# Short-term Road Traffic Prediction based on Deep Cluster at Large-scale Networks

Lingyi Han, Student Member, IEEE, Kan Zheng, Senior Member, IEEE, Long Zhao, Senior Member, IEEE, Xianbin Wang, Fellow, IEEE, and Xuemin Shen, Fellow, IEEE

Abstract—Short-term road traffic prediction (STTP) is one of the most important modules in Intelligent Transportation Systems (ITS). However, network-level STTP still remains challenging due to the difficulties both in modeling the diverse traffic patterns and tacking high-dimensional time series with low latency. Therefore, a framework combining with a deep clustering (DeepCluster) module is developed for STTP at largescale networks in this paper. The DeepCluster module is proposed to supervise the representation learning in a visualized way from the large unlabeled dataset. More specifically, to fully exploit the traffic periodicity, the raw series is first split into a number of sub-series for triplets generation. The convolutional neural networks (CNNs) with triplet loss are utilized to extract the features of shape by transferring the series into visual images. The shape-based representations are then used for road segments clustering. Thereafter, motivated by the fact that the road segments in a group have similar patterns, a model sharing strategy is further proposed to build recurrent NNs (RNNs)based predictions through a group-based model (GM), instead of individual-based model (IM) in which one model are built for one road exclusively. Our framework can not only significantly reduce the number of models and cost, but also increase the number of training data and the diversity of samples. In the end, we evaluate the proposed framework over the network of Liuli Bridge in Beijing. Experimental results show that the DeepCluster can effectively cluster the road segments and GM can achieve comparable performance against the IM with less number of models.

Index Terms—Short-term traffic prediction, Large-scale networks, Deep representation learning, Shape-based features

#### I. Introduction

The short-term road traffic prediction (STTP) technique has been studied in achieving efficient route planning and traffic control in Intelligent Transportation Systems (ITS) recently [1]. The main idea of STTP is to predict the road traffic state (i.e., flow, speed and density) in the next five to thirty minutes by analyzing historical data [2]. However, existing STTP studies mainly focused on one road segment, or a small-scale network containing several adjacent road segments, which is opposite to the effective route planning that requires a global perspective based on the information of the

L. Han, K. Zheng, and L. Zhao are with the Intelligent Computing and Communication ( $\rm IC^2$ ) Lab, Wireless Signal Processing and Networks Lab (WSPN) Key Lab of Universal Wireless Communications, Ministry of Education, Beijing University of Posts & Telecommunications, Beijing, China, 100088. Contact email: zkan@bupt.edu.cn

X. Wang is with Department of Electrical and Computer Engineering, University of Western Ontario, London, ON N6A 5B9, Canada.

X. Shen is with Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON N2L 3G1, Canada.

Corresponding author: Kan Zheng.

whole network [3]–[5]. Besides, the majority of existing STTP algorithms are limited to a single scenario such as freeway, arterial or corridor, which are difficult to be generalized to a heterogeneous road network. The past STTP method for large-scale road network is to develop a specific model for each road segment termed as individual-based model (IM), or a general model for all road segments termed as whole-based model (WI). Since the multiplicity and heterogeneity of the large-scale network, neither of the two models is appropriate for the large-scale networks. Firstly, too many IMs will take up lots of storage resources in ITS. Secondly, a WI is not competent for modeling the whole network with different types of traffic patterns. Moreover, the development of ITS over the city increases the number of traffic data in terms of time span and granularity [6]. Making full use of the big traffic data to improve the performance of the prediction becomes a challenge. Therefore, a feasible STTP at large-scale network needs to be studied.

Generally, representation learning, a.k.a. dimension reduction, is used to transform the raw data into a good representation that makes the subsequent tasks easy. It plays an important role in time series clustering, because time series are essentially high-dimensional and susceptible to noise. Hence, clustering directly with raw series is computationally expensive and distance measures are highly sensitive to the distortions. Recently, deep learning (DL) has been developed with great success in many areas, including computer vision, speech recognition and natural language processing due to its theoretical function approximation properties [7] and demonstrated feature learning capabilities [8]. Therefore, deep representation is used for traffic series clustering.

In this paper, a feasible framework composed of a deep clustering module and several prediction models is proposed for STTP at large-scale networks. More specifically, a shape-based representation learning method is developed for road segments clustering. On the other hand, several predictions are combined to achieve the STTP at the network. The main contributions of the paper are summarized as follows:

- By fully exploiting the periodicity of traffic patterns, we propose a method to generate triplets from unlabeled dataset. The raw traffic series are divided into sub-series by periods, three of which are selected to generate a triplet according to a specific criterion. The dimension of sub-series used for representation learning is significantly reduced, compared to raw series.
- A supervised deep clustering module termed as DeepCluster, is developed. Unlike the existing hand-

craft features, such as the frequency transformation, wavelet transformation, Shapelets  $et\ al.$ , a pure data driven method is proposed to learn the shape-based representations of traffic series in a visualized way. A rasterization strategy is first designed to transform the traffic series into traffic images. A convolutional neural network (CNN) with triplet loss is then used for representation learning. At last, the representations are used to cluster the network into K groups by traditional clustering methods.

• Based on the idea of model sharing, K group-based models (GMs) that are constituting a prediction at network is proposed to achieve a good tradeoff between the quantity of models and the performance of predictions. Specifically, all road segments in one group share one prediction model and each GM allows the training samples generated by the road segment from the same group to be aggregated to learn the model. Model sharing increases the number and the diversity of the training samples, which is beneficial for DNNs training. The experiment results validate that the GM has stronger generalization ability than IM. We also analyze the impact of input interval on performance by experiments.

The rest of paper is organized as follows. Section II reviews the related works. In Section III, the data used throughout the paper is described. Section IV formulates the STTP problem at large-scale network. In Section V, the DL methodologies are introduced. The proposed framework of STTP including Deep-Cluster and DeepPrediction is then proposed in Section VI. In Section VII, simulation results demonstrating the performance of the proposed framework are given, before concluding the paper in Section VIII.

# II. RELATED WORKS

# A. Time Series Representation Learning

A wide colorvariety of methods had been developed for time series representation learning in clustering [9]–[11], such as spectral transformation [12], wavelets transformation [12], eigenvalue analysis techniques [13], piecewise linear approximation (PLA) [14], adaptive piecewise constant approximation (APCA) [15], symbolic approximation (SAX) [16], piecewise aggregate approximation (PAA) [17], perceptually important point (PIP) [18] *et al.* However, all these methods are hand-craft features, which are designed to describe specific time series pattern and heavily rely on the database.

A new trend appears with artificial neural networks (ANNs), especially deep NNs (DNNs) based representation learning in clustering, which are data-driven and capable of learning a powerful representation from raw data through a high-level and non-linear mapping. Therefore, some works have used the deep representation learning to improve clustering performance. C. Song *et al.* in [19] integrated *K*-means algorithm into a stacked auto-encoder (SAE) by minimizing the reconstruction error as well as the distance between data

points and corresponding clusters. It alternatively learned the representations and updated cluster centers. In [20], [21], the k-means algorithm used the nonlinear representations that are learned by DNNs for clustering. J. Xie et al. in [22] proposed a deep embedded clustering that simultaneously learned the representations and cluster assignments by defining a centroidbased probability distribution and minimizing its Kullback-Leibler (KL) divergence to an auxiliary target distribution. K. Tian et al. in [23] improved the existing works by proposing a general flexible framework that integrated traditional clustering methods into different DNNs. The framework is optimized by alternating direction of multiplier method (ADMM). However, the above methods all worked with the static data that is simple and low dimensional compared with time series data in general. On the other hand, there is less research on the deep representation learning of time series in clustering. Therefore, an efficient time series representation learning algorithm dedicated for clustering needs to be developed.

# B. Short Term Traffic Prediction

There are numerous researches on single-point STTP [3], such as autoregressive integrated moving average (ARIMA) family of models, regression models, Markov models, Kalman filters, Bayesian networks, traffic flow theory-based simulation models and ANNs. Obviously, single-point models predict the future traffic state for a target road segment only using its own historical data, which ignores the relations between the target road segment and adjacent segments. Consequently, some researches have focused on predicting one or multiple segments by taking the spatio-temporal interrelations between adjacent road segments into account [24]–[27]. However, the above network-level STTP researches are restricted to small regions that containing several adjacent road segments.

Recently, a few literatures begin to pay attention to the predictions at the large-scale networks. In [29], [30], dynamic simulator based on traffic flow theory was used for STTP at the whole network with limited traffic data. [31]-[34] only predicted the traffic state of the representative road subset to achieve the prediction at the whole network by utilizing data compression technologies. However, the performance of prediction was poor resulted from compression and reconstruction errors. Min et al. in [35] considered a road network consists of about 500 road segments. However, they developed a custom model for the test area, which is not practical. M. Asif et al. in [36] performed prediction for each individual road segment with support vector regression (SVR) algorithm over a large network containing 5,000 road segments. Then K-means algorithm was used to cluster the road segments to analyze the spatial prediction performance. But the prediction method may not work well, since the performances differed greatly among clusters and the mean error of one cluster is up to 17.18% of five-minute prediction. Besides, STTP for each individual road segment is hard to implement on largescale networks in practice. X. Ma et al. in [37] proposed a CNN-based method that arranges the traffic data into 2D (2dimensional) matrices as inputs to predict the large-scale traffic speeds. However, they only built one model and expected it to

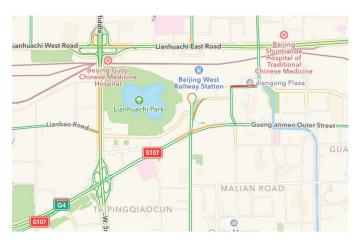


Fig. 1. The topology of the network at Liuli Bridge, Beijing.

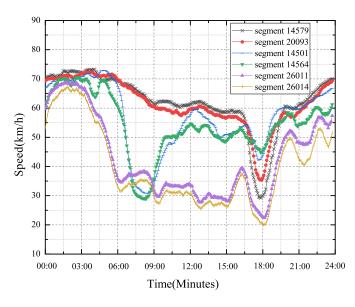


Fig. 2. The every five-minute average traffic speeds of six road segments from September, 2017 to November, 2017.

fit for all segments without considering the fact that the whole network is heterogeneous with different type of segments. Therefore, these attempts are hard to be implemented on largescale networks with high accuracy.

# III. THE DATA

The traffic data used throughout the paper is described in this section. The topology of Liuli Bridge is shown in Fig. 1. The network consists of about 1,000 road segments with a diverse level of road functions including express way, arterial road, access road, side road  $\it et~al.$  In addition, the dataset collected by Beijing Transportation Institute contains the traffic speed data from September, 2017 to November, 2017 with five-minute sampling interval. Hence, it has totally  $90\times288$  measured data, where 90 means the total number of days and 288 means the number of values collected in each day. The data is measured by vehicles that are equipped with GPS such as taxis and buses.

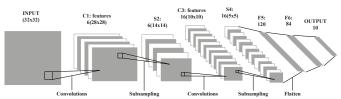


Fig. 3. Architecture of LeNet-5. Above the rectangles are the number of channels and its size in parenthesis.

# IV. FORMULATION OF STTP PROBLEM

Consider a large-scale network  $\Phi$  consisting of  $N_r$  road segments, i.e.,  $\Phi = \{\boldsymbol{x}^{(r)}\}_{r=1}^{N_r}$ , where  $\boldsymbol{x}^{(r)} = [x_1^{(r)}, x_2^{(r)}, \dots, x_{N_t}^{(r)}]$  is a time series of  $N_t$  measurements at segment r. We denote a sub traffic series by

$$\mathbf{x}_{t:L:l} = [x_t, x_{t+l}, \dots, x_{t+(L-1)l}],$$
 (1)

where  $\boldsymbol{x}_{t:L:l}$  is a set of L continuous measured values with intervals l from a time series  $\boldsymbol{x}$ , that starts at position t with  $1 \le t \le N_t$ ,  $1 \le l \le N_t$  and  $1 \le L \le N_t$ .  $\boldsymbol{x}_{t:L:1}$  is abbreviated as  $\boldsymbol{x}_{t:L}$  for simplicity.

Let  $\hat{x}_{t+N_0}$  be the forecast of traffic state of the prediction horizon  $N_0$ , given the corresponding  $N_i$  historical measurements up to time t. The goal of STTP is to construct a mapping  $f(\cdot)$  between the historical traffic state and the future one, i.e.,

$$\hat{x}_{t+N_o} = f(x_{t-N_i}, x_{t-N_i+1}, \dots, x_t)$$

$$= f(x_{t-N_i:N_i}).$$
(2)

As stated above, IM and WM are both inappropriate for the large-scale networks, because they not only consist of a large number of road segments, but also a variety of types of road segments as shown in Fig. 2. On one hand, it's unpractical to construct and store massive amounts of IMs in ITS. Besides, the number of training samples collected from one segment is insufficient to learn a robust DL model. On the other hand, it's impossible to build a model for the whole network with different types of traffic pattens. In addition, the model is vulnerable to the curse of dimensionality by taking historical data from all segments as inputs. Then how to make a proper utilization of the tremendous traffic data to achieve the effective and practical STTP is still a problem.

To tackle this problem, we cluster the road segments into groups, each of which has a typical traffic pattern. Within each group, the traffic patterns of all road segments are highly similar in shape. Based on that, a STTP model is built for a group, rather than a segment or whole network. The challenges in our problem include *i*) representation learning of the traffic series that are high-dimensional and sensitive to distortion, and *ii*) representation learning from unlabeled traffic data that are beneficial to cluster task.

## V. DEEP LEARNING FOR TRAFFIC PREDICTIONS

In this section, we deals with the tremendous traffic series by means of the DL technologies, including CNNs and recurrent NNs (RNNs), which will be explained in this section.

# A. Convolutional Neural Networks

The key aspect of CNNs is that the features are not designed by human engineers, but are learned from data using a general-purpose learning procedure [8]. Fig. 3 shows the architecture of a typical CNN, named LeNet-5. CNNs can take any form of arrays, such as 1D series, 2D images and 3D videos as inputs. A CNN is made up of layers, where two main types of layers different with the regular ANNs are convolutional layers (C layers in Fig. 3) and subsampling layers (S layers in Fig. 3).

In the l-th convolutional layer  $C^{(l)}$ , the outputs of the previous layer are fed to convolve with several convolutional kernels. After that, the outputs are added by biases and activated by a nonlinear function to form new representations (features in Fig. 3) for the next layer. Assuming the current layer accept an input volume  $O^{(l-1)}$  of size  $W_o^{(l-1)} \times H_o^{(l-1)} \times D_o^{(l-1)}$ . Formally, the output  $O^{(k,l)}$  of size  $W_o^{(l)} \times H_o^{(l)} \times 1$  filtered by the k-th kernel  $K^{(k,l)}$  of size  $W_k^{(l)} \times H_k^{(l)} \times 1$  with stride s is given by

$$O^{(k,l)} = g(K^{(k,l)} \otimes O^{(l-1)} + b^{(k,l)}), k = 1, 2, \dots, N_k^{(l)}$$
 (3)

where  $N_{\mathbf{k}}^{(l)}$  is the number of kernels and  $b^{(k,l)}$  is a bias of layer  $\mathbf{C}^{(l)}$ , respectively.  $\otimes$  represents a discrete convolution operator .  $g(\cdot)$  is a activation function such as  $\tanh$  function, relu function et al. By concatenating  $\mathbf{O}^{(k,l)}$  along the last dimension, the output  $\mathbf{O}^{(l)}$  for layer  $\mathbf{C}^{(l)}$  of the size  $W_{\mathbf{0}}^{(l)} \times H_{\mathbf{0}}^{(l)} \times D_{\mathbf{0}}^{(l)}$  can be derived, and both of which can be calculated by

$$W_{\rm o}^{(l)} = \lfloor \frac{W_{\rm o}^{(l-1)} - W_{\rm k}^{(l)}}{s} \rfloor + 1,$$
 (4)

$$H_{\rm o}^{(l)} = \lfloor \frac{H_{\rm o}^{(l-1)} - H_{\rm k}^{(l)}}{s} \rfloor + 1,$$
 (5)

$$D_{\rm o}^{(l)} = N_{\rm k}^{(l)},\tag{6}$$

where  $\lfloor \cdot \rfloor$  represents rounded down. With parameter sharing, there are  $N_{\rm w}^{(l)}$  learnable weights of layer  $C^{(l)}$  in total,

$$N_{\rm w}^{(l)} = W_{\rm k}^{(l)} \times H_{\rm k}^{(l)} \times D_{\rm o}^{(l-1)} \times N_{\rm k}^{(l)} + N_{\rm k}^{(l)}. \tag{7}$$

In the (l+1)-th subsampling layer  $S^{(l+1)}$ , the spatial resolution of representations is reduced to increase the level of distortion invariance. After layer  $C^{(l)}$ , the layer  $S^{(l+1)}$  accepts a volume of size  $W^{(l)}_{\rm o} \times H^{(l)}_{\rm o} \times D^{(l)}_{\rm o}$  as input. Specifically, representations in the previous layer are pooled over neighborhood within a rectangular region of  $W^{(l+1)}_{\rm s} \times H^{(l+1)}_{\rm s}$ , by either a max-pooling function

$$\begin{split} O_{i,j,k}^{(l+1)} &= \max_{\substack{i \leq p \leq i + W_{s}^{(l+1)} \\ j \leq q \leq j + H_{s}^{(l+1)}}} (O_{p,q,k}^{(l)}), \\ &i \leq W_{o}^{(l+1)}, j \leq H_{o}^{(l+1)}, k \leq D_{o}^{(l+1)} \end{split} \tag{8}$$

or an average-pooling function, where  $O^{(l+1)}$  is the output of size  $W_{\rm o}^{(l+1)} imes H_{\rm o}^{(l+1)} imes D_{\rm o}^{(l+1)}$ , and both of them can be calculated by

$$W_{o}^{(l+1)} = \lfloor \frac{W_{o}^{(l)} - W_{s}^{(l+1)}}{s} \rfloor + 1, \tag{9}$$

$$H_{\rm o}^{(l+1)} = \lfloor \frac{H_{\rm o}^{(l)} - H_{\rm s}^{(l+1)}}{s} \rfloor + 1,$$
 (10)

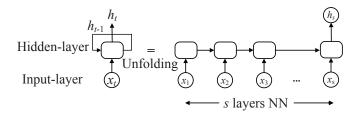


Fig. 4. Architecture of a basic three-layer RNN.

$$D_{0}^{(l+1)} = D_{0}^{(l)}. (11)$$

The convolutional and subsampling operators make the new representations more invariance to the distortion compared to the raw data. Besides, the parameter sharing make the CNNs capable of processing high-dimensional inputs. The aforementioned characteristics allow to adopt the CNNs for time series representation learning. In this section, we explore an efficient deep CNN architecture, FaceNet [39] to learn the deep representations of the raw time series.

#### B. Recurrent Neural Networks

Unlike the regular ANNs, RNNs are capable of exhibiting the temporal correlations of time series, which makes them applicable to tasks such as language modeling, speech recognition or time series forecasting.

Assuming the duration of the temporal correlations (defined as time step) is s, a three-layer RNN can be regarded as a s-layer feed-forward NN by unfolding it through time, As shown in Fig. 4. The RNN reads a series  $x_{1:s}$  one by one and each RNN block takes a value at one time as input. The current hidden state  $h_t$  at time t is computed from the current input  $x_t$  and the previous hidden state  $h_{t-1}$  by

$$h_t = g(x_t, h_{t-1})$$
  
=  $g(x_t, g(x_{t-1}, h_{t-2}))$   
=  $\cdots$  (12)

where  $h_{t-2}$  is the hidden state of the last two RNN blocks.  $g(\cdot)$  is the activation function of the hidden layer. The key idea of RNNs is to imitate a sequential dynamic behavior with a chain-like structure that allows the information to be passed from previous layer to the current one. In this paper, RNNs are used to model the temporal correlations of traffic series.

# VI. PROPOSED FRAMEWORK FOR STTP

In this section, a framework dedicated for STTP at large-scale networks is described in details. The architecture of this framework is shown in Fig. 5. It consists of two major components, i.e., DeepCluster and DeepPrediction. The inputs are historical traffic states with fixed interval coming from different road segments, while the outputs are predictions for a given time period. The inputs are fed to the DeepCluster module, and are divided into several groups. Afterwards, the DeepPrediction module performs the predictions for the network.

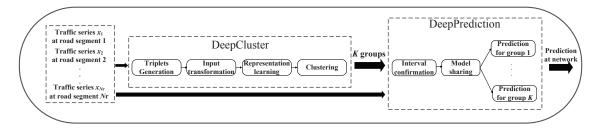


Fig. 5. The block diagram of representation learning.

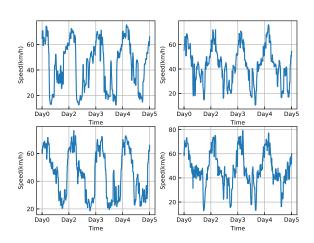


Fig. 6. The every five-minute average traffic speeds of four road segments on different days of the week from September, 2017 to November, 2017.

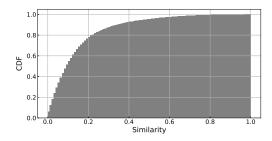


Fig. 7. The CDF of similarity with real traffic speed.

#### A. DeepCluster

Traffic series clustering method at large-scale networks is first proposed, which is implemented via deep representation learning. Before developing the clustering algorithm, the problem of clustering at large-scale networks is formally defined as follows:

Definition 1: Given a large-scale network  $\Phi$  consists of  $N_r$  traffic series, i.e.,  $\Phi = \{\boldsymbol{x}^{(r)}\}_{r=1}^{N_r}$  the process of partitioning of  $\Phi$  into K groups  $\{\boldsymbol{C}^{(1)}, \boldsymbol{C}^{(2)}, \ldots, \boldsymbol{C}^{(K)}\}$ , is called *traffic series clustering*. In such a way that homogenous traffic series are grouped together based on a certain similarity measure.

In contrast to the traditional extrinsic hand-craft features, human brains can seize the intrinsic visual-based features easily, which is why they can quickly distinguish different types of the time series under the help of high abstraction

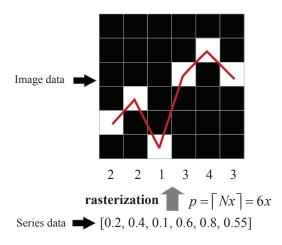


Fig. 8. The schematic diagram of the inputs transformation.

ability. Moreover, compared with raw time series, the intrinsic visual-based features are more steady. They are less affected by the distortions and the scale of samples. To address the issues of the raw data-based or hand-craft-based clustering methods, we use the deep representation learning for series clustering. The DNN is employed to learn a mapping from the raw high-dimensional traffic series to the low-dimensional representations that are used for clustering.

The DeepCluster module includes triplets generation, inputs transformation, representation learning and clustering. Details of each step are given below.

1) Triplets Generation. As can be seen in Fig. 6, the traffic patterns follow the same trend among days. In order to study the traffic periodic pattern in a day, we calculate the traffic similarity defined in [38]. The traffic similarity is defined as the normalized gaps between each pair of measurements in two consecutive days from one road segment. As stated in Sec. IV, traffic speeds are collected every 5 minutes. Since one day has 288 time intervals, the traffic similarity SIM<sub>t\*</sub> at segment r in time slot t\* can be calculated by

$$\operatorname{SIM}_{t^*}^{(r)} = \frac{|x_{t^*}^{(r)} - x_{t^*+288}^{(r)}|}{\max_{1 \le t \le N_t - 288} |x_t^{(r)} - x_{t+288}^{(r)}|}.$$
 (13)

The cumulative distribution function (CDF) of  $SIM_{t^*}^{(r)}$  is shown in Fig. 7. We can see that more than 80%  $SIM_{t}^{(r)}$  are smaller than 0.2, which indicates that periodic pattern exists in traffic series at most read segments.

To fully exploit the traffic temporal features and periodic patterns, we split the traffic series into sub-series by periods, and generate triplets for representation learning. Given  $N_{\rm r}$  traffic series with period  $N_{\rm p}$  measured from  $N_{\rm r}$  road segments, we split the series into sub-series by periods, termed as periodic sub-series. Thus, we have  $d=N_{\rm t}/N_{\rm p}$  periodic sub-series for each segment,

$$\boldsymbol{x}_{1:N_{p}}^{(r)}, \boldsymbol{x}_{N_{p}+1:N_{p}}^{(r)}, \dots, \boldsymbol{x}_{(d-1)N_{p}+1:N_{p}}^{(r)},$$
 (14)

here  $\boldsymbol{x}_{jN_{\mathrm{p}}+1:N_{\mathrm{p}}}^{(r)}=[x_{jN_{\mathrm{p}}+1}^{(r)},x_{jN_{\mathrm{p}}+2}^{(r)},\ldots,x_{(j+1)N_{\mathrm{p}}}^{(r)}]$  is the (j+1)-th periodic sub-series at segment r with  $0\leq j\leq d-1$ . A triplet is made up by randomly choosing two different periodic sub-series from one segment, and one sub-series from another segment,

$$\{\boldsymbol{x}_{iN_{\rm p}+1:N_{\rm p}}^{(r_i)}, \boldsymbol{x}_{jN_{\rm p}+1:N_{\rm p}}^{(r_i)}, \boldsymbol{x}_{kN_{\rm p}+1:N_{\rm p}}^{(r_j)}\}.$$

$$0 \le i, j, k \le d-1, i \ne j, r_i \ne r_j$$

$$(15)$$

2) **Inputs transformation.** In order to extract the features of shape, a rasterization strategy is designed to visualize the series into images shown in Fig. 8. The transformed images can reveal the shape information of series well, such as bulge, sink and so on. Let the series  $\mathbf{x} = [x_1, x_2, \dots, x_N]$  be standardized by min-max normalization to keep values between 0 and 1. A series is transformed to a matrix by expanding each element to a vector. For the *i*-th element  $x_i$ , the position  $p_i$  at the *i*-th column of the matrix is,

$$p_i = [Nx_i].$$
  $p_i \in \{1, 2, \dots, N\}$  (16)

The matrix  $X_{N,N}$  corresponding to the series x can be written as:

$$X_{NN} = [255_N(p_1), 255_N(p_2), \dots, 255_N(p_N)], (17)$$

where  $255_N(p_i)$  is a N-dimensional vector with the pixel value of 255 at its i-th entry standing for white and 0 standing for black elsewhere. The transformed image is shown in Fig. 9. The matrixes are used as the inputs to the representation learning. The sub-image corresponding to the sub-series  $\boldsymbol{x}_{iN_p+1:N_p}^{(r)}$  is represented as  $\boldsymbol{X}_i^{(r)}$ . Therefore, the triplet becomes,

$$\{\boldsymbol{X}_{i}^{(r_{i})}, \boldsymbol{X}_{j}^{(r_{i})}, \boldsymbol{X}_{k}^{(r_{j})}\}.$$

$$0 \le i, j, k \le d - 1, i \ne j, r_{i} \ne r_{j}$$
(18)

3) Representation learning and clustering. DNNs with triplet loss from [39] is employed to strive for a representations over a triplet, from an image space into a feature space. The triplet loss encourages the representations of a pair of sub-images from one segment to be close to each other in the feature space, and the those from different segments to be far away. The representation of x is denoted by f(x). Thus, the triplet loss that is being minimized is,

$$||f(\boldsymbol{X}_{i}^{(r_{i})}) - f(\boldsymbol{X}_{j}^{(r_{i})})||_{2}^{2} - ||f(\boldsymbol{X}_{i}^{(r_{i})}) - f(\boldsymbol{X}_{k}^{(r_{j})})||_{2}^{2},$$

$$\forall \{\boldsymbol{X}_{i}^{(r_{i})}, \boldsymbol{X}_{j}^{(r_{i})}, \boldsymbol{X}_{k}^{(r_{j})}\} \in \Gamma$$
(19)

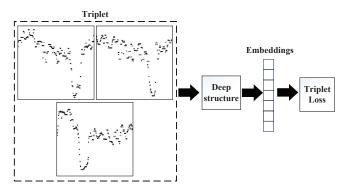


Fig. 9. The block diagram of representation learning.

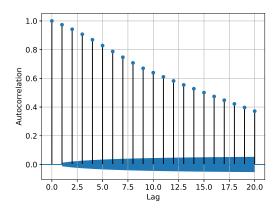


Fig. 10. The ACF of the traffic speeds of a random segment from lag 1 to lag 20. The shaded area represents the 95% confidence intervals, which is used to determine whether the autocorrelation coefficients is significantly different from zero.

where  $\Gamma$  is the set of all possible triplets. The structure of DNNs with triplet loss is shown in Fig. 9, where the outputs of the last layer are the representations used for clustering. The dimension of the representations in clustering is lower than the raw series. For example, considering a traffic series with five-minute interval during 90 days. The length of whole series is  $288 \times 90$ , while the length of daily sub-series is 288. If we use 32-dimensional representations in clustering, the ratio of reduction in dimension is about  $\frac{32}{288 \times 90} \approx 0.1\%$ . Subsequently, we average all the representations from one road segment, and cluster the representations into K groups, where K is much less than  $N_{\rm r}$ . Therefore,  $N_{\rm r}$  road segments are clustered into K groups,

$$C^{(k)} = \{x^{(r,k)}\}, 1 \le k \le K, 1 \le r \le N_{\rm r}$$
 (20)

where  $C^{(k)}$  denotes the k-th group.  $x^{(r,k)}$  represents the r-th road segment in network  $\phi$ , which is clustered into group  $C^{(k)}$ .

# B. DeepPrediction

After partitioning the network into K groups, we build a prediction model for a group in the DeepPrediction module. Some definitions and statements are first given.

Definition 2: Given two functions  $g^{(1)}: \mathbb{R} \to \mathbb{R}$  and  $g^{(2)}: \mathbb{R} \to \mathbb{R}$ , if  $\hat{g}^{(1)}$  coincides with  $g^{(2)}$  within a specified measurement range after horizontal translation,  $g^{(1)}$  is homogeneous with  $g^{(2)}$ .

Statement 1: Given two homogeneous functions  $g^{(1)}$ :  $\mathbb{R} \to \mathbb{R}$  and  $g^{(2)}$ :  $\mathbb{R} \to \mathbb{R}$ , for simplicity assuming  $g^{(1)}$  has coincided with  $g^{(2)}$ , and N distinct successive samples  $(x_i,y_i^{(1)}) \in \mathbb{R} \times \mathbb{R}$  generated from  $g^{(1)}$ . Construct a mapping between historical y values and the future y value:  $f^{(1)}:[y_1^{(1)},y_2^{(1)},\ldots,y_{N-1}^{(1)}]\to y_N^{(1)}$ . Similarly get N successive samples from  $g^{(2)}$  at same x values and construct the mapping  $f^{(2)}:[y_1^{(2)},y_2^{(2)},\ldots,y_{N-1}^{(2)}]\to y_N^2$ . It is obvious that  $f^{(2)}$  is equal to  $f^{(1)}$ .

Based on the Statement 1, we propose an idea of model sharing that all road segments within a group can share a prediction model. The implementation of the DeepPrediction is elaborated as follows:

1) **Interval confirmation.** According to the periodicity, it is intuitive to use the measurements in a period to predict the next traffic state. In order to measure the autocorrelation between current and past traffic values, we calculate the autocorrelation function (ACF) at lag *i*, which is the correlation between series values that are *i* intervals apart. As shown in Fig. 10, the measurements are linearly correlated with the contiguous measurements. The high autocorrelations imply that importing all measurements in a period will result in information redundancy. We calculate the input interval *l* by

$$l = \max_{p_i > p, i \ge 1} \{i\},\tag{21}$$

where  $p_i$  denotes the ACF at lag i, and p is the given threshold that is determined by experiments. Therefore, The input series from  $x_{t-N_i:N_i}$  becomes

$$\boldsymbol{x}_{t-N_i:N_i:l}. \tag{22}$$

The length of the input reduces from  $N_{\rm p}$  to  $N_{\rm i} = \lceil N_{\rm p}/l \rceil$  correspondingly, where  $\lceil \cdot \rceil$  represents the operation of rounded-up.

 Model sharing. Within each group, we train a model for all road segments, which is known as group-based model (GM). We generate the training samples for each group as

$$\boldsymbol{x}^{(r,k)} \to < \boldsymbol{x}_{t:N_i:l}^{(r,k)}, x_{t+N_i+N_0}^{(r,k)} >, \boldsymbol{x}^{(r,k)} \in \boldsymbol{C}^{(k)}$$
 (23)

where  $\boldsymbol{x}_{t:N_i:l}^{(r,k)}$  and  $\boldsymbol{x}_{t+N_i+N_o}^{(r,k)}$  denote the input and output of model, respectively. After that, we aggravate the samples within a group to train a GM  $f^{(k)}(\cdot)$  for group  $\boldsymbol{C}^{(k)}$ .

$$x_{t+N_i+N_o}^{(r,k)} = f^{(k)}(\boldsymbol{x}_{t:N_i:l}^{(r,k)}), \ \boldsymbol{x}^{(r,k)} \in \boldsymbol{C}^{(k)}$$
 (24)

Then the aggregated STTP model  $f(\cdot)$  at the large-scale network can be written as:

$$x_{t+N_{i}+N_{o}}^{(r,k)} = f(\boldsymbol{x}_{t:N_{i}:l}^{(r,k)})$$

$$= f^{(k)}(\boldsymbol{x}_{t:N_{i}:l}^{(r,k)}), \quad k \in \{1, 2, \dots, K\},$$
(25)

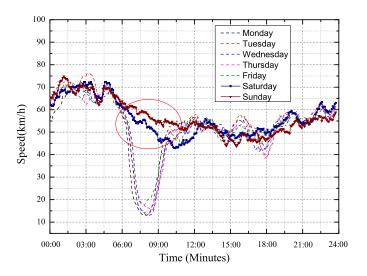


Fig. 11. The every five-minute average traffic speeds of one random segment on weekdays versus to the ones at weekends from September, 2017 to November, 2017.

#### VII. PERFORMANCE EVALUATION

In this section, we evaluate the proposed framework on the network mentioned in Section III. 27 road segments are chosen for simplicity. The network, experimental settings and performance metrics are described at first. Then, we analyze the performance over different metrics.

# A. Experiment Settings

For DeepCluster module, we split the traffic series into 90 daily sub-series of length 288 for each segment. Fig. 11 shows that the traffic patterns on weekdays are different from the ones at weekends between six and ten o'clock in the morning, since most people do not work at weekends (The circular region). Besides, the traffic patterns behave abnormally during the National Day than usual, as shown in Fig. 12. As a result, 60 daily sub-series are chosen by getting rid of the ones at weekends and during the National Day. Then we transfer the sub-series of size  $1 \times 288$  into images of size  $288 \times 288$ . As discussed in Section VI-A, we generate triplets by the daily sub-series from 27 road segments, which are used for representation learning. The deep structure of FaceNet used in this paper is the Inception\_ResNet, the configuration of which is the same with [39]. As a segment's representative, the average representations of the sub-series is used for clustering by Kmeans method. K is confirmed by Silhouette coefficient [40].

For DeepPrediction module, we use the state of the art RNNs, i.e., long short term memory (LSTM) [41] for STTP. The input span of traffic series is chosen to be a day. Then the length of the input is  $N_{\rm i} = \lceil 288/l \rceil$  that is confirmed by the experiments discussed later. We split the data into training set and testing set for each road segment, and aggregate the training set belonging to the same group to train the LSTM. In the end, K GMs are aggregated.

The key parameters of the relevant DNNs are listed in Table VII-A. If not mentioned specifically, all models are trained on eighty percent of data while tested on the remaining

TABLE I
THE CONFIGURATIONS OF THE RELEVANT DNNS.

Module	Network	Parameter	Size	
		Image size	160	
DeepCluster		Batch size	12	
	<sup>a</sup> FaceNet	Segments per batch	6	
		Images per segment	9	
		Embedding size	32	
DeepPrediction	<sup>b</sup> LSTM	Time steps	$\lceil 288/s \rceil$	
		LSTM1	$[1 \times 50]$	
		LSTM2	$[50 \times 25]$	
		Dense1	$[25 \times 200]$	
		Dense2	$[200 \times 1]$	

<sup>&</sup>lt;sup>a</sup>The implications of the parameters given in this table are explained exactly in [39].

data. 10-fold cross-validation is adopted over training dataset. The K-means method is implemented using the Scikit-learn Python 3.6.5. The NNs are conducted with a NVIDIA p2000 GPU, TensorFlow r1.8, CUDA 9.0 and CuDNN 9.0. Moreover, four performance metrics includes relative error (RE), mean relative error (MRE), max mean relative error (MARE) and minimum mean relative error (MIRE) are used for evaluation, which are defined as

$$e_{\text{RE}}^{(r,k)} = \frac{\left| x_{t+N_o}^{(r,k)} - \hat{x}_{t+N_o}^{(r,k)} \right|}{x_{t+N}^{(r,k)}}, 1 \le k \le K$$
 (26)

$$e_{\text{MRE}}^{(k)} = \frac{1}{|C^{(k)}|} \sum_{\boldsymbol{x}^{(r,k)} \in C^{(k)}} e_{\text{RE}}^{(r,k)}, 1 \le k \le K$$
 (27)

$$e_{\text{MARE}}^{(k)} = \max_{\boldsymbol{x}^{(r,k)} \in \boldsymbol{C}^{(k)}} \{e_{\text{RE}}^{(r,k)}\}, 1 \le k \le K$$
 (28)

$$e_{\text{MIRE}}^{(k)} = \min_{\boldsymbol{x}^{(r,k)} \in \boldsymbol{C}^{(k)}} \{e_{\text{RE}}^{(r,k)}\}, 1 \le k \le K \tag{29}$$

where  $e_{\mathrm{RE}}^{(r,k)}$  denotes the RE of r-th segment in network clustered into group  $C^{(k)}$  with  $x_{t+N^o}^{(r,k)}$  being the true speed and  $\hat{x}_{t+N^o}^{(r,k)}$  being the prediction.  $|C^{(k)}|$  is the number of road segments in the group k. Besides,  $e_{\mathrm{MRE}}^{(k)}$ ,  $e_{\mathrm{MARE}}^{(k)}$  and  $e_{\mathrm{MIRE}}^{(k)}$  are MRE, MARE and MIRE of group k, respectively. The performance metrics for road network can be similarly calculated.

# B. Simulation Results

Three experiments are conducted, including road segments clustering, interval confirmation and STTP at network.

1) **Road segments clustering.** All 27 road segments are clustered into 3 groups by DeepCluster as shown from

TABLE II The Performance under Different Input Intervals

Input interval l	MRE of Training(%)	MRE of Testing(%)		
1	3.70	5.77		
3	4.18	5.25		
5	4.37	5.48		
7	5.83	8.77		

Fig. 13a to Fig. 13c. It can be found that the series in a group are in general *homogeneous* with the other series defined at Section VI-A, which demonstrates the proposed DeepCluster's capacity of extracting the shape-based features. For example, the segments in cluster 1 have a breakdown in traffic speed during the evening peak period, followed by speed recovery. The cluster 2 have a breakdown during the morning peak, and start to swing at the middle speed back-and-forth. The segments in cluster 3 have some slight resemblances to cluster 1 during the evening peak period. However there is a stable condition holding the middle speed after six o 'clock in the morning.

- 2) Interval confirmation. This part investigates the effect of input interval on predictive performance and determines the threshold p of the ACF defined in Section VI-B. The LSTM is performed to predict the next five-minute speed under different input intervals l over the 3 random segments. From the performance listed in Table II, the MRE of training increases with the decrease of l. However, the performance improvements are insignificant when  $l \leq 5$ , such as the training MRE at 3.7% and 4.4% when l=1 and l=5. Besides, the testing MRE at l=1 is slightly larger than that at l=5. This is because the capacity of model becomes stronger as input interval decreases, leading to overfitting. From this result, the threshold is empirically set to 0.8. In the end, the input interval is set to 5 corresponding to twenty-five minutes for all other simulations.
- 3) STTP at network. For the performance comparison, we construct an IM for a segment by the same configuration of LSTM under different prediction horizon  $N_o$ . Simulation results are listed in Table III. The IMs have lower training MRE than the GMs due to the fact that the capacity of the IMs is highly stronger than that of the GMs. However, the GMs can get lower gaps between training MRE and testing MRE in all tests, since increasing the number and diversity of the training samples can improve generalization capability of the model. On the contrary, the IMs are constrained by the problem of overfitting resulted from modeling the noise. As shown in Fig. 14, the gaps of GMs are close to 0 while the gaps of IMs are around 2%.

From Table III, we can observe that the GMs perform better than IMs in terms of testing error in a relatively simple task of five-minute forecasting. The testing MRE of the GMs and IMs are 4.12% and 5.05% for group 1, 4.07% and 4.94% for group 2, 5.00% and 5.37%

<sup>&</sup>lt;sup>b</sup>The inputs and outputs size are described in  $[rows \times cols]$ .

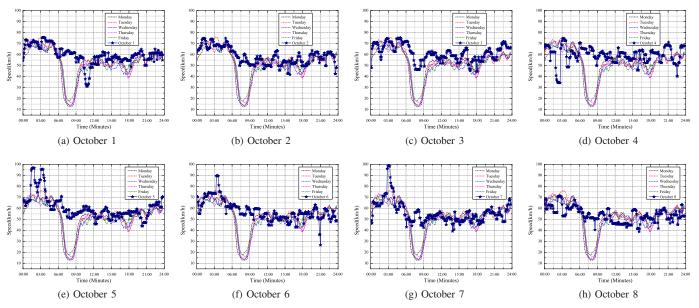


Fig. 12. The every five-minute average traffic speeds of a random segment on weekdays from September, 2017 to November, 2017 versus the ones during the National Day.

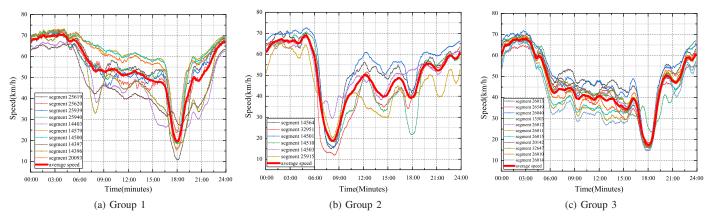


Fig. 13. The every five-minute average traffic speeds of the road segments on weekdays from September, 2017 to November, 2017 in different groups. The thicker red lines represent the centers of the corresponding clusters.

for group 3, respectively. However, as the task becomes complex, the capacity of GMs become insufficient. For example, the testing MRE of GMs are around 1% more than that of IMs when  $N_{\rm o}=2$ , while the testing MRE of GMs are around 2% more than that of IMs when  $N_{\rm o}=3$ .

As shown in Fig. 15, the GM can predict the trends of traffic speed well, but the performance gets worse with the increase of the prediction horizon. It also shows that the model does not work well of 10 and 15 minutes forecasting in rush hours (The dash area in Fig. 15), that the traffic speed switching sharply.

The proposed framework is scalable that can be applied for the large-scale networks easily by reducing the number of models significantly, and can reach the compromise of the number of models and prediction performance. Compared to the traditional 27 IMs, the number of prediction models has been reduced up to

 $\frac{(27-3)}{27} \approx 88\%$  with about 0.7% - 1.9% performance degradation, in terms of network MRE in our test, as shown in Fig. 14. In conclusion, the performance of the framework is comparable to that of customized IMs, which validates the ability for STTP at large-scale networks.

# VIII. CONCLUSION

The characteristics of the multiplicity and heterogeneity make STTP at large-scale network a challenging and important problem. By exploiting the characteristic of traffic patterns, a DL framework for STTP at large-scale networks is proposed in this paper. The key point of the framework is the combination of the DeepCluster and the DeepPrediction, as well as the model sharing strategy. We analytically evaluate the proposed framework over a real large-scale network of Liuli Bridge in Beijing and some insights into generic DL models are obtained. Despite that the prediction performances of the GMs

TABLE III
THE GROUP PERFORMANCE OF THE PROPOSED FRAMEWORK

Prediction Horizon	Group	Algorithm	MRE of Training(%)	MRE of Testing(%)	<b>Gap</b> (%)	MARE of Testing(%)	MIRE of Testing(%)
1(five-minute)	1	GM	3.97	4.12	0.1	6.87	2.65
		IM	3.20	5.05	1.9	9.80	3.03
	2	GM	4.08	4.07	0	5.76	3.39
		IM	3.59	4.94	1.3	6.04	3.67
	3	GM	4.96	5.00	0	6.54	4.39
		IM	3.72	5.37	1.6	6.86	4.40
2(Ten-minute)	1	GM	5.92	6.04	0.1	10.07	3.70
		IM	3.80	5.77	2.0	9.94	3.57
	2	GM	6.22	6.22	0	9.86	5.29
		IM	3.90	5.67	1.8	6.70	4.84
	3	GM	7.16	7.24	0	9.56	6.27
		IM	4.20	6.17	2.0	7.73	4.91
<b>3</b> (Fifteen-minute)	1	GM	7.08	7.35	0.3	11.82	4.61
		IM	4.01	5.82	1.8	9.21	3.97
	2	GM	7.71	7.93	0.2	11.54	6.68
		IM	4.12	5.66	1.5	7.00	4.69
	3	GM	8.36	8.43	0.1	12.00	7.07
		IM	4.71	6.38	1.7	8.49	4.87

GM: Group-based Model. IM: Individual-based Model.

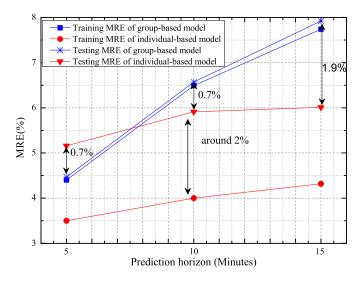


Fig. 14. The network MRE of GM and IM.

are slightly worse than that of IMs in most tests, the GMs have a better generalization ability. For five-minute prediction, the GM gets 0.7% error lower than IM. We also discuss the effect of input interval on the prediction performance, which guides the framework on how to select the effective input interval. Furthermore, we use only 3 models to achieve the STTP at network, while the traditional way needs to construct 27 models.

## IX. ACKNOWLEDGEMENT

This work is funded in part by the National Natural Science Foundation of China under Grant 61731004.

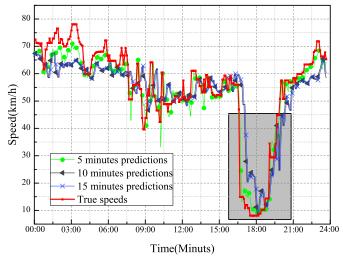


Fig. 15. The every five-minute predictions using GM model versus true speeds of a random road segment in group 1 on November 31, 2017.

## REFERENCES

- [1] M. Wang, H. Shan, R. Lu, R. Zhang, X. Shen, and F. Bai, "Real-time path planning based on hybrid-VANET-enhanced transportation system," *IEEE Transactions on Vehicular Technology*, vol. 64, no. 5, pp. 1664-1678, 2015.
- [2] S. Sun, and C. Zhang, "The selective random subspace predictor for traffic flow forecasting," *IEEE Transactions on Intelligent Transportation Systems*, vol. 8, no. 2, pp. 367-373, June 2007.
- [3] E. I. Vlahogianni, M. G. Karlaftis, and J. C. Golias, "Short-term traffic forecasting: Where we are and where we' re going," *Transportation Research Part C: Emerging Technologies*, vol. 43, pp. 3-19, June 2014.
- [4] A. Ermagun, D. Levinson, "Spatiotemporal traffic forecasting: review and proposed directions," *Transport Reviews*, vol. 38, no. 6, pp. 786-814, February 2018.
- [5] I. Lana, J. D. Ser, M. Velez, and E. Vlahogianni, "Road Traffic

- Forecasting: Recent Advances and New Challenges," *IEEE Intelligent Transportation Systems Magazine*, vol. 10, no. 2, pp. 93-109, April 2018.
- [6] W. Xu, H. Zhou, N. Cheng, F. Lyu, W. Shi, J. Chen, and X. Shen, "Internet of vehicles in big data era" *IEEE/CAA Journal of Automatica Sinica*, vol. 5, no. 1, pp. 19-35, 2018.
- [7] K. Hornik, "Approximation capabilities of multilayer feedforward networks," *Neural Networks*, vol. 4, no. 2, pp. 251-257, March 1991.
- [8] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436-444, May 2015.
- [9] A. Bagnall, E. Keogh, S. Lonardi S, and G. Janacek, "A bit level representation for time series data mining with shape based similarity," *Data Mining and Knowledge Discovery*, vol. 13, no. 1, pp. 11-40, May 2006.
- [10] X. Wang, A. Mueen, H. Ding, G. Trajcevski, P. Scheuermann, and E. Keogh, "Experimental comparison of representation methods and distance measures for time series data," *Data Mining and Knowledge Discover*, vol. 26, no. 2, pp. 275-309, March 2013.
- [11] S. Aghabozorgi, A. S. Shirkhorshidi, and T. Y. Wah, "Time-series clustering-A decade review," *Information Systems*, vol. 53, pp. 16-38, October 2015.
- [12] R. Agrawal, C. Faloutsos, and A. Swami, "Efficient similarity search in sequence databases," in *Proceedings of 4th International Conference* on Foundations of Data Organization and Algorithms, Heidelberg, June 1993, pp. 69-84.
- [13] F. Korn, H. V. Jagadish, and C. Faloutsos, "Efficiently supporting ad hoc queries in large datasets of time sequences," in *Proceedings of ACM SIGMOD International Conference on Management of Data*, Tucson, Arizona, USA, May 1997, pp. 289-300.
- [14] E. J. Keogh, and M. J. Pazzani, "An Enhanced Representation of Time Series Which Allows Fast and Accurate Classification, Clustering and Relevance Feedback," in *Proceedings of 4th International Conferencee* on Knowledge Discovery and Data Mining (KDD), NY, August 1998, pp. 239-247.
- [15] E. Keogh, K. Chakrabarti, M. Sharad, and M. Pazzani, "Locally adaptive dimensionality reduction for indexing large time series databases," in Proceedings of ACM SIGMOD International Conference on Management of data, Santa Barbara, California, USA, May 2001, pp. 151-162.
- [16] E. Keogh, S. Lonardi, and C. A. Ratanamahatana, "Towards parameter-free data mining," in *Proceedings of 10th ACM SIGMOD International Conference on Knowledge Discovery and Data Mining*, Seattle, WA, USA, August 2004, pp. 206-215.
- [17] B. K. Yi, and C. Faloutsos, "Fast time sequence indexing for arbitrary Lp norms," in *Proceedings of 26th International Conference on Very Large Data Bases (VLDB)*, Cairo, Egypt, September 2000, pp. 385-394.
- [18] F. L. Chung, T. C. Fu, R. W. P. Luk, and V. T. Y. Ng, "Flexible time series pattern matching based on perceptually important points," in Proceedings of International Joint Conference on Artificial Intelligence Workshop (Learning From Temporal and Spatial Data), Seattle, WA, August 2001, pp. 1-7.
- [19] C. Song, F. Liu, Y. Huang, L. Wang, and T. Tan, "Auto-encoder based data clustering," in *Iberoamerican Congress on Pattern Recognition*, Springer, 2013, pp. 117-124.
- [20] F. Tian, B. Gao, Q. Cui, E. Chen, and T. Liu, "Learning deep representations for graph clustering," in *Proceedings of 28th AAAI Conference on Artificial Intelligence*, 2014, pp. 1293-1299.
- [21] J. Xu, P. Wang, G. Tian, B. Xu, J. Zhao, and F. Wang, "Short Text Clustering via Convolutional Neural Networks," in *Proceedings* of Conference on North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT), Denver, Colorado, May 2015, pp. 62-69.
- [22] J. Xie, R. Girshick, and A. Farhadi, "Unsupervised deep embedding for clustering analysis," in *Proceedings of the 33rd International Conference* on *Machine Learning*, San Juan, PR, USA, May 2016, pp. 478-487.
- [23] K. Tian, S. Zhou, and J. Guan, "DeepCluster: A General Clustering Framework Based on Deep Learning," in *Proceedings of Joint Euro*pean Conference on Machine Learning and Knowledge Discovery in Databases, Springer, Cham, December 2017, pp. 809-825.
- [24] S. Sun, R. Huang, and Y. Gao, "Network-scale traffic modeling and forecasting with graphical lasso and neural networks," *Journal of Trans*portation Engineering, vol. 138, no. 11, pp. 1358-1367, 2012.
- [25] B. Yu, X. L. Song, F. Guan, Z. M. Yang, and B. Z. Yao, "k-Nearest neighbor model for multiple-time-step prediction of short-term traffic condition," *Journal of Transportation Engineering*, vol. 142, no. 6, pp. 1-10, 2016.
- [26] G. N. Polson, O. V. Sokolov, "Deep learning for short-term traffic flow prediction," *Transportation Research Part C: Emerging Technologies*, vol. 79, pp. 1-17, June 2017.

- [27] Y. Li, R. Yu, C. Shahabi, and Y. Liu, "Diffusion convolutional recurrent neural network: Data-driven traffic forecasting," in *Proceedings of In*ternational Conference of Learning Representation (ICLR), Vancouver, Canada, April 2018.
- [28] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, Boston, MA, USA, June 2015, pp. 815-823.
- [29] M. Ben-Akiva, M. Bierlaire, H. Koutsopoulos, and R. Mishalani, "Dyna-MIT: a simulation-based system for traffic prediction," in *Proceedings of DACCORD Short Term Forecasting Workshop*, Delft, Netherlands, Feb. 1998
- [30] A. Abadi, T. Rajabioun, and P. A. Ioannou, "Traffic flow prediction for road transportation networks with limited traffic data," *IEEE Transac*tions on Intelligent Transportation Systems, vol. 16, no. 2, pp. 653-662, April 2015.
- [31] T. Djukic, J. W. C. Van Lint, and S. P. Hoogendoorn, "Application of principal component analysis to predict dynamic origindestination matrices," *Transportation Research Record: Journal of the Transportation Research Board*, vol. 2283, no. 1, pp. 81-89, January 2012
- [32] N. Mitrovic, M. T. Asif, U. Rasheed, J. Dauwels, and P. Jaillet, "CUR decomposition for compression and compressed sensing of large-scale traffic data," in *Proceedings of 16th International Conference on Intelligent Transportation Systems (ITSC)*, Hague, Netherlands, October 2013, pp. 1475-1480.
- [33] M. T. Asif, S. Kannan, J. Dauwels, and P. Jaillet, "Data compression techniques for urban traffic data," in *Proceedings of IEEE Symposium* on Computational Intelligence in Vehicles and Transportation Systems (CIVTS), Singapore, April 2013, pp. 44-49.
- [34] N. Mitrovic, M. T. Asif, J. Dauwels, and P. Jaillet, "Low-Dimensional Models for Compressed Sensing and Prediction of Large-Scale Traffic Data," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 5, pp. 2949-2954, April 2015.
- [35] W. Min, and L. Wynter, "Real-time road traffic prediction with spatiotemporal correlations," *Transportation Research Part C: Emerging Tech*nologies, vol. 19, no. 4, pp. 606616, 2011.
- [36] M. T. Asif, J. Dauwels, C. Y. Goh, A. O. E. Fathi, M. Xu, M. M. Dhanya, N. Mitrovic, and P. Jaillet, "Spatiotemporal patterns in large-scale traffic speed prediction," *IEEE Trans. Intell. Transp. Syst.*, vol. 15, no. 2, pp. 794-804, April 2014.
- [37] X. Ma, Z. Dai, Z. He, J. Ma, Y. Wang, and Y. P. Wang, "Learning Traffic as Images: A Deep Convolutional Neural Network for Large-Scale Transportation Network Speed Prediction," *Sensors*, vol. 17, no. 4. pp. 1-16, April 2017.
- [38] K. Xie, L. Wang, X. Wang, G. Xie, J. Wen, G, Zhang, J. Cao, and D. Fang, "Accurate Recovery of Internet Traffic Data: A Sequential Tensor Completion Approach," *IEEE/ACM Transactions on Networking (TON)*, vol. 26, no.2, pp. 793-806, April 2018.
- [39] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of Conference on Computer Vision and Pattern Recognition (CVPR)*, Boston, MA, USA, June 2015, pp. 815-823.
- [40] P. J. Rousseeuw, "Silhouettes: a graphical aid to the interpretation and validation of cluster analysis," *Journal of Computational and Applied Mathematics*, vol. 20, pp. 53-65, November 1987.
- [41] S. Hochreiter, and J. Schmidhuber, "Long short-term memory," Neural Computation, vol. 9, no. 8, pp. 17351780, November 1997.

Lingyi Han received the B.S. degree from Beijing University of Posts and Telecommunications, China, in 2016. She is currently pursuing the Ph.D. degree with the Intelligent Computing and Communication Laboratory, Key Laboratory of Universal Wireless Communications, Ministry of Education, Beijing University of Posts and Telecommunications. His research interests include data mining and artificial intelligence in Internet-of-Things.

Kan Zheng received the B.S., M.S., and Ph.D. degrees from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 1996, 2000, and 2005, respectively. He is currently a Professor at BUPT. He is the author of more than 200 journal articles and conference papers in the field of resource optimization in wireless networks, M2M networks, VANET, and so on. He has rich industry experiences on the standardization of the new emerging technologies. Dr. Zheng holds editorial board positions for several journals. He has organized several special issues in famous journals, including the IEEE COMMUNICATIONS ON SURVEYS AND TUTORIALS and TRANSACTIONS ON EMERGING TELECOMMUNICATIONS TECHNOLOGIES.

Long Zhao received the B.S. degree from Nanjing University of Posts and Telecommunications, Nanjing, China, in 2008 and the M.S. degree from Harbin Engineering University, Harbin, China, in 2011. He is currently working toward the Ph.D. degree with the Beijing University of Posts and Telecommunications, Beijing. From April 2014 to March 2015, he was a Visiting Scholar at the Department of Electrical Engineering, Columbia University, supervised by Prof. X. Wang. His research interests include wireless communications and signal processing.

Xianbin Wang (S'98-M'99-SM'06-F'17) received the Ph.D. degree in electrical and computer engineering from the National University of Singapore in 2001.

From January 2001 to July 2002, he was a System Designer with STMicroelectronics, where he was responsible for the system design of DSL and Gigabit Ethernet chipsets. He was with the Communications Research Centre Canada (CRC) as a Research Scientist/Senior Research Scientist from 2002 to 2007. He is currently a Professor and the Canada Research Chair with Western University, Canada. His current research interests include 5G technologies, Internet-of-Things, communications security, and locationing technologies. He has over 300 peer-reviewed journal and conference papers, in addition to 26 granted and pending patents, and several standard contributions.

Dr. Wang is an IEEE Distinguished Lecturer. He received many awards and recognition, including the Canada Research Chair, the CRC Presidents Excellence Award, the Canadian Federal Government Public Service Award, the Ontario Early Researcher Award, and five IEEE Best Paper Awards. He was involved in a number of IEEE conferences, including GLOBECOM, ICC, VTC, PIMRC, WCNC, and CWIT, in different roles, such as the Symposium Chair, a Tutorial Instructor, the Track Chair, the Session Chair, and the TPC Co-Chair. He currently serves as an Editor/Associate Editor of the IEEE TRANSACTIONS ON COMMUNICATIONS, the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY. He was also an Associate Editor of the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS from 2007 to 2011, and the IEEE WIRELESS COMMUNICATION LETTERS from 2011 to 2016.

**Xuemin** (Sherman) Shen (M'97-SM'02-F'09) received the B.Sc. degree from Dalian Maritime University, Dalian, China, in 1982 and the M.Sc. and Ph.D. degrees from Rutgers University, Piscataway, NJ, USA, in 1987 and 1990, all in electrical engineering.

From 2004 to 2008, he was the Associate Chair for Graduate Studies with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON, Canada. He is currently a Professor and the University Research Chair with the Department of Electrical and Computer Engineering, University of Waterloo. He is a coauthor or editor of six books and the author of several papers and book chapters in wireless communications and networks, control, and filtering. His research interests include resource management in interconnected wireless/ wired networks, wireless network security, wireless body area networks, and vehicular ad hoc and sensor networks.

Dr. Shen served as the Technical Program Committee Chair for the 2010 Fall IEEE Vehicular Technology Conference (IEEE VTC10 Fall); the Symposia Chair for the 2010 IEEE International Conference on Communications (IEEE ICC10); the Tutorial Chair for IEEE VTC11 Spring and IEEE ICC08; the Technical Program Committee Chair for the 2007 IEEE Global Communications Conference; the General Co-Chair for the 2007 IEEE International Conference on Communications and Networking in China and the 2006 Third International Conference on Quality of Service in Heterogeneous Wired/ Wireless Networks; and the Chair for IEEE Communications Society Technical Committee on Wireless Communications and Peer-to-Peer Communications and Networking. He also serves/served as the Editor-in-Chief for IEEE NETWORK, Peer-to-Peer Networking and Application, and IET Communications; as a Founding Area Editor for IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS; as an Associate Editor for IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, Computer Networks, and ACM Wireless Networks; and as the Guest Editor for IEEE JOUR-NAL ON SELECTED AREAS IN COMMUNICATIONS, IEEE WIRELESS COMMUNICATIONS, IEEE COMMUNICATIONS MAGAZINE, and ACM Mobile Networks and Applications. He is a registered Professional Engineer of Ontario, Canada; a Fellow of the Canadian Academy of Engineering and the Engineering Institute of Canada; and a Distinguished Lecturer of the IEEE Vehicular Technology and Communications Societies.