

Improving short-term bike sharing demand forecast through an irregular convolutional neural network

Xinyu Li^a, Yang Xu^{a,b,*}, Xiaohu Zhang^c, Wenzhong Shi^{a,d}, Yang Yue^e, Qingquan Li^e

^aDepartment of Land Surveying and Geo-Informatics, The Hong Kong Polytechnic University

^bThe Hong Kong Polytechnic University Shenzhen Research Institute

^cDepartment of Urban Planning and Design, The University of Hong Kong

^dSmart Cities Research Institute, The Hong Kong Polytechnic University

^eDepartment of Urban Informatics, School of Architecture and Urban Planning, Shenzhen University

Abstract

As an important task for the management of bike sharing systems, accurate forecast of travel demand could facilitate dispatch and relocation of bicycles to improve user satisfaction. In recent years, many deep learning algorithms have been introduced to improve bicycle usage forecast. A typical practice is to integrate convolutional (CNN) and recurrent neural network (RNN) to capture spatial-temporal dependency in historical travel demand. For typical CNN, the convolution operation is conducted through a kernel that moves across a “matrix-format” city to extract features over spatially adjacent urban areas. This practice assumes that areas close to each other could provide useful information that improves prediction accuracy. However, bicycle usage in neighboring areas might not always be similar, given spatial variations in built environment characteristics and travel behavior that affect cycling activities. Yet, areas that are far apart can be relatively more similar in temporal usage patterns. To utilize the hidden linkage among these distant urban areas, the study proposes an irregular convolutional Long-Short Term Memory model (IrConv+LSTM) to improve short-term bike sharing demand forecast. The model modifies traditional CNN with irregular convolutional architecture to extract dependency among “semantic neighbors”. The proposed model is evaluated with a set of benchmark models in five study sites, which include one dockless bike sharing system in Singapore, and four station-based systems in Chicago, Washington, D.C., New York, and London. We find that IrConv+LSTM outperforms other benchmark models in the five cities. The model also achieves superior performance in areas with varying levels of bicycle usage and during peak periods. The findings suggest that “thinking beyond spatial neighbors” can further improve short-term travel demand prediction of urban bike sharing systems.

Keywords: bike sharing, deep learning, travel demand forecast, spatial-temporal analysis, irregular convolution

1. Introduction

Shared bicycles have received increasing attention in urban transportation during the past few decades [1, 2]. As a green transport option for short-distance travel in cities, bike-sharing services can reduce carbon emissions and enhance last-mile connectivity to public transit [3, 4, 5]. During COVID-19 pandemic, bike-sharing is found to be a more resilient mode that can mitigate the fear of overcrowding in public transit [6, 7, 8]. Given the importance of bike-sharing services in urban transportation, accurate demand forecasting is crucial for effective rebalancing in daily operations.

*Corresponding author: yang.ls.xu@polyu.edu.hk

Many studies attempted to develop frameworks that can accurately estimate the bicycle demand throughout the city by applying traditional and machine learning models [9, 10, 11, 12, 13, 14, 15, 16]

In recent years, deep learning approaches have been widely used for predicting short-term traffic demand [17, 18, 19, 20, 21, 22, 23, 24, 25]. A critical task is to model the spatial-temporal dependency in travel demand. Two mainstream architectures, Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN), are often integrated to capture the spatial and temporal information of traffic demand [26, 27, 28, 29, 30]. Typically, CNN exploits regular convolutional kernels scanning through the input features (e.g., images) to extract spatial characteristics of travel demand [31]. RNN leverages the extracted temporal dynamic behavior from the past elements of the sequence to predict the next element [31, 32]. To better capture spatial-temporal information, several hybrid deep learning frameworks incorporating both CNN and RNN architectures are developed and these models achieved good performance in various traffic prediction tasks [17, 21, 22, 23].

However, CNN has certain shortcomings when it is employed to capture spatial-temporal information of bike sharing demand. CNN achieves desirable performance in object detection for images, because adjacent pixels of the same object are often highly correlated. Unlike images, bicycle usage in neighboring urban areas can be quite different due to spatial variations of travel behavior and built environmental characteristics [33, 34]. On the other hand, for certain areas that are far apart, the bicycle usage patterns may exhibit similar temporal rhythms. Given the regular shape of convolutional kernels, deep learning models with typical CNN architecture are not able to capture the similarities of bike usage patterns among distant urban areas. If such similarities can be captured and incorporated into the prediction models, it may further enhance the accuracy and reliability of bike sharing demand forecast.

To bridge the research gap, this paper introduces an irregular convolutional Long Short-Term Memory model (IrConv+LSTM) to improve short-term demand forecast for urban bike-sharing systems. The model employs irregular convolutional architecture to capture the dependency of bicycle usage among distant urban areas. Given the areas being forecast, an irregular convolution operation is performed over their *semantic neighbors*, which refer to places that show similar temporal bicycle usage patterns. Two measures, namely Pearson Correlation Coefficient (IrConv+LSTM:P) and Dynamic Time Warping (IrConv+LSTM:D), are used as similarity metrics to identify semantic neighbors for the areas being forecast. The two variants of the proposed model (IrConv+LSTM:P and IrConv+LSTM:D) and several benchmark models are evaluated and compared over bike-sharing systems in five cities, including one dockless bike-sharing system in Singapore and four station-based systems in Washington D.C., Chicago, New York, and London, respectively.

The remainder of this paper is organized as follows. Section 2 provides a review of relevant literature. Study sites and related terminologies are introduced in Section 3. Section 4 and 5 illustrate the methodologies and experimental results. In Section 6, we conclude the study and discuss possible future research directions.

2. Literature Review

Many studies have been conducted on traffic demand prediction using either parametric or non-parametric models [35]. In parametric models, data series is modelled as a dynamic variation from a systemic basis. Most parametric models adopt filters to estimate parameters that capture the system characteristics for predicting future status. Autoregressive Integrated Moving Average (ARIMA) model and Kalman Filter are two typical parametric models [36]. ARIMA model uses Autoregressive (AR) or Moving Average(MA) methods to simulate the temporal autocorrelation of a smoothed sequence. The applications of ARIMA in traffic prediction include forecasting traffic accidents, traffic status from the perspective of speed, volume and travel time [13, 14, 15, 16]. Kalman Filtering is an algorithm for optimal system state estimation by using status equations of linear systems. Similar to ARIMA, Kalman Filtering not only predicts traffic status but also assists traffic management

and controls [37, 38, 39]. Such parametric models filter a wealth of information based on the strong assumptions of the data. Also, spatial information is not considered in such models. Therefore, the parametric models cannot fully model the spatial-temporal characteristics from the historical traffic data.

With the advancement of computing power, machine learning algorithms have been increasingly adopted for travel demand forecast. These algorithms include Support Vector Machine (SVM), Random Forest, Bayesian Network, Markov Model, Neural Network, and hybrid deep learning models [40, 41, 42, 43]. In particular, deep learning models have been attracting much attention in the last decade. Many models are applied to predict the traffic status, such as Convolutional Neural Network (CNN), Recurrent Neural Network (RNN) with its variants, encoder-decoder and attention mechanisms for sequential prediction, and Graph Convolutional Network (GCN) for graphic knowledge learning [26, 29, 44, 45, 46, 47]. For example, Ma *et al.* establishes a space-time matrix representing traffic sensors ordered by road directions and time, then adopts CNN to capture spatial and temporal information for predicting [26]. In [28, 29, 48], RNN and LSTM are employed to predict the traffic status, including traffic speed and congestion. To incorporate spatial and temporal characteristics, many scholars also propose several hybrid models to forecast the traffic demands or flows. Zhang *et al.* employs three residual CNNs to capture spatial-temporal information of historical pedestrian flows for predicting citywide crowd flows [20]. In [49], CNN and attention LSTM are adopted to forecast passenger flow in urban rail transit. Ren *et al.* adopts residual CNN and LSTM blocks for extracting high-level spatial-temporal information of pedestrian flow volumes to forecast citywide pedestrian volumes [23]. In general, deep learning models has achieved better performance in predicting traffic status than the parametric models.

There are several studies for predicting bike-sharing demand using deep learning approaches. Several studies adopt GCN to capture spatial characteristics and employ attention mechanisms or fully connected networks to extract sequential information in SBSS for predicting bike-sharing demand [45, 50, 51]. Ai *et al.* adopts convolutional LSTM to predict the bicycle demand in DBSS [17]. However, the characteristics of bike-sharing systems in cities are not considered in the models, affecting the models' performance in predicting bike-sharing demand. Besides, few studies tested the proposed deep learning model in different city contexts to predict both SBSS and DBSS. Such types of cross-city and cross-system studies are essential to examine the generalizability of proposed deep learning architectures.

3. Study Areas and Data Preprocessing

In this study, we use bike-sharing datasets in five different cities to assess the performance of our proposed model. These datasets include one dockless bike-sharing system (DBSS) in Singapore, and four station-based systems (SBSS) in Chicago, Washington D.C., New York, and London, respectively. The datasets of SBSS are collected from the bicycle trip records that document the departure and arrival time and stations. The dataset of DBSS in Singapore is collected from raw GPS coordinates of starting and ending locations of a trip. The GPS trajectories are preprocessed to remove GPS drifts and outliers using the approach from a prior study in Singapore [34].

As shown in Fig. 1, we use regular grids to summarize the usage pattern of shared bicycles in the five cities for both DBSS and SBSS. A city area is divided into a regular $w * h$ grid map under a specific spatial resolution. During a time interval, the grid values represent the number of bicycle pick-ups distributed in the city. As shown in Table 1, we set $1km$ as the spatial resolution and 1 hour as the temporal resolution for this study. At the k^{th} time interval, the definition of a grid map ($X_k(w, h)$) is shown in Eq. 1. $x_k(i, j)$ denotes the number of pick-ups in the cell located in the i^{th}

Table 1: Description of Research Areas

Para.	City	Singapore	Chicago	Washington D.C.	New York	London
Spatial Resolution				1 Kilometre		
Temporal Resolution				1 Hour		
Training Period		16/06-02/08, 2017	01/06-30/09, 2019	01/06-30/09, 2019	01/06-30/09, 2019	01/06-30/09, 2019
Validation Period		03/08-31/08, 2017	01/10-25/10, 2019	01/10-31/10, 2019	01/10-25/10, 2019	01/10-31/10, 2019

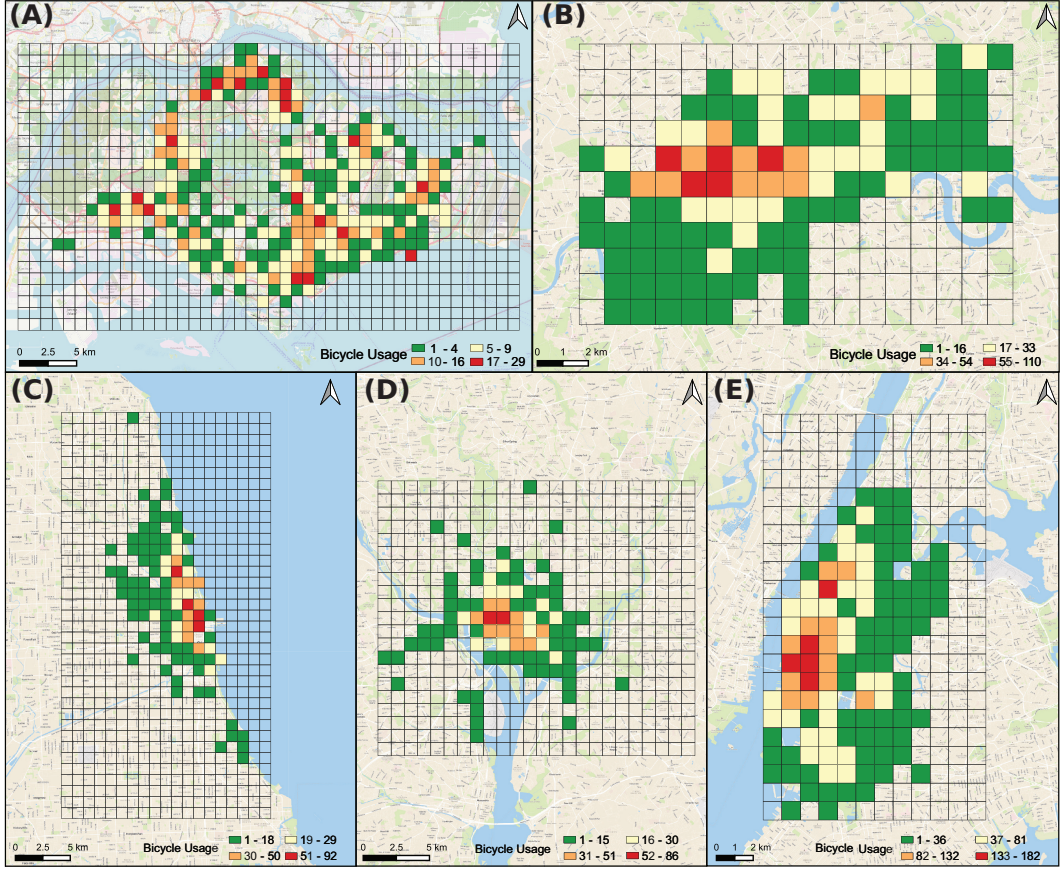


Figure 1: Spatial distribution of bicycle usage in five cities. (A) Singapore during 5-6 PM on August 25th, 2017; (B) London during 5-6 PM on August 24th, 2019; (C) Chicago during 1-2 PM on August 25th, 2019; (D) Washington, D.C. during 5-6 PM on August 25th, 2019; (E) New York during 5-6 PM on August 25th, 2019.

row and the j^{th} column at the k^{th} time interval, defined in Eq. 2.

$$X_k(w, h) = \begin{bmatrix} x_k(1, 1) & x_k(1, 2) & \cdots & x_k(1, h-1) & x_k(1, h) \\ x_k(2, 1) & x_k(2, 2) & \cdots & x_k(2, h-1) & x_k(2, h) \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_k(w, 1) & x_k(w, 2) & \cdots & x_k(w, h-1) & x_k(w, h) \end{bmatrix} \quad (1)$$

$$x_k(i, j) = |\{T \in \mathbb{T}_k \mid T(O) \in (i, j) \wedge T(D) \notin (i, j)\}| \quad (2)$$

Here \mathbb{T}_k denotes the trajectories occurred in k^{th} time slot; $T(O)$ and $T(D)$ represent the departure and arrival locations of a trajectory T , respectively; $T(O) \in (i, j) \wedge T(D) \notin (i, j)$ denotes a trajectory starts from the cell (i, j) but ends in another cell except the cell (i, j) ; $|\cdot|$ denotes the cardinality of a set. In this study, trips that start and end in the same cell are excluded because they do not affect the overall balance of bicycle supply and demand within a cell.

4. Methodology

4.1. Overall architecture of the proposed model

Fig. 2 illustrates the overall architecture of the proposed model. The model consists of three separate modules with the same structure. Each module takes a specific set of historical observations (bike-sharing demand) as input. Instead of using all the historical observations to train the model, we identify key periods with different levels of recency to the target period (for which the prediction is made) and feed them into the three modules. The approach effectively reduces the time complexity of the training model by also mitigating the negative effect of redundant information in historical data. This practice has proved to achieve better performance than models trained using full historical observations [20, 23, 21, 52]. The definitions of the three key periods (*trend*, *period* and *closeness*) will be elaborated in Section 4.3. As shown in Fig. 2, each module adopts three layers of irregular

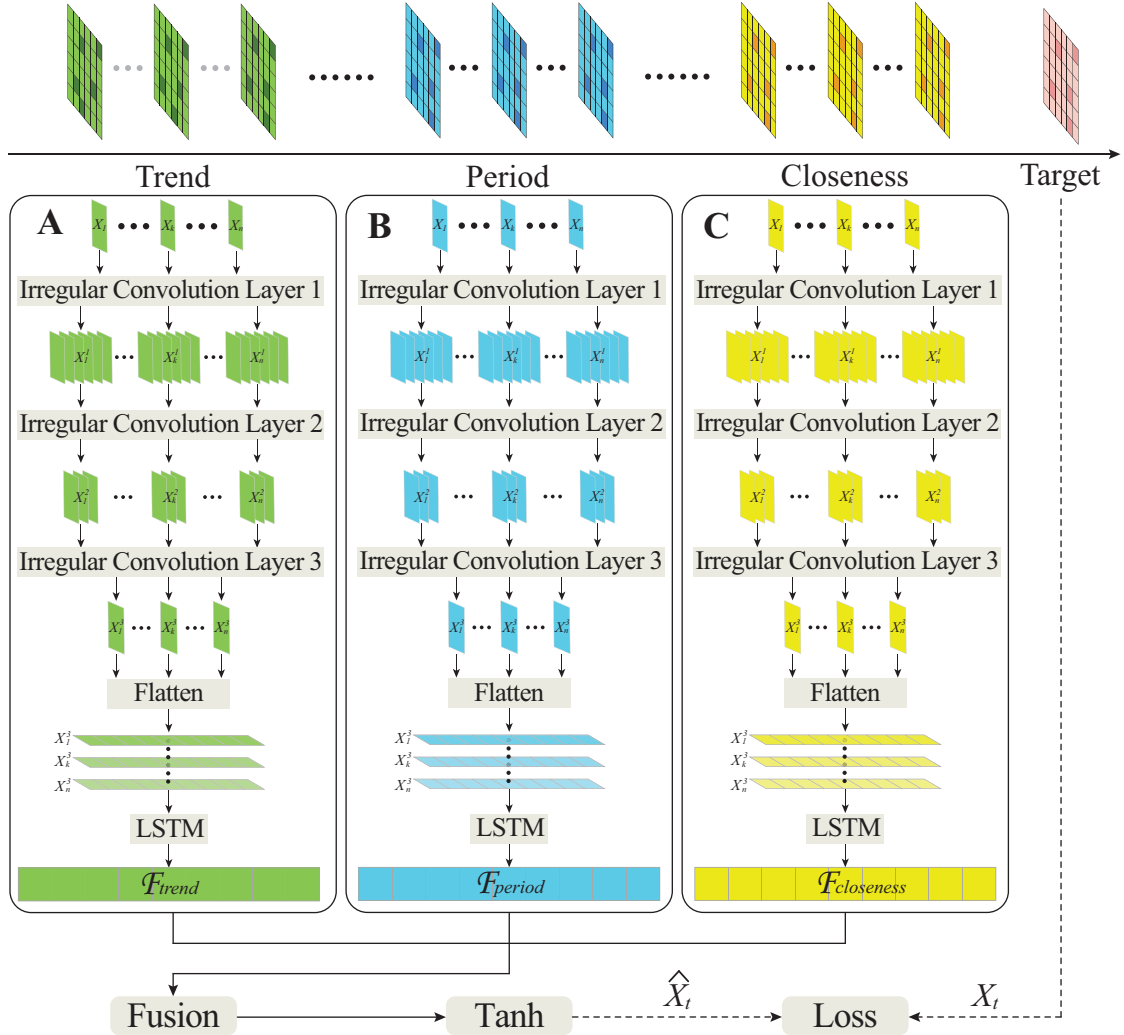


Figure 2: Overall architecture of the proposed model.

convolutional architecture to capture the characteristics of bicycle demand among urban areas. The vector sequence formed by flattening the output of the irregular convolution is used as the input to the LSTM model¹ to extract the temporal information in the sequence. The outputs of three hybrid modules are fed into a feature fusion layer. The output of the feature fusion layer is activated by a non-linear function generating the predicted value. The predicted value with its corresponding actual usage value participates loss estimation and backpropagation to update parameters in the model. In the following section, we formally introduce the architecture of irregular convolutional network.

4.2. Irregular convolutional neural network

The major difference between irregular convolution and traditional convolution lies in the cells involved in the convolutional operation. Generally, the number of cells involved in the convolution is known as the convolutional kernel size. Among cells corresponding to each kernel, the cell being forecast is called the central cell, and the other cells involved in convolution are called neighbors. Taking the convolutional kernel size of nine as an example, the neighbors involved in traditional convolution are spatially adjacent to the central cell, as shown in Fig. 3A. In contrast, the neighbors can be located anywhere in the study area for irregular convolution (Fig. 3B). In this study, for each central cell, we identify the top eight cells which show similar temporal bicycle usage patterns observed from the historical observation data. We call these cells as *semantic neighbors*. Compared to the traditional convolution, irregular convolution is more flexible to exploit cells with similar temporal usage patterns to the central cell.

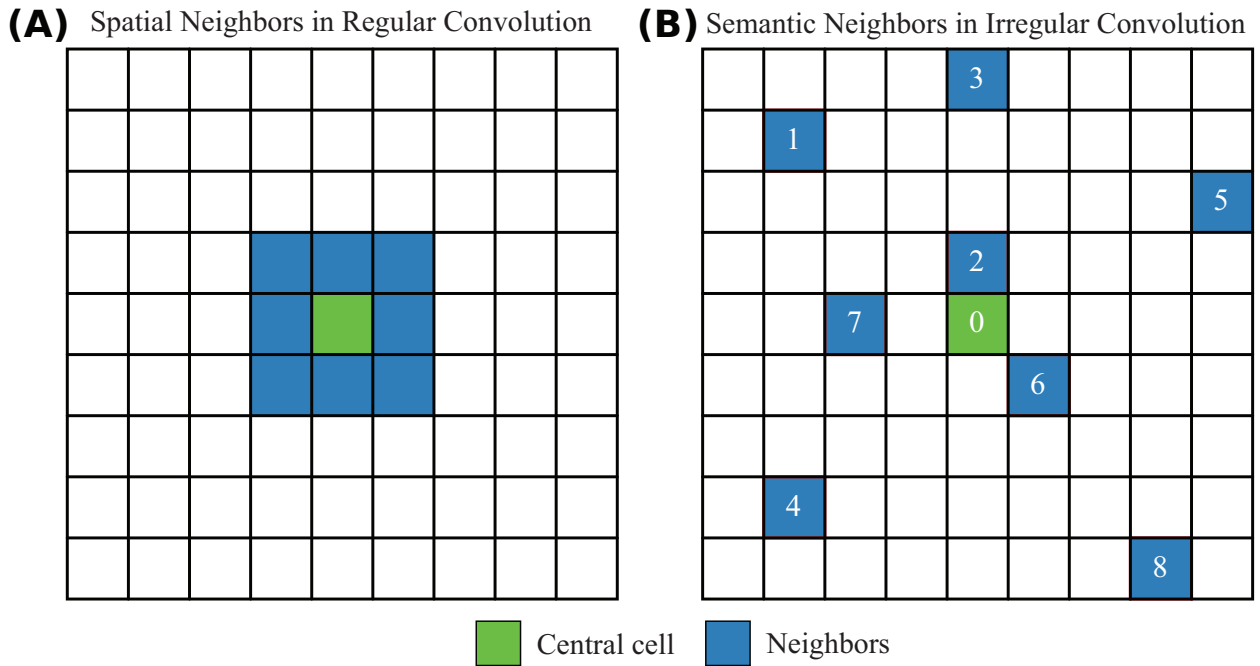


Figure 3: Illustration of traditional and irregular convolution with a kernel size of 9: (A) spatial neighbors are adjacent to the central cell in the regular convolution; (B) in the irregular convolution, semantic neighbors are identified based on the top 8 cells with the highest similarity of temporal bicycle usage patterns to the central cell.

This study uses two metrics, namely, the Pearson Correlation Coefficient [53, 54] and Dynamic Time Warping (DTW) [55, 56], to quantify the similarity of temporal bicycle usage patterns between the cells. Because these two metrics measure temporal similarity from different perspectives, we aim to assess which metrics tends to result into a better prediction accuracy.

¹The description of the LSTM model is provided in Appendix A.

The Pearson correlation coefficient measures the strength of a linear association between two sequences. The coefficient is the ratio between covariance of two variables and the product of their standard deviations. The two sequences are more positively correlated when the ratio is closer to 1. Eq. 3 shows the calculation of the Pearson correlation coefficient. Fig. 4A shows an example of the Pearson correlation coefficient of two bicycle usages sequences within 24 hours.

$$r_{xy} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}} \quad (3)$$

where \bar{X} and \bar{Y} denote the means of sequence X and Y , respectively. n denotes the length of sequence.

For DTW, given two bicycle usage sequence $X = x_1, x_2, \dots, x_i, \dots, x_{|X|}$ and $Y = y_1, y_2, \dots, y_j, \dots, y_{|Y|}$, the optimization objective is to find a shortest warp distance between two sequences $dist(W)$ (typically Euclidean distance), as shown from Eq. 4 to Eq. 6:

$$dist(W) = \underset{W}{\text{minimum}} \left(\sum_{k=1}^{K} \text{distance}(w_{ki}, w_{kj}) \right) \quad (4)$$

subject to :

$$W = w_1, w_2, \dots, w_K, \max(|X|, |Y|) \leq K < |X| + |Y| \quad (5)$$

$$w_k = (i, j), w_{k+1} = (i', j'), i \leq i' \leq i + 1, j \leq j' \leq j + 1 \quad (6)$$

where K is the length of warp path, w_k denotes the k^{th} element of the warp path, i and j represent the index of a certain record in two sequences X and Y , respectively. $\text{distance}(w_{ki}, w_{kj})$ is the distance between two records (the former from X and the latter from Y) in the k^{th} element of the warp path. Thus, a cost matrix D is constructed based on the distance from any record in one of sequences (X/Y) to any record in another sequence (Y/X). DTW adopts a greedy search from $D(|X|, |Y|)$ to $D(1, 1)$ to find the minimum distance between two sequences. In this study, we use historical observations of each cell during the whole training period to estimate the similarity. The length of the evaluated sequences are the same for two mentioned metrics ($n = |X| = |Y|$). Fig. 4B shows an example of finding the shortest-distance warp path based on the cost matrix D for two bicycle usage sequences.

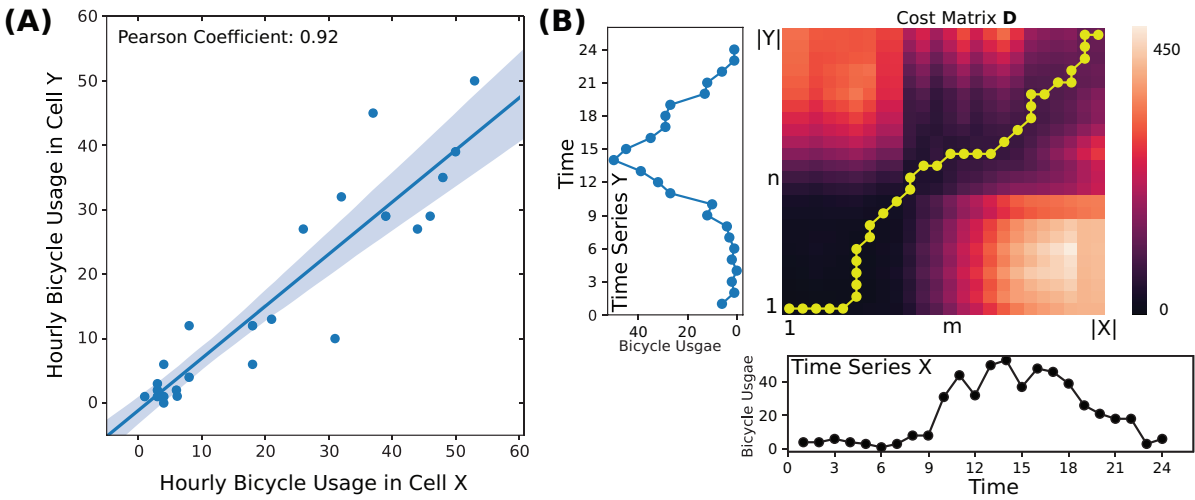


Figure 4: Two similarity metrics adopted in this study: (A) Pearson correlation coefficient; (B) Dynamic time warping (DTW).

The irregular convolutional computation is analogous to the traditional convolution:

$$y(i, j) = b(i, j) + \sum_{c=1}^{C_{in}} \sum_{s=1}^S x_c^s(i, j) w_c^s \quad (7)$$

where C_{in} denotes the number of channels in input $x(i, j)$; S denotes the size of convolutional kernel; $x_c^s(i, j)$ represents the semantic neighbor s associated with central cell $x(i, j)$ in the channel c ; w_c^s denotes the weight in the convolutional kernel corresponding to the semantic neighbor $x_c^s(i, j)$; $b(i, j)$ denotes the learnable bias.

4.3. Definitions of Trend, Period and Closeness

As mentioned in 4.1, given a target period (for which the prediction is made), we identify three key periods from historical observations according to the levels of recency to the target period. We name them as *Closeness*, *Period*, and *Trend* (Fig. 2). The purpose of selecting observations from these periods as inputs is to reduce the negative impact of redundant information in the whole training data on model performance while lowering the training time complexity. For example, the usage patterns of shared bicycles during peak hours may be similar during weekdays, and bicycle usage might be associated with that at the same time in the previous week. This strategy of identifying key periods as training input has proved to achieve better performance than models trained using full historical observations [20, 23, 21, 52]. In this study, the definitions of such three key periods are given in Eq. 8 to Eq. 10.

$$\mathcal{X}_t^{closeness} = \{X_{t-l_c}, X_{t-(l_c-1)}, \dots, X_{t-1}\} \quad (8)$$

$$\mathcal{X}_t^{period} = \{X_{t-24*l_p}, X_{t-24*(l_p-1)}, \dots, X_{t-24}\} \quad (9)$$

$$\mathcal{X}_t^{trend} = \{X_{t-7*24*l_q}, X_{t-7*24*(l_q-1)}, \dots, X_{t-7*24}\} \quad (10)$$

We select l_c as 24 to capture historical shared bicycle usage data in the past 24 hours. The value of l_p is set as 7 to select the usage data at the same time for each day in the past week. l_q is chosen as 2 to provide historical information at the same time in the past week and the week before the past. The time complexity of training a model using the above mentioned key periods will be lower than that of training with a lengthy period of data.

The outputs of three separate modules are fused for the final forecast. We adopt the weighted element-wise addition method to merge the three spatial-temporal features, including \mathcal{F}_{trend} , \mathcal{F}_{period} and $\mathcal{F}_{closeness}$. Then, the feature map is activated by a *tanh* function to generate the prediction values \hat{X}_t that participates in the loss and backpropagation with the actual bicycle usage X_t . The computation of the prediction values \hat{X}_t is shown Eq. 11.

$$\hat{X}_t = \tanh(\mathcal{F}) = \tanh(W_t \circ \mathcal{F}_{trend} + W_p \circ \mathcal{F}_{period} + W_c \circ \mathcal{F}_{closeness}) \quad (11)$$

where W_t , W_p and W_c denote the learnable parametric vectors with the same shape size of the corresponding feature, \mathcal{F} is the feature map after fusion, and *tanh* denotes the activation function.

4.4. Hyperparameter settings and benchmark models

Since we adopt two metrics to quantify the similarity of temporal usage patterns, the performance of two variants of IrConv+LSTM is evaluated in this study. We name them as IrConv+LSTM:P and IrConv+LSTM:D, respectively. The hyperparameters for both variants are the same and refer to several existing studies [46, 57, 47]. Three irregular convolutional layers in each separate module operation (shown in Fig. 2) are adopted with 32 filters, 16 filters, and 1 filter, respectively. The last filter in the irregular architecture is to aggregate the high-dimensional information into one channel. The convolutional kernel size for two variants of our model is set as nine.

Four baseline models are adopted for performance comparison with two IrConv+LSTM variants, including one parametric model (ARIMA) and three deep learning models (LSTM, STRN, and CNN+LSTM):

- **ARIMA:** Auto-Regressive Integrated Moving Average model is a parametric model widely used for time series forecasting. ARIMA is a combination of the differenced autoregressive model (AR) with the moving average model (MA) [13]. ARIMA can handle non-stationary sequences by replacing the data values with the difference between their values and the previous values to obtain stationary sequences. Thus, ARIMA is widely adopted to predict traffic status that is dynamically changed over time, such as traffic accidents, traffic speed, and traffic volume [14, 15, 16].
- **LSTM:** Long Short-Term Memory is a widely used deep learning architecture in the field of traffic prediction [29, 28, 48]. As a variant of Recurrent Neural Network (RNN), one of the advantages of LSTM is that it can capture the information for both long and short periods of time by introducing gate theory. LSTM is effective for processing long sequences of input data because it tackles the gradient vanishing and explosion problems that can be encountered when training traditional RNN models.
- **STRN:** Spatial-Temporal Residual Network is a hybrid deep learning prediction model initially proposed to forecast short-term pedestrian volume [20]. It consists of three Residual Convolutional Neural Networks (ResNet). The spatial features captured by ResNets from multiple fragments of historical data are fused together for the final prediction. The strategy of adopting fragments of historical data to train deep learning prediction models has been referenced in many studies [58, 23, 21, 30].
- **CNN+LSTM:** CNN+LSTM is a hybrid deep learning model for spatial-temporal prediction. This model couples a traditional convolution and an LSTM model to extract historical spatial-temporal information for the prediction. It has been adopted for traffic prediction by several studies [59, 60]. The difference between this model and our proposed model is that CNN+LSTM adopts traditional convolution involving spatial neighbors.

To make the prediction results of IrConv+LSTM and CNN+LSTM comparable, we set similar hyperparameters and structures for both of them. Specifically, three hybrid modules in CNN+LSTM extract spatial-temporal information from their respective key historical periods for modeling. The definitions of such three key periods in CNN+LSTM is the same as them in our model. Each separate module also contains three traditional convolution layers with 32 filters in the 1st layer, 16 filters in the 2nd layer, 1 filter in the 3rd layer. The convolutional kernel size adopted in CNN+LSTM is also set as nine. The hyperparameter settings of STRN are referred to the settings in article [20]. The parameters of all benchmarks have been calibrated to achieve optimal prediction results.

We use the first 80% of the hourly usage data as the training data and the last 20% data to validate the performance of models in each city, as shown in Table 1. The loss function of all deep learning models in this study is the MSE Loss function. Eq. 12 defines the MSE Loss function that is used to measure the errors between the prediction value and the ground truth [61]. The optimization algorithm for updating parameters in deep learning models is important for backpropagation. This study employs RMSProp (Root Mean Square Prop) as the optimization algorithm across all deep learning models [62]. Moreover, three indicators are employed to evaluate the performance of models, including Mean Absolute Percentage Error (MAPE), Root Mean Square Error (RMSE), and Mean Absolute Error (MAE) [20, 23, 22]. The definitions of them are shown from Eq. 13 to Eq. 15.

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (12)$$

$$MAPE = \frac{100\%}{n} \sum_{i=1}^n \left| \frac{\hat{y}_i - y_i}{y_i} \right| \quad (13)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (14)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (15)$$

Here \hat{y}_i denotes the prediction result and y_i denotes the actual value of bicycle usage. All deep learning frameworks are constructed on the Pytorch platform [63]. Also, the models are built on a server with NVIDIA Tesla V100 and a workstation with NVIDIA RTX 2070 Super Graphics Card.

5. Analysis Results

5.1. Overall accuracy of the proposed model and benchmark models

In this section, we compare the overall accuracy between our proposed model and the benchmark models. Table 2 shows the overall accuracy of two variants of our proposed model and other benchmark models. Given a specific indicator, the model with the best performance is marked with * in Table 2.

Table 2: Overall accuracy of all models across five cities

City	Index	ARIMA	LSTM	STRN	CNN+ LSTM	IrConv- LSTM:P	IrConv- LSTM:D
Singapore	MAPE	0.8488	0.6715	0.7026	0.6696	0.5617*	0.5638
	MAE	2.6267	2.2411	2.2488	2.0971	1.9727	1.9655*
	RMSE	4.0656	3.4403	3.3602	3.2243	3.0911	3.0764*
Chicago	MAPE	1.3000	0.8615	0.7131	0.7584	0.7240	0.6028*
	MAE	4.6156	3.6955	2.5619	2.8373	2.4036	2.2047*
	RMSE	12.6296	10.1246	4.7583	6.5034	4.7802	4.3356*
Washington, D.C.	MAPE	0.9672	0.7133	0.6858	0.6492	0.5599*	0.5675
	MAE	3.7649	2.3535	2.3327	2.4026	1.9758	1.9697*
	RMSE	7.2917	4.0227	4.0657	4.0718	3.3669	3.3590*
New York	MAPE	2.9407	0.7479	0.7770	0.7153	0.6243	0.6201*
	MAE	16.9310	7.5721	7.2278	6.1815	5.9865	5.7769*
	RMSE	31.7418	15.3383	13.7237	11.3954	11.3385	10.8141*
London	MAPE	1.7924	0.7483	0.6753	0.6862	0.5852	0.5523*
	MAE	8.7771	4.8125	4.2348	4.0423	3.7829	3.5785*
	RMSE	17.8831	9.4463	7.1533	7.2279	6.6326	6.2960*

Generally, our proposed model achieves better performance than other baseline models across five cities based on the prediction accuracy. As a typical parametric model, ARIMA is hard to process the non-stationary sequence and cannot leverage any spatial dependency of bike-sharing usage. Thus, it achieves the lowest accuracy of forecasting bicycle demand across five cities. Unlike ARIMA, LSTM model is effective to extract temporal information of bicycle usage in each cell

from non-stationary sequence by using the mechanism of deep learning technology. Hence, the results of LSTM are better than that of ARIMA. However, LSTM still cannot utilize the spatial dependency of bicycle usage among cells for prediction. Based on the results, the hybrid deep learning model that couples convolution architectures and LSTM model to extract spatial-temporal features of bicycle usage generally outperforms LSTM in all cities. Moreover, STRN also performs better than LSTM. Although STRN and CNN+LSTM achieve good prediction results, the performance of our proposed model is still better than them in all indicators. Notably, our proposed model achieves an improvement of MAPE from 8% (in Washington, D.C.) to 12% (in London), compared to the model with the best performance in four benchmarks.

Compared to CNN+LSTM, our proposed model only replaces the spatial neighbors involved in regular convolution with the semantic neighbors adopted in irregular convolution. However, the prediction accuracy of our model is much higher than CNN+LSTM across five cities. Specifically, among five cities, our approach achieves an improvement of MAPE by 8%-15% compared to CNN+LSTM. Such results imply that the semantic neighbors are more effective than spatial neighbors for predicting bike-sharing usage. The semantic neighbors are identified according to the similarity of temporal bicycle usage to their corresponding central cells. Although the semantic neighbors are not always spatially adjacent to their central cells, they can provide essential information for spatial-temporal modeling than the spatial neighbors. Therefore, the prediction results suggest that involving semantic neighbors in the irregular convolution is a good strategy for predicting shared-bicycle usage across all study areas.

We also find that the performance of the variant with DTW metric is better than the variant with Pearson measure in most cities. The variant with Pearson correlation coefficient only performs slightly better in two cities (Singapore and Washington D.C.) from the perspective of MAPE. Based on the characteristics of Pearson correlation coefficient, the usage variations of two measured sequences are simply reflected by their respective standard deviations, and the temporal offsets between such two sequences cannot be quantified by the coefficient. In contrast, DTW metric is effective to measure the similarity of usage variations between two sequences by calculating the distance of each record in one sequence to all other observations in the other. Also, the shortest warping distance between the sequences searched by DTW metric reflects the similarity without the effects of temporal offsets between the sequences. In other words, for two sequences that are not synchronized, DTW metric is able to quantify their similarity more precisely than Pearson correlation coefficient. Especially for quantifying the similarity of temporal bicycle usage patterns, DTW metric, which considers the temporal offsets, better reflects the similarity of travel behavior among areas. In sum, our proposed model achieves the best overall performance compared to other baseline models. Also, DTW metric is an effective way to select semantic neighbors involved in the irregular convolution.

5.2. Performance of models over cells with varying levels of bicycle usage

When operating bike-sharing systems, one of the important tasks is to satisfy users' travel needs. However, bike-sharing demand is not evenly distributed in urban areas. The number of areas with high user demand is relatively small but most users intend to ride bicycles in such areas. In contrast, areas with low demand require fewer bicycles, but they are widely distributed in the city. Therefore, there is a trade-off when allocating bicycles across urban areas. If the travel demand in areas with different usage can be precisely predicted, it is helpful for the deployment of bicycles.

We separate the cells into five quantiles in each city according to their respective hourly usage to assess the models' performance over these quantiles. Fig. 5 shows the MAE distributions of each quantile of bicycle usage in the five cities. We find that the performance of our model is better than baseline models across five quantiles of bicycle usage in most cities. Particularly, for these high-demand cells, the prediction accuracy of our proposed model is higher than other benchmarks. For example, the average MAE of our approach is less than 17 bicycles in Singapore, lower than that of the benchmarks. Also, the maximum MAE in Chicago is close to 10 bicycles, much smaller

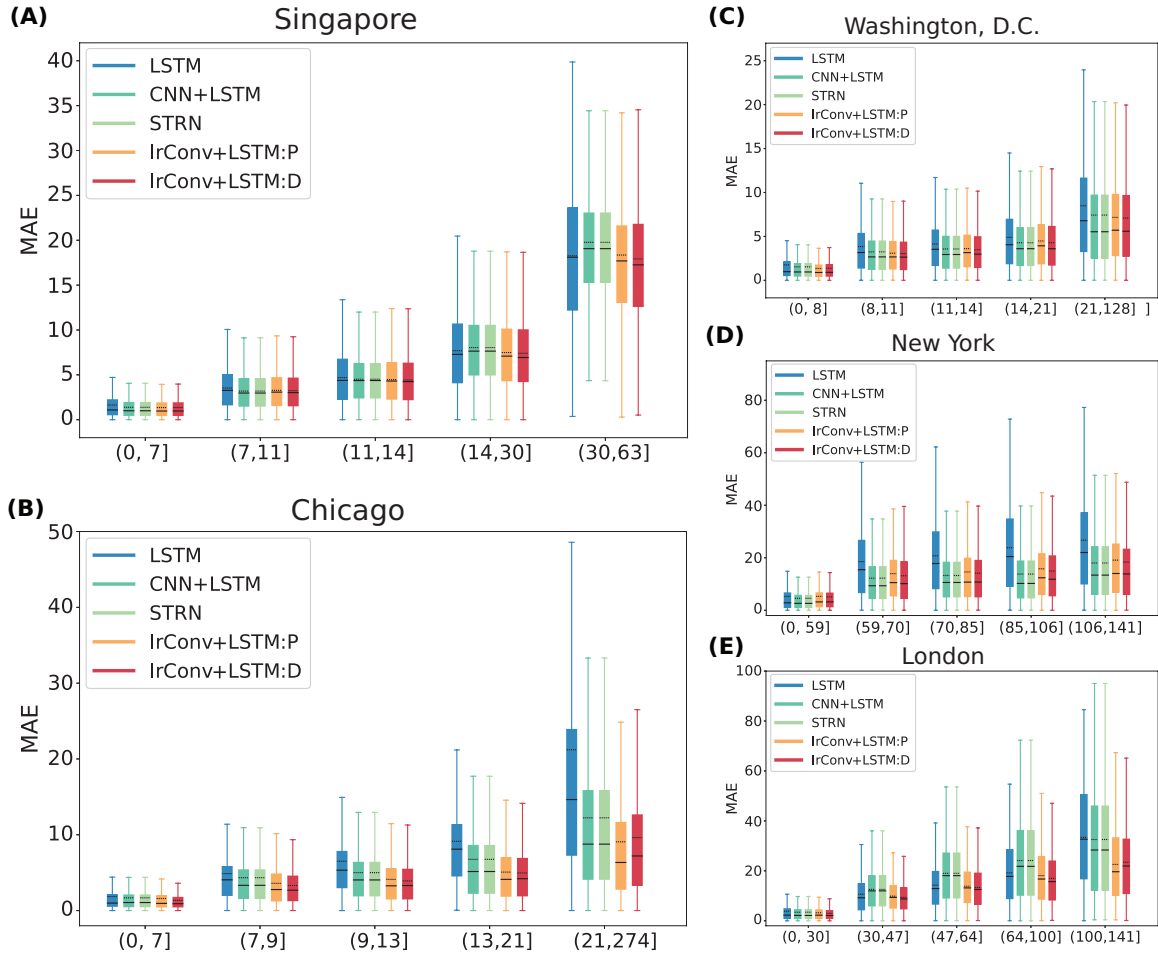


Figure 5: Performance of models in cells with various levels of bicycle usage. The cells in each city are divided into five quantiles according to their respective hourly usage, in order from low to high usage shown in horizontal axes.

than other benchmarks. However, in Washington, D.C., and New York, the average MAE of our approach is similar to that of other baseline models. This fact indicates that our approach has a similar prediction ability to the other models in such two cities with large bicycle usage. Additionally, our proposed model achieves better performance on the other quantiles of bicycle usage. Even in low-demand areas in most cities, our proposed model’s average MAE is lower than the other baseline models. To sum up, IrConv+LSTM has achieved better performance over areas with varying levels of bicycle usage.

5.3. Performance of models during peak hours

Another important task in operating bike-sharing systems is meeting the users’ needs during the morning and evening commuting hours. For example, in some cities, there are many users who prefer to ride bicycles to address first- and/or last-mile problems during the morning and evening peak hours [34, 1]. The operators need to allocate enough available bicycles for users during such peak hours. Accurate bicycle demand forecasting can help them develop proper scheduling plans to satisfy users with as little cost as possible.

We select the morning peak (7:00-10:00 AM) and evening peak (5:00-8:00 PM) observations from the validation datasets to assess the models’ performance in each city. Fig. 6 represents the prediction accuracy of each model during peak hours in five cities based on the two indicators (MAPE and MAE). The prediction error of our proposed model is the lowest across five cities during both peak

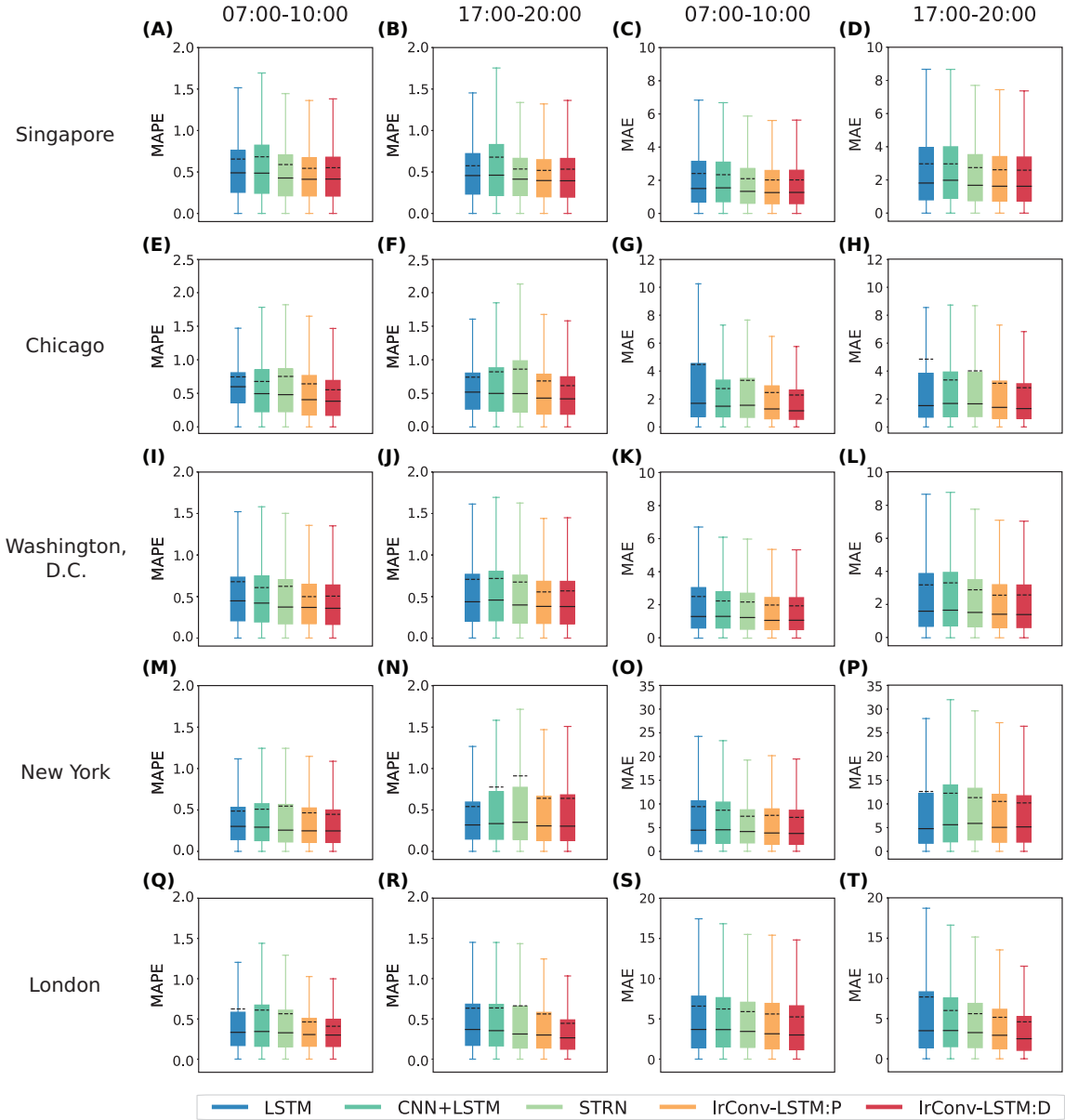


Figure 6: Performance of models during peak hours: (A-D) Singapore; (E-H) Chicago; (I-L) Washington D.C.; (M-P) New York; and (Q-T) London.

hours. Specifically, in London, the MAE of IrConv+LSTM:D is smaller than five during the evening period, better than the benchmarks. In Chicago, the MAPE of our approach is close to 50% during both peak periods, which is smaller than other baseline models. Although IrConv+LSTM performs similarly to CNN+LSTM in Singapore during both peak hours, the performance of IrConv+LSTM is stable and better across all five cities. Comparing the two variants of our proposed model, the performance of the variant adopted DTW metric is more robust in most cities than the variant with the Pearson measure. To sum up, our model outperforms the benchmarks during peak hours in five cities.

5.4. Comparative analysis between semantic and spatial neighbors

As mentioned in Section 5.1, our approach involving semantic neighbors in irregular convolution has a notable improvement in model performance compared to the benchmark (CNN+LSTM) in-

corporating spatial neighbors in traditional convolution. Here we perform an additional analysis to gain insights into the spatial relationship between cells’ semantic and spatial neighbors.

First, we analyze the difference in temporal similarity between the central cell and spatial neighbors vs. the central cell and semantic neighbors. As shown in Fig. 7, the similarity of temporal usage patterns between central cells and their semantic neighbors is generally higher than that between central cells and their spatial neighbors. Combined with the prediction accuracy of our model and CNN+LSTM, the similarity of semantic neighbors is more effective in enhancing the accuracy of deep learning models. However, the traditional convolution fails to leverage such similarity due to its structural limitations. Moreover, the similarity differences between the two types of neighbors imply that directly applying traditional convolution involving spatial neighbors for shared bicycle usage prediction might not always achieve desirable performance. In the five study sites evaluated in this research, CNN+LSTM based on regular convolution architecture performs worse than our proposed model.

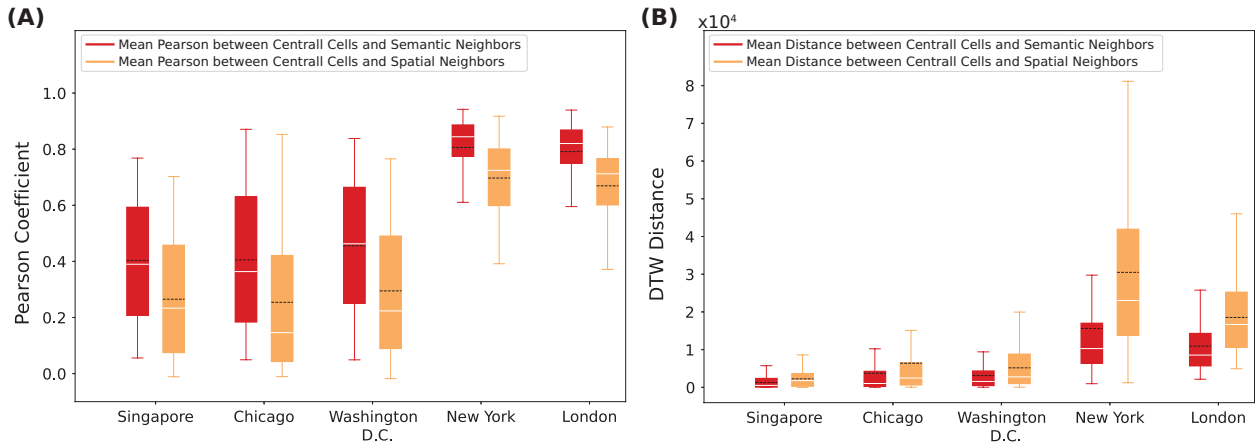


Figure 7: **(A)** Distribution of average Pearson correlation coefficients between the central cell and spatial neighbors vs. semantic neighbors; **(B)** Distribution of average DTW distance between the central cell and spatial neighbors vs. semantic neighbors.

From a spatial point of view, we evaluate the extent to which semantic neighbors overlap with spatial neighbors. Fig. 8 shows the distribution of overlap cells between spatial and semantic neighbors in five cities. We find that more than 70% of central cells have their spatial and semantic neighbors sharing less than two cells, which can be observed under both Pearson and DTW metrics. This fact indicates that the semantic neighbors are mainly distributed at non-spatially adjacent areas to the central cells, further illustrating the difference between semantic and spatial neighbors. In other words, some areas in a city, though far apart, could show relatively similar temporal bicycle usage patterns than spatially adjacent areas. This is potentially attributed to built environment characteristics, urban functions, and other factors that shape users’ travel behavior across urban areas.

6. Conclusion and Discussion

This paper proposes a hybrid deep learning model coupling irregular convolution architectures and LSTM modules to predict bicycle usage demand across one dockless bike-sharing system in Singapore and four station-based systems in Chicago, Washington D.C., New York, and London. The irregular convolution is applied over semantic neighbors that refer to places with temporal usage patterns similar to those of the areas being forecast. To measure the cells’ similarity and build the semantic neighbors, one variant of our model uses the Pearson correlation coefficient while the other adopts Dynamic Time Warping (DTW). Our proposed model and four benchmark models are

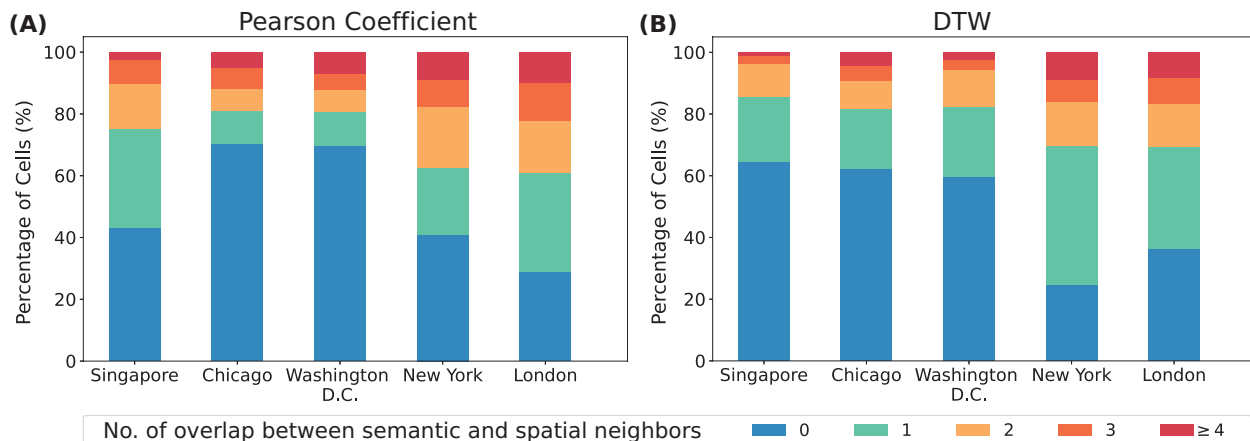


Figure 8: Number of overlapping cells between a central’s spatial and semantic neighbors: (A) results based on Pearson correlation coefficient; (B) results based on Dynamic Time Warping (DTW).

evaluated in five cities. Based on the prediction results, the proposed model generally outperforms the benchmarks across five cities. The prediction accuracy of the proposed model is also higher than that of the benchmarks in areas with different bicycle usage levels and during peak hours. Comparing the two variants of our model, the performance of the variant with the DTW metric is better than that of the variant with the Pearson coefficient metric in most of the cities.

We find that the semantic neighbors adopted in the irregular convolution are quite different from the spatial neighbors involved in the traditional convolution. Specifically, the similarity of temporal usage patterns between the central cells and their semantic neighbors is generally higher than that between the central cells and their spatial neighbors under both the Pearson and DTW metrics. Unlike the spatial neighbors, the semantic neighbors are mainly distributed in areas that are not spatially adjacent to their central cells. The model comparison suggests that relating areas that share similar temporal usage patterns through irregular convolution is helpful to improving the prediction accuracy. The findings also indicate that the neighbors involved in the convolution and the metrics for quantifying similarity (e.g., Pearson correlation vs. DTW) among urban areas tend to influence the performance of the prediction model.

The study suggests that “thinking beyond spatial neighbors” can inspire new solutions that improve short-term travel demand prediction. The implications go beyond applications of bike-sharing systems. In general, a reliable travel demand or traffic prediction model requires a good understanding on the spatial-temporal dependency of historical travel demand. Although the subjects being studied could vary (e.g., vehicles, cyclist, pedestrians and goods), it is reasonable to assume that there exists spatial autocorrelation in travel patterns. This partially explains why deep learning architecture with typical CNN (i.e., by capturing spatial dependency) is helpful to improving travel demand prediction in different transportation applications. However, given that travel patterns are affected by spatial variations of built environment characteristics, urban functions and socioeconomic characteristics, areas close to each other might show different travel dynamics that are sometimes not correlated over time. Therefore, the typical CNN architecture can sometimes introduce “noise” into the prediction process. The irregular convolution with semantic neighbors is one example of many possible strategies to overcome this limitation.

In this study, the semantic neighbors are identified solely based on the temporal similarity in historical bicycle usage patterns. These areas with similar temporal travel patterns could signify potential similarities in built environment characteristics and other factors that shape cycling behaviors. However, these environmental and socioeconomic characteristics are not directly utilized in our proposed model. A possible direction for future research is to consider both travel patterns and

environmental factors for defining similarities and identifying semantic neighbors. Combing these static (e.g., built environment characteristics) and dynamic features (i.e., bicycle usage patterns) may further improve the robustness of the prediction model.

Appendix A. The Structure of Long Short-Term Memory Network

We adopt the Long Short-Term Memory (LSTM) model to extract temporal information from the features captured by the irregular convolutional network [27, 28, 30, 32]. The structure of LSTM is shown in Fig. A.9. As a variant of Recurrent Neural Network, LSTM adopts the gate theory to control the information captured from both long and short periods of time for avoiding gradient descent or explosion. The computations of LSTM are shown from Eq. A.1 to Eq. A.6.

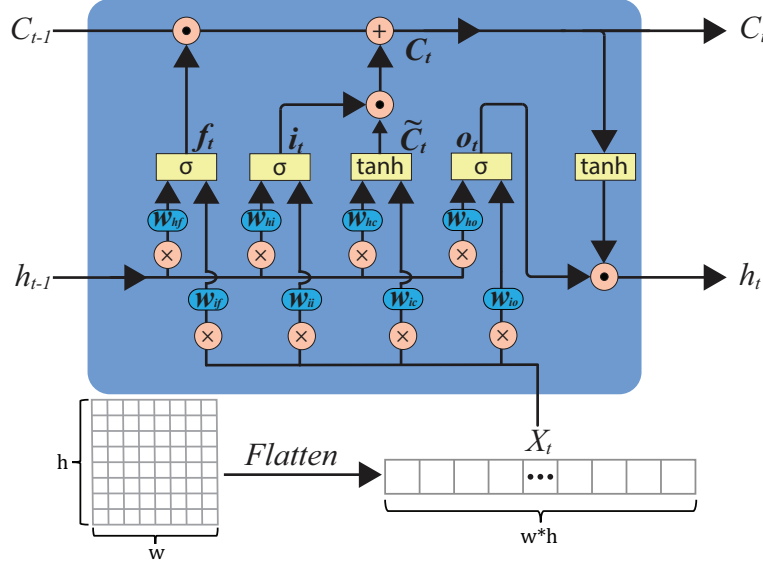


Figure A.9: The structure of Long Short-Term Memory Network

$$f_t = \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \quad (\text{A.1})$$

$$i_t = \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) \quad (\text{A.2})$$

$$\tilde{C}_t = \tanh(W_{ic}x_t + b_{ic} + W_{hc}h_{t-1} + b_{hc}) \quad (\text{A.3})$$

$$o_t = \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) \quad (\text{A.4})$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad (\text{A.5})$$

$$h_t = o_t \odot \tanh(C_t) \quad (\text{A.6})$$

Here f_t , i_t and o_t are the forget, input and output gates, respectively. W_* and b_* are respectively learnable weights and bias. \odot is the Hadamard product, and σ denotes the sigmoid function. Particularly, the input x_t is the feature captured by irregular convolutional network at time t ; C_t denotes the cell state at time t aggregating the information from forget, input gates, cell state and hidden state of the previous layer; h_t represents the hidden state that incorporates the information from C_t and the output gate. According to the equations of LSTM, the hidden state is essential for controlling short-term temporal information, while the cell state stores the long-term temporal information. Therefore, LSTM achieves good performance in prediction by capturing short- and long-term temporal information. Besides, LSTM only accepts vectors as the input in each time interval. As a result, the information captured by the irregular convolution are flattened as a vector then input to LSTM model (shown as Fig. A.9).

References

- [1] P. Midgley, Bicycle-sharing schemes: enhancing sustainable mobility in urban areas, United Nations, Department of Economic and Social Affairs 8 (2011) 1–12.
- [2] J. Larsen, Bike-sharing programs hit the streets in over 500 cities worldwide, Earth Policy Institute Washington, DC, 2013.
- [3] T. Litman, D. Burwell, Issues in sustainable transportation, *International Journal of Global Environmental Issues* 6 (4) (2006) 331–347.
- [4] L. Steg, R. Gifford, Sustainable transportation and quality of life, *Journal of transport geography* 13 (1) (2005) 59–69.
- [5] H. Haghshenas, M. Vaziri, Urban sustainable transportation indicators for global comparison, *Ecological Indicators* 15 (1) (2012) 115–121.
- [6] J. Jobe, G. P. Griffin, Bike share responses to covid-19, *Transportation Research Interdisciplinary Perspectives* 10 (2021) 100353.
- [7] S. Hu, C. Xiong, Z. Liu, L. Zhang, Examining spatiotemporal changing patterns of bike-sharing usage during covid-19 pandemic, *Journal of transport geography* 91 (2021) 102997.
- [8] K. Kim, Impact of covid-19 on usage patterns of a bike-sharing system: case study of seoul, *Journal of transportation engineering, Part A: Systems* 147 (10) (2021) 05021006.
- [9] T. Raviv, O. Kolka, Optimal inventory management of a bike-sharing station, *Iie Transactions* 45 (10) (2013) 1077–1093.
- [10] M. Dell’Amico, M. Iori, S. Novellani, A. Subramanian, The bike sharing rebalancing problem with stochastic demands, *Transportation research part B: methodological* 118 (2018) 362–380.
- [11] M. Dell’Amico, E. Hadjicostantinou, M. Iori, S. Novellani, The bike sharing rebalancing problem: Mathematical formulations and benchmark instances, *Omega* 45 (2014) 7–19.
- [12] D. Singhvi, S. Singhvi, P. I. Frazier, S. G. Henderson, E. O’Mahony, D. B. Shmoys, D. B. Woodard, Predicting bike usage for new york city’s bike sharing system, in: *Workshops at the twenty-ninth AAAI conference on artificial intelligence*, 2015.
- [13] K. Kumar, V. K. Jain, Autoregressive integrated moving averages (arima) modelling of a traffic noise time series, *Applied Acoustics* 58 (3) (1999) 283–294.
- [14] R. Avuglah, K. Adu-Poku, E. Harris, Application of arima models to road traffic accident cases in ghana, *International journal of statistics and applications* 4 (5) (2014) 233–239.
- [15] D. Billings, J. S. Yang, Application of the arima models to urban roadway travel time prediction—a case study, in: *2006 IEEE International Conference on Systems, Man and Cybernetics, Vol. 3, IEEE, 2006*, pp. 2529–2534.
- [16] S. Lee, D. B. Fambro, Application of subset autoregressive integrated moving average model for short-term freeway traffic volume forecasting, *Transportation Research Record* 1678 (1) (1999) 179–188.
- [17] Y. Ai, Z. Li, M. Gan, Y. Zhang, D. Yu, W. Chen, Y. Ju, A deep learning approach on short-term spatiotemporal distribution forecasting of dockless bike-sharing system, *Neural Computing and Applications* 31 (5) (2019) 1665–1677.

- [18] Y. Pan, R. C. Zheng, J. Zhang, X. Yao, Predicting bike sharing demand using recurrent neural networks, *Procedia computer science* 147 (2019) 562–566.
- [19] Y. Li, R. Yu, C. Shahabi, Y. Liu, Diffusion convolutional recurrent neural network: Data-driven traffic forecasting, *arXiv preprint arXiv:1707.01926* (2017).
- [20] J. Zhang, Y. Zheng, D. Qi, Deep spatio-temporal residual networks for citywide crowd flows prediction, in: *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [21] X. Li, Y. Xu, Q. Chen, L. Wang, X. Zhang, W. Shi, Short-term forecast of bicycle usage in bike sharing systems: A spatial-temporal memory network, *IEEE Transactions on Intelligent Transportation Systems* (2021).
- [22] Y. Zhang, T. Cheng, Y. Ren, K. Xie, A novel residual graph convolution deep learning model for short-term network-based traffic forecasting, *International Journal of Geographical Information Science* 34 (5) (2020) 969–995.
- [23] Y. Ren, H. Chen, Y. Han, T. Cheng, Y. Zhang, G. Chen, A hybrid integrated deep learning model for the prediction of citywide spatio-temporal flow volumes, *International Journal of Geographical Information Science* 34 (4) (2020) 802–823.
- [24] V. Sathishkumar, J. Park, Y. Cho, Using data mining techniques for bike sharing demand prediction in metropolitan city, *Computer Communications* 153 (2020) 353–366.
- [25] Y. Yang, A. Heppenstall, A. Turner, A. Comber, Using graph structural information about flows to enhance short-term demand prediction in bike-sharing systems, *Computers, Environment and Urban Systems* 83 (2020) 101521.
- [26] X. Ma, Z. Dai, Z. He, J. Ma, Y. Wang, Y. Wang, Learning traffic as images: a deep convolutional neural network for large-scale transportation network speed prediction, *Sensors* 17 (4) (2017) 818.
- [27] W. Xiangxue, X. Lunhui, C. Kaixun, Data-driven short-term forecasting for urban road network traffic based on data processing and lstm-rnn, *Arabian Journal for Science and Engineering* 44 (4) (2019) 3043–3060.
- [28] R. Fu, Z. Zhang, 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)*, IEEE, 2016, pp. 324–328.
- [29] Z. Zhao, W. Chen, X. Wu, P. C. Chen, J. Liu, Lstm network: a deep learning approach for short-term traffic forecast, *IET Intelligent Transport Systems* 11 (2) (2017) 68–75.
- [30] B. Du, H. Peng, S. Wang, M. Z. A. Bhuiyan, L. Wang, Q. Gong, L. Liu, J. Li, Deep irregular convolutional residual lstm for urban traffic passenger flows prediction, *IEEE Transactions on Intelligent Transportation Systems* 21 (3) (2019) 972–985.
- [31] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, *nature* 521 (7553) (2015) 436–444.
- [32] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural computation* 9 (8) (1997) 1735–1780.
- [33] Y. Du, F. Deng, F. Liao, A model framework for discovering the spatio-temporal usage patterns of public free-floating bike-sharing system, *Transportation Research Part C: Emerging Technologies* 103 (2019) 39–55.

- [34] Y. Xu, D. Chen, X. Zhang, W. Tu, Y. Chen, Y. Shen, C. Ratti, Unravel the landscape and pulses of cycling activities from a dockless bike-sharing system, *Computers, Environment and Urban Systems* 75 (2019) 184–203.
- [35] A. M. Nagy, V. Simon, Survey on traffic prediction in smart cities, *Pervasive and Mobile Computing* 50 (2018) 148–163.
- [36] D. Xu, Y. Wang, L. Jia, Y. Qin, H. Dong, Real-time road traffic state prediction based on arima and kalman filter, *Frontiers of Information Technology & Electronic Engineering* 18 (2) (2017) 287–302.
- [37] C. Antoniou, M. Ben-Akiva, H. N. Koutsopoulos, Nonlinear kalman filtering algorithms for on-line calibration of dynamic traffic assignment models, *IEEE Transactions on Intelligent Transportation Systems* 8 (4) (2007) 661–670.
- [38] M. W. Szeto, D. C. Gazis, Application of kalman filtering to the surveillance and control of traffic systems, *Transportation Science* 6 (4) (1972) 419–439.
- [39] H. van Lint, T. Djukic, Applications of kalman filtering in traffic management and control, in: *New Directions in Informatics, Optimization, Logistics, and Production*, informs, 2012, pp. 59–91.
- [40] Y. Zhang, Y. Liu, Traffic forecasting using least squares support vector machines, *Transportmetrica* 5 (3) (2009) 193–213.
- [41] Y. Hou, P. Edara, Y. Chang, Road network state estimation using random forest ensemble learning, in: *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, IEEE, 2017, pp. 1–6.
- [42] S. Sun, C. Zhang, G. Yu, A bayesian network approach to traffic flow forecasting, *IEEE Transactions on intelligent transportation systems* 7 (1) (2006) 124–132.
- [43] S. Sun, C. Zhang, Y. Zhang, Traffic flow forecasting using a spatio-temporal bayesian network predictor, in: *International conference on artificial neural networks*, Springer, 2005, pp. 273–278.
- [44] Z. Wang, X. Su, Z. Ding, Long-term traffic prediction based on lstm encoder-decoder architecture, *IEEE Transactions on Intelligent Transportation Systems* (2020).
- [45] W. Zi, W. Xiong, H. Chen, L. Chen, Tagcn: Station-level demand prediction for bike-sharing system via a temporal attention graph convolution network, *Information Sciences* 561 (2021) 274–285.
- [46] Q. Liu, B. Wang, Y. Zhu, Short-term traffic speed forecasting based on attention convolutional neural network for arterials, *Computer-Aided Civil and Infrastructure Engineering* 33 (11) (2018) 999–1016.
- [47] C. Zhang, H. Zhang, D. Yuan, M. Zhang, Citywide cellular traffic prediction based on densely connected convolutional neural networks, *IEEE Communications Letters* 22 (8) (2018) 1656–1659.
- [48] D. Shin, K. Chung, R. Park, Prediction of traffic congestion based on lstm through correction of missing temporal and spatial data, *IEEE Access* 8 (2020) 150784–150796.
- [49] J. Zhang, F. Chen, Z. Cui, Y. Guo, Y. Zhu, Deep learning architecture for short-term passenger flow forecasting in urban rail transit, *IEEE Transactions on Intelligent Transportation Systems* (2020) 1–11doi:10.1109/TITS.2020.3000761.

- [50] T. S. Kim, W. K. Lee, S. Y. Sohn, Graph convolutional network approach applied to predict hourly bike-sharing demands considering spatial, temporal, and global effects, *PloS one* 14 (9) (2019) e0220782.
- [51] L. Lin, Z. He, S. Peeta, Predicting station-level hourly demand in a large-scale bike-sharing network: A graph convolutional neural network approach, *Transportation Research Part C: Emerging Technologies* 97 (2018) 258–276.
- [52] Y. Wu, H. Tan, L. Qin, B. Ran, Z. Jiang, A hybrid deep learning based traffic flow prediction method and its understanding, *Transportation Research Part C: Emerging Technologies* 90 (2018) 166–180.
- [53] J. Benesty, J. Chen, Y. Huang, I. Cohen, Pearson correlation coefficient, in: *Noise reduction in speech processing*, Springer, 2009, pp. 1–4.
- [54] L. Kristoufek, Measuring correlations between non-stationary series with dcca coefficient, *Physica A: Statistical Mechanics and its Applications* 402 (2014) 291–298.
- [55] M. Müller, Dynamic time warping, *Information retrieval for music and motion* (2007) 69–84.
- [56] J. Taylor, X. Zhou, N. M. Rouphail, R. J. Porter, Method for investigating intradriver heterogeneity using vehicle trajectory data: A dynamic time warping approach, *Transportation Research Part B: Methodological* 73 (2015) 59–80.
- [57] D. Yang, S. Li, Z. Peng, P. Wang, J. Wang, H. Yang, Mf-cnn: traffic flow prediction using convolutional neural network and multi-features fusion, *IEICE TRANSACTIONS on Information and Systems* 102 (8) (2019) 1526–1536.
- [58] Y. Liu, Z. Liu, R. Jia, Deepff: A deep learning based architecture for metro passenger flow prediction, *Transportation Research Part C: Emerging Technologies* 101 (2019) 18–34.
- [59] T. Y. Kim, S. B. Cho, Predicting residential energy consumption using cnn-lstm neural networks, *Energy* 182 (2019) 72–81.
- [60] M. Cao, V. O. Li, V. W. Chan, A cnn-lstm model for traffic speed prediction, in: *2020 IEEE 91st Vehicular Technology Conference (VTC2020-Spring)*, IEEE, 2020, pp. 1–5.
- [61] S. Wang, W. Liu, J. Wu, L. Cao, Q. Meng, P. J. Kennedy, Training deep neural networks on imbalanced data sets, in: *2016 international joint conference on neural networks (IJCNN)*, IEEE, 2016, pp. 4368–4374.
- [62] A. Graves, Generating sequences with recurrent neural networks, *arXiv preprint arXiv:1308.0850* (2013).
- [63] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala, Pytorch: An imperative style, high-performance deep learning library, in: H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 32, Curran Associates, Inc., 2019, pp. 8024–8035.