crowd Flow Prediction in Metropolis," *Proceedings of the AAAI Conference on Artificial Intelligence* 33, no. 1 (2019): 1020–1027, <https://doi.org/10.1609/aaai.v33i01.33011020>.
49. T. Chen, L. Nie, J. Pan, L. Tu, B. Zheng, and X. Bai, "Origin-Destination Traffic Prediction Based on Hybrid Spatio-Temporal Network," in *2022 IEEE International Conference on Data Mining (ICDM)* (2022), 879–884, <https://doi.org/10.1109/ICDM54844.2022.00101>.
50. F. Li, J. Feng, H. Yan, et al., "Dynamic Graph Convolutional Recurrent Network for Traffic Prediction: Benchmark and Solution," *ACM Transactions on Knowledge Discovery From Data* 17, no. 1 (2023): 1–21, <https://doi.org/10.1145/3532611>.
51. S. Wang, H. Miao, H. Chen, and Z. Huang, "Multi-Task Adversarial Spatial-Temporal Networks for Crowd Flow Prediction," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management* (2020), 1555–1564.
52. Y. Zhang, S. Wang, B. Chen, J. Cao, and Z. Huang, "TrafficGAN: Network-Scale Deep Traffic Prediction With Generative Adversarial Nets," *IEEE Transactions on Intelligent Transportation Systems* 22, no. 1 (2021): 219–230, <https://doi.org/10.1109/tits.2019.2955794>.
53. K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), 770–778, <https://doi.org/10.1109/CVPR.2016.90>.
54. G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely Connected Convolutional Networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2017), 2261–2269, <https://doi.org/10.1109/CVPR.2017.243>.
55. City of New York. NYC OpenData (2020).