

# Short-term traffic flow forecasting with spatial-temporal correlation in a hybrid deep learning framework

Wu Yuankai<sup>\*1</sup> and Tan Huachun<sup>2</sup>

<sup>1</sup>School of Mechanical Engineering, Beijing Institute of Technology

<sup>2</sup>School of Mechanical Engineering, Beijing Institute of Technology

December 6, 2016

## Abstract

Deep learning approaches have reached a celebrity status in artificial intelligence field, its success have mostly relied on Convolutional Networks (CNN) and Recurrent Networks. By exploiting fundamental spatial properties of images and videos, the CNN always achieves dominant performance on visual tasks. And the Recurrent Networks (RNN) especially long short-term memory methods (LSTM) can successfully characterize the temporal correlation, thus exhibits superior capability for time series tasks. Traffic flow data have plentiful characteristics on both time and space domain. However, applications of CNN and LSTM approaches on traffic flow are limited. In this paper, we propose a novel deep architecture combined CNN and LSTM to forecast future traffic flow (CLTFP). An 1-dimension CNN is exploited to capture spatial features of traffic flow, and two LSTMs are utilized to mine the short-term variability and periodicities of traffic flow. Given those meaningful features, the feature-level fusion is performed to achieve short-term traffic flow forecasting. The proposed CLTFP is compared with other popular forecasting methods on an open datasets. Experimental results indicate that the CLTFP has considerable advantages in traffic flow forecasting. in additional, the proposed CLTFP is analyzed from the view of Granger Causality, and several interesting properties of traffic flow and CLTFP are discovered and discussed .

**Traffic flow forecasting, Convolutional neural network, Long short-term memory, Feature-level fusion**

## 1 Introduction

The accurate and reliable forecasting of short-term traffic flow is the precursor of a multitude of intelligent transportation systems (ITS) such as proactive dynamic traffic control, intelligent route guidance and intelligent location-based service, thus it is always a hot topic in ITS. Recent

---

<sup>\*</sup>Corresponding Author: 5433809@bit.edu.cn

developments in information collection and transmission have introduced the notion of big data in the field of ITS (Zheng et al., 2016), which has directed many researches toward data-driven forecasting approaches.

What are frequently used in data-driven forecasting are the two different approaches which are statistics and neural networks (Karlaftis and Vlahogianni, 2011). The statistics such as autoregressive integrated moving average (ARIMA) (Min and Wynter, 2011), Markov chain (Qi and Ishak, 2014) and Bayesian network (Wang, Deng, and Guo, 2014), generally focus on finding the spatial-temporal pattern from a probabilistic perspective and then use that predictive information for forecasting. They can provide some insights on the probabilistic mechanisms generating the traffic data and capture the uncertainty within traffic flow. However, the statistics frequently fail when dealing with nonlinearity within traffic flow because the linear architecture is often used, and they always suffer from curse of dimensionality, which is a common phenomenon in big data era.

Compared with those classical statistical models, the neural networks have several advantages. First, neural networks use tens of thousands of neuron activities to simulate the unknown relationship, and thus they are non-parametric approaches, which are more flexible with input variables. Second, the neural networks are more capable of handling nonlinearity with the help of nonlinear activation functions. Because of those advantages of neural networks and complexity of traffic flow, neural networks are widely used in short-term traffic flow forecasting (Chan et al., 2012).

Recently, neural networks with deep architectures have proven to be very successful in image, video, audio and language learning tasks (LeCun, Bengio, and Hinton, 2015). In short-term traffic forecasting area, though traditionally shallow neural networks are generally adopted, the deep neural networks have also aroused enormous researches' interests. Deep multi-layer fully connected networks are frequently employed in current short-term traffic forecasting (Huang et al., 2014; Lv et al., 2015), and pre-training strategies with unsupervised learning algorithms such as Restricted Boltzmann machine (RBM) (Hinton, Osindero, and Teh, 2006) and Stacked AutoEncoder (SAE) (Vincent et al., 2008) are often used. Though pre-training strategies significantly promoted the performance of fully-connected networks (Hinton and Salakhutdinov, 2006), as each neuron in fully-connected layer is connected to every neuron in the previous layer, which makes that fully-connected networks are expensive in terms of memory and computation. Moreover, there is no assumptions about the features in the fully-connected architecture, thus it is difficult for a fully-connected neural networks to capture representative features from data with plentiful characteristics.

Also, like frequently studied data in machine learning area such as video and audio, traffic flow data have plentiful characteristics in both time and space domain (Tan et al., 2016a, 2013). There are some obvious characteristics, for example, in space domain, traffic flow patterns on some location are more likely to have stronger dependencies on nearby locations (topological locality); in time domain, the traffic flow several weeks/days before even has a long-term impact on current traffic flow because of people's traveling habits (long-term memory). A representative characterization of those spatial-temporal features is the key to successful traffic flow forecasting. In recent years, one of the most successful deep neural networks to model topological locality is the convolutional

neural network (CNN) (Krizhevsky, Sutskever, and Hinton, 2012), it uses filters to find relationships between neighboring inputs, which can make it easier for the network to converge on the correct solution. And one of the most successful architecture to characterize long-term memory is long short-term memory network (LSTM) (Graves, Jaitly, and Mohamed, 2013), which learns both short-term and long-term memory by enforcing constant error flow through the designed cell state (Hochreiter and Schmidhuber, 1997).

Motivated by the success of CNN and LSTM, and with consideration of the spatial-temporal characteristics of traffic flow, we propose a novel short-term traffic flow prediction method based on the combination of CNN and LSTM (CLTFP). A deep convolution neural network is utilized to mine the space features of traffic flow data. LSTMs are employed to learn features of both short-term time variation and long-term periodicity. Then we feed those spatial-temporal features into a linear regression layer to predict future traffic flow. In feature-level based data fusion, it is natural to assume that a small portion of features have strong impact (Bishop, 2006). Thus in order to strengthen the sparsity of features, we add a  $l_1$  norm constraint on weights of the regression layer. Finally, we train the neural network end-to-end. Our method is evaluated on traffic flow of a freeway corridor collected from an open datasets, the results show that our method exhibits better performance than state-of-arts. Moreover, we analyze the features captured by CLTFP in terms of incremental predictability, the analysis graphically demonstrate how black-box typed CLTFP understand causality between future-past traffic flow.

## 2 Related works

The deep architecture for short-term traffic flow forecasting has recently been studied in ITS, but mainly on deep fully-connected architecture. Huang et al. (2014) used a deep belief network to capture the spatial-temporal features of traffic flow and proposed a multi-task learning architecture to perform exit station flow and road flow forecasting. Similarly, Lv et al. (2015) proposed a stacked autoencoder model based short-term traffic flow forecasting. Tan et al. (2016b) investigated the impact of pre-traning with different deep belief networks on the DNN based short-term traffic flow forecasting. Chen et al. (2016) developed a Stack denoise Autoencoder to learn hierarchical representation of urban traffic flow. Polson and Sokolov (2016) used a deep neural network architecture to forecast traffic flows during special events. The above approaches have been able to accurately forecast future traffic flow to some degree, but they did not exploit the topological locality and long-term memory of traffic flow, which hindered their predictive power. Motivated by the success of LSTM, Ma et al. (2015) applied LSTM to short-term traffic forecasting, they claimed that LSTM can capture long-term memory of traffic flow. However, despite the efficient usage of long temporal dependency, the spatial dependency is not fully utilized in their work.

Several attempts have been made to combine CNN and LSTM architectures especially for connection of computer vision and natural language processing. For example, several methods have made use of CNN features and LSTM for image/video description generation (Vinyals et al., 2015;

(Yao et al., 2015; Peris et al., 2016). The combination of CNN and LSTM have also been successfully applied to visual activity recognition (Donahue et al., 2015), sentiment analysis (Wang et al., 2016) and video classification (Wu et al., 2015). As traffic flow data share some common properties with visual data (Yang, Kalpakis, and Biem, 2014) and language, the success of combining CNN and LSTM on computer vision and language processing indicates the potential of such combination on traffic flow forecasting.

### 3 Model

In this section, we describe our CNN-LSTM based short-term traffic flow prediction method (CLTFP). See Fig. 1 for the graphical illustration of the proposed model. It can be found from Fig.1 that CLTFP consists of a 1D CNN, two LSTM RNNs and a fully-connected layer, we will describe each part in the following subsection.

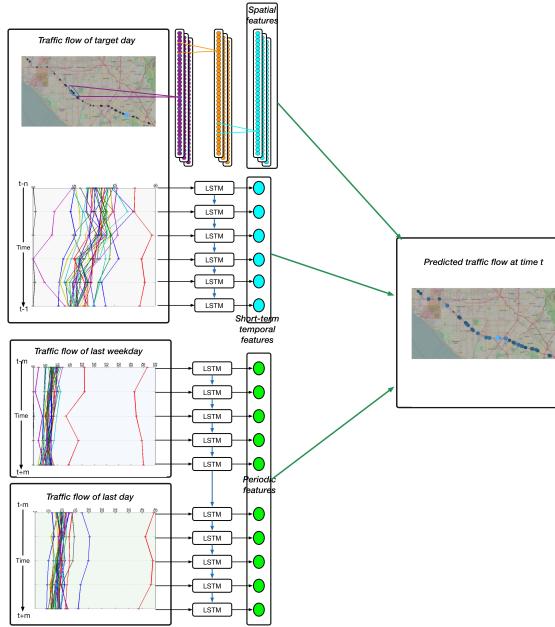


Figure 1: CLTFP, our model, is based end-to-end on a neural network consisting of a 1D CNN (capture spatial features), two LSTM RNNs (capture short-term and periodic features) and followed by a fully connected layer to fuse all features to forecast traffic flow at target time point  $t$ .

#### 3.1 Spatial features captured by CNN

Suppose we need to forecast traffic flow of  $p$  locations  $\{s_i\}_{i=1}^p$  in  $(t, t+1, \dots, t+h)$ , in which  $h$  is the prediction horizon. The historical traffic flow data of  $p$  locations  $\{s_i\}_{i=1}^p$  in  $(t-n, t-1)$  are traditionally used as inputs for generating prediction in  $(t, t+1, \dots, t+h)$ . If we put the historical

data together, we can get a data matrix

$$\mathbf{S} = \begin{bmatrix} S_1 \\ S_2 \\ \vdots \\ \vdots \\ S_p \end{bmatrix} = \begin{bmatrix} s_1(t-n) & s_1(t-n+1) & \cdots & s_1(t-1) \\ s_2(t-n) & s_2(t-n+1) & \cdots & s_2(t-1) \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ s_p(t-n) & s_p(t-n+1) & \cdots & s_p(t-1) \end{bmatrix}. \quad (1)$$

The traffic flow usually depends on traffic flow of that location and its neighbors. As Convolutional Neural Networks (CNNs) have been successful in handling data representation with a locality structure, we naturally adapt a 1-dimensional CNN to capture spatial features of traffic flow. Our 1D CNN does not attack time mode. Instead, the time dimension is treated as channels of an image, which means that we only perform convolution on vectors  $T_q = [s_1(t-n+q), s_2(t-n+q), \dots, s_p(t-n+q)]^T$  ( $0 \leq q \leq n-1$ ) of matrix  $\mathbf{S}$  in Eq.(1).

For locations of a freeway corridor given in Fig.1, the traffic travel through from upstream  $s_1$  to downstream  $s_p$ , the conventional 1D CNN is naturally exploited to capture spatial features of such a transportation network, where the  $k$ -th feature map is obtained as follows

$$h_q^k = o_c(w_q^k * T_q^k + b_q^k), \quad (2)$$

where  $w_q^k$  is the weights vector,  $b_q^k$  is the bias,  $o_c$  denotes a nonlinear activation and  $*$  denotes the convolution. The pooling layers are not employed in our model. Because it is evident that an all convolution net achieves better performance on small images recognition (Springenberg et al., 2014), and the space dimension of traffic data on short-term traffic forecasting task is always limited.

There are many more complex transportation networks than the one given in Fig.1, e.g. a transportation network of a big city. In these cases, the conventional 1D CNN can not be used without any modifications. However, the traffic flow in any transportation network always has some graph structures (Shahsavari and Abbeel, 2015), Thus the CNN on graph-structured data proposed by Henaff, Bruna, and LeCun (2015) can be an alternative for such transportation networks.

### 3.2 Short-term temporal features

As stated in Ma et al. (2015), the traditional forecasting models mainly suffer from two drawbacks in handling time mode information of traffic flow: (1) Traditional methods especially traditional RNNs are difficult to train if the traffic flow series has long time lags, which means traditional RNNs will provide poor performance if the time window size  $n$  in Eq.(1) is too large. (2) It is difficult to find the optimal time window size  $n$ , as the correlation between traffic flow of different time points is affected by many complex factors such as weather, speed and unpredictable incident. The LSTM is one of the more practical ways to tackle these issues, thus we propose to use LSTM to capture the

time mode information of traffic flow. Different from the work of Ma et al. (2015), which directly used LSTM to generate predictions of traffic flow, we use LSTM to generate time features at each time point and build a deeper and more complex traffic flow forecasting model.

Similar to the traditional RNNs, a LSTM structure is composed of one input layer, one or several hidden layers and one output layer. The core idea of LSTM is the memory cell in hidden layers, which is designed to avoid the gradient vanish and explosion in traditional RNNs. As shown in Fig. 2, a memory cell contains four main parts: an input gate, a neuron with a self-recurrent connection, a forget gate and an output gate.

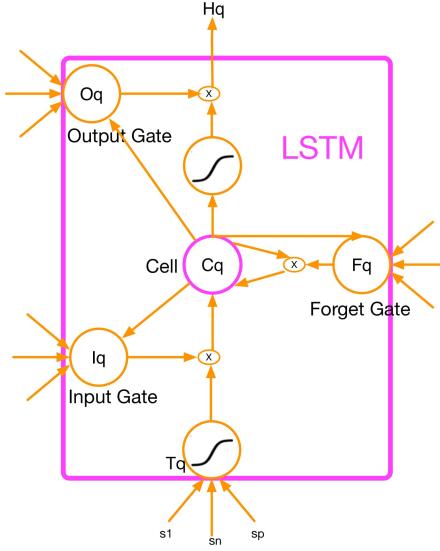


Figure 2: LSTM structure

For the generation of short-term temporal features, the inputs of LSTM is denoted as  $T = (T_0, T_1, \dots, T_{n-1})$  where  $T_q = [s_1(t-n+q), s_2(t-n+q), \dots, s_p(t-n+q)]^T$  in Eq (1), and the output temporal features in each historical time point is denoted as  $H = (H_0, H_1, \dots, H_{n-1})$ ,  $n$  is the time window size. The generated temporal features are iteratively calculated by following equations:

$$I_q = \sigma(W_i T_q + U_i H_{q-1} + W_{ci} \cdot C_{q-1} + b_i), \quad (3)$$

$$F_q = \sigma(W_f T_q + U_f H_{q-1} + W_{cf} \cdot C_{q-1} + b_f), \quad (4)$$

$$C_q = I_q \cdot \sigma_h(W_c T_q + U_c H_{q-1} + b_c) + F_q \cdot C_{q-1}, \quad (5)$$

$$O_q = \sigma(W_o T_q + U_o H_{q-1} + V_o \cdot C_q + b_o), \quad (6)$$

$$H_q = O_q \cdot \sigma_h(C_q). \quad (7)$$

where  $\cdot$  represents the Hadamard product of two vectors, and  $\sigma(\cdot)$  and  $\sigma_h(\cdot)$  are activation functions.  $\sigma(\cdot)$  is traditionally set to be a function of domain  $[0, 1]$  to control the information flow through

time. And  $\sigma_h(\cdot)$  is often set to be a centered activation function.

### 3.3 Periodic features

People are used to repeating some behaviors on a same time period of day, e.g. people routinely go to work in the morning and go home in the evening. This is why we can observe a strong periodicity within traffic flow. The periodicity of traffic flow have been identified as a major contributing factor for the traffic flow forecasting. A desirable model that successfully characterize the periodicities can accurately forecast future traffic flow. In this subsection, we propose to use LSTM to capture the features of periodicities for traffic flow forecasting. The inputs for periodicities at time point  $t$  are given as following:

$$\mathbf{S}^d = \begin{bmatrix} s_1(t^d - n^d) & s_1(t^d - n^d + 1) & \cdots & s_1(t^d + n^d) \\ s_2(t^d - n^d) & s_2(t^d - n^d + 1) & \cdots & s_2(t^d + n^d) \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ s_p(t^d - n^d) & s_p(t^d - n^d + 1) & \cdots & s_p(t^d + n^d), \end{bmatrix}. \quad (8)$$

$$\mathbf{S}^w = \begin{bmatrix} s_1(t^w - n^w) & s_1(t^w - n^w + 1) & \cdots & s_1(t^w + n^w) \\ s_2(t^w - n^w) & s_2(t^w - n^w + 1) & \cdots & s_2(t^w + n^w) \\ \vdots & \vdots & & \vdots \\ \vdots & \vdots & & \vdots \\ s_p(t^w - n^w) & s_p(t^w - n^w + 1) & \cdots & s_p(t^w + n^w), \end{bmatrix}. \quad (9)$$

where  $t^d$  and  $t^w$  denote the same time point of  $t$  in last day and last weekday respectively,  $n^d$  and  $n^w$  denote the time lags of daily periodicity and weekly periodicity respectively.

For features of daily periodicity, the input of LSTM is denoted as  $T^d = (T_0^d, T_1^d, \dots, T_{2n^d}^d)$  where  $T_q^d = [s_1(t^d - n^d + q), s_2(t^d - n^d + q), \dots, s_p(t^d - n^d + q)]^T$  in Eq (8). The same input can be easily adapted to weekly periodicity. It is obvious that there exists correlation between  $\mathbf{S}^d$  and  $\mathbf{S}^w$ , so as illustrated in Fig.1, the connection between LSTMs for daily periodicity and weekly periodicity are added. With the LSTM architecture given in Fig.1, we can capture features  $H_0^w, H_1^w, \dots, H_{2n^w}^w$  of weekly periodicity and features  $H_0^d, H_1^d, \dots, H_{2n^d}^d$  of daily periodicity to forecast future traffic flow at time point  $t$ .

### 3.4 Feature-level fusion

As shown in Fig. 1, the proposed CLTFP captures spatial features  $h^1, h^2, \dots, h^{nc}$ , short-term temporal features  $H_0, H_1, \dots, H_{n-1}$ , weekly periodic features  $H_0^w, H_1^w, \dots, H_{2n^w}^w$ , and daily periodic features  $H_0^d, H_1^d, \dots, H_{2n^d}^d$ . In order to fuse them to perform short-term traffic flow forecasting, we

concatenate all the features sequentially into a feature vector, then add a regression layer to perform forecasting. The objective function of regression is the sum of square errors between predicted value  $s_1^\epsilon(t), s_2^\epsilon(t), \dots, s_n^\epsilon(t)$  and future value  $s_1(t), s_2(t), \dots, s_n(t)$ . There may be redundancies between the features captured by CLTFP. To handle the feature redundancy problem, we add a sparsity regularization in weights of fully connection layer thus our model is likely to assign a weight close to zero to redundant feature.

## 4 Experiments

In this section, we use traffic flow data from PeMS(<http://pems.eecs.berkeley.edu/>) to evaluate the proposed CLTFM. The CLTFM is compared with several state-of-the-art forecasting methods with deep architecture: LSTM (Ma et al., 2015), SAE (Lv et al., 2015), a shallow neural network and the gradient boosting regression tree (GBRT) method (Zhang and Haghani, 2015). In additional, we perform a study to discover the performance contribution of different features of CLTFM from the view of Granger Causality. All experiments are performed by a PC (CPU: Intel Xeon(R) E5-2620 2.1GHz, 64GB memory, GPU: NVIDIA Tesla K40C).

### 4.1 Datasets

The peculiar traffic flow data from PeMS throughout North-bound I-405 trip are used for our experiments. Traffic flow of 33 locations given in Fig.3 on this trip are used for our study. The particular time period used in this paper is from 01/04/2014 to 30/06/2015. The traffic volume are aggregated every 5 min. Thus, one detector preserves 288 data points per day. We use earlier 110000 past-future pairs to train all models, and the rest pairs are used as test data.

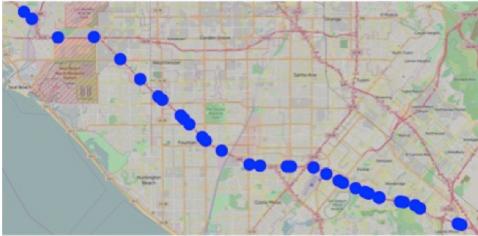


Figure 3: Traffic flow locations studied in this paper

### 4.2 Experimental setup

For all methods, time window size  $n$  of  $\mathbf{S}$  is set as 15, which means that 75 min historical data are used to perform forecasting of next 5 min. The time lags of daily periodicity  $n^d$  and weekly periodicity  $n^w$  for long-term inputs of CLTFP are set as 6, which means that the traffic flow before and after 30 minutes in previous day and weekday are used to generate forecasting.

For 1D CNN structure of CLTFP, a 3-layer fully convolution structure is used, there are 30 filters in each layer, the filter lengths of first 2 layers are set as 3, the filter length of last layer is set as 2, SReLU (Jin et al., 2015) is used as the activation function of CNN. For the LSTM capturing short-term temporal features  $H_0, H_1, \dots, H_{n-1}$ , the feature dimension of each time point is set as 40. For the LSTM capturing long-term features  $H_0^d, H_1^d, \dots, H_{2n^d}^d$  and  $H_0^w, H_1^w, \dots, H_{2n^w}^w$ , the feature dimension of each time point is set as 25. For the regression layer of CLTFP, the  $l_1$  norm regularizer on weights is set as 0.002.

CLTFP are trained based on Adamax optimizer (Kingma and Ba, 2014), we randomly select 10% training data as validation dataset to control earlystopping. The architecture of CLTFP are built upon Keras framework (Chollet, 2015). The structures and parameters for other methods are set according to the reports on corresponding papers.

### 4.3 Comparison results

In this paper, the mean absolute percentage error (*MAPE*) is used to compare the performance of traffic forecasting. The *MAPE* will be lower if the traffic volumes are higher. In observance of this, this paper also applies the mean absolute error (*MAE*) as a complementary measure for *MAPE*,

$$MAE = \frac{1}{n_p} \sum_{t=1, s=1}^n |z_{st} - N_{st}|, \quad (10)$$

$$MAPE = \frac{1}{n_p} \sum_{t=1, s=1}^n \frac{|z_{st} - N_{st}|}{N_{st}} \times 100\%, \quad (11)$$

where  $z_{st}$  = predicted traffic flow at time point  $t$  on location  $s$ ;  $N_{st}$  = actual traffic flow;  $n_p$  = number of predictions. The aim of indexes *MAE* and *MAPE* is to measure the errors between predicted values and actual values. The forecasting correctness of spatial distribution is also an important index for this comparison as we perform prediction on multiple locations, thus we define an average correlation error (ACE) to measure the ability of spatial distribution forecasting:

$$ACE = \frac{1}{n_t} \sum_{t=1}^n Corr(z_{:t}, N_{:t}), \quad (12)$$

where  $z_{:t}$  = predicted traffic flow vector at time point  $t$ ;  $N_{:t}$  = actual traffic flow vector;  $n_t$  = number of prediction steps.

Table.1 gives quantitative results of CLTFP, LSTM, SAE, shallow NN and GBRT, it can be found that CLTFP achieves better performance than other methods in terms of prediction accuracy and spatial distribution. The reason is that CLTFP makes full use of spatial distribution, short-term temporal variability and long-term periodicities.

Table 1: The quantitative results of different methods

features	MAE	MAPE	ACE
CLTFP	19.37	7.36%	0.9263
LSTM	21.53	8.55%	0.9137
SAE	20.36	8.07%	0.9198
NN	20.61	8.31%	0.9174
GBRT	22.52	8.52%	0.9109

#### 4.4 Analysis of Features

One constant criticism of using neural networks on transportation area has been that they are black box models, with little understanding of how the networks work and what knowledge the networks find from data. Recently, it is found that Granger Causality, which characterizes the causality based on incremental predictability, can be adopted to understand black-box typed prediction approaches (Li et al., 2015). In this subsection, we focus on analyzing proposed CLTFP from the view of Granger Causality.

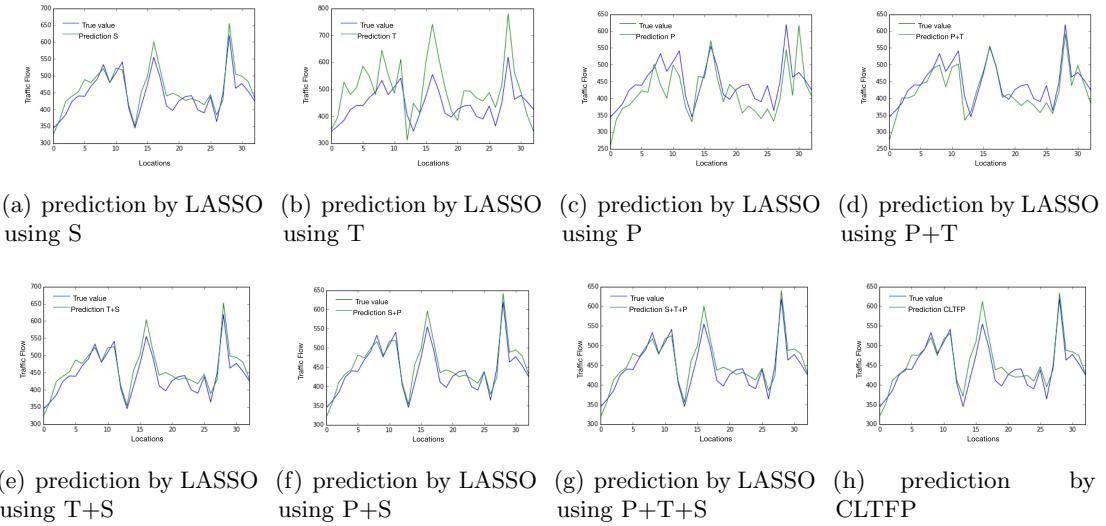


Figure 4: Forecasting results of different LASSO models and proposed CLTFP at one time point

The experimental analysis is conducted as following: we leverage available spatial features (S), short-term temporal features (T) and periodic features (P) generated from well-trained model CLTFP to train several Lasso models, and then conduct a comparison between forecasting results of those Lasso models with different combinations of features. All Lasso models are fit with Least Angle Regression, the penalty terms of  $l_1$  priors are all set as 0.002.

The experimental results are given in Table. 2. This table is quite revealing in several ways. 1. We can generally conclude that more types of feature help to build better prediction results. As the

Table 2: The quantitative results of different features and their combinations (S: spatial features, T: Short-term temporal features, P: Periodicity features)

methods	MAE	MAPE	ACE
S	19.83	7.91%	0.9302
T	54.43	28.40%	0.8532
P	40.79	18.06%	0.7996
P+T	35.07	14.59%	0.8763
S+T	19.61	7.49%	0.9314
S+P	19.47	7.38%	0.9312
S+T+P	19.32	7.29%	0.9323

improved predictability can be achieved by adding those features, we can draw that future traffic flow are dependent on all those information from the view of Granger Causality. 2. The model with feature S significantly outperforms the model with feature T and P, it suggests that future traffic flow of the studied corridor is heavily dependent on spatial information of near-term traffic flow though the temporal information in both near-term and long-term (last week and last day) have some influence. 3. The model with feature P achieves lower errors than model with feature T, however, it has weaker predictability on spatial distribution of future traffic flow. It suggests that the travel habits have more influence on the total traffic flow, but near-term traffic flow on transportation network is more related to the future distribution of traffic flow. 4. The model with features S, T and P even outperforms our well-trained CLTFP in terms of *MAE*, *MAPE* and *ACE*. It indicates that the performance of traffic flow forecasting can be promoted by using a proper regression model on features generated from a well-trained neural networks. The similar phenomena can be also found in Fig. 4, which gives quantitative visualization of forecasting results at one time point.

## 5 Conclusions and future work

A novel deep learning based short-term traffic flow forecasting method CLTFP combined with CNN and LSTM is proposed in this paper, the forecasting results of CLTFP are encouraging, it indicates the potential of CNN and LSTM on transportation applications. Moreover, incremental predictability is applied to analyze the black-box typed forecasting method, the analysis shows that neural network based forecasting method can provide many meaningful knowledges of traffic flow.

This proposed CLTFP admits many improvements and extensions:

1. The features captured by LSTM achieves only modest forecasting accuracy, some more complex structures, for example, convolutional LSTM structure (Xingjian et al., 2015) can be an alternative.
2. Traffic flow are affected by many other factors such as weather, social event and state of the roads. How to exploit those information as auxiliary information is a future direction.

3. The applications of our model on general transportation network and similar spatial-temporal data on other domain are straightforward extensions.

## References

- Bishop, C. M. 2006. Pattern recognition. *Machine Learning* 128.
- Chan, K. Y.; Dillon, T. S.; Singh, J.; and Chang, E. 2012. Neural-network-based models for short-term traffic flow forecasting using a hybrid exponential smoothing and levenberg–marquardt algorithm. *IEEE Transactions on Intelligent Transportation Systems* 13(2):644–654.
- Chen, Q.; Song, X.; Yamada, H.; and Shibasaki, R. 2016. Learning deep representation from big and heterogeneous data for traffic accident inference. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- Chollet, F. 2015. keras. <https://github.com/fchollet/keras>.
- Donahue, J.; Anne Hendricks, L.; Guadarrama, S.; Rohrbach, M.; Venugopalan, S.; Saenko, K.; and Darrell, T. 2015. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2625–2634.
- Graves, A.; Jaitly, N.; and Mohamed, A.-r. 2013. Hybrid speech recognition with deep bidirectional lstm. In *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*, 273–278. IEEE.
- Henaff, M.; Bruna, J.; and LeCun, Y. 2015. Deep convolutional networks on graph-structured data. *arXiv preprint arXiv:1506.05163*.
- Hinton, G. E., and Salakhutdinov, R. R. 2006. Reducing the dimensionality of data with neural networks. *Science* 313(5786):504–507.
- Hinton, G. E.; Osindero, S.; and Teh, Y.-W. 2006. A fast learning algorithm for deep belief nets. *Neural computation* 18(7):1527–1554.
- Hochreiter, S., and Schmidhuber, J. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- Huang, W.; Song, G.; Hong, H.; and Xie, K. 2014. Deep architecture for traffic flow prediction: deep belief networks with multitask learning. *IEEE Transactions on Intelligent Transportation Systems* 15(5):2191–2201.
- Jin, X.; Xu, C.; Feng, J.; Wei, Y.; Xiong, J.; and Yan, S. 2015. Deep learning with s-shaped rectified linear activation units. *arXiv preprint arXiv:1512.07030*.

- Karlaftis, M. G., and Vlahogianni, E. I. 2011. Statistical methods versus neural networks in transportation research: Differences, similarities and some insights. *Transportation Research Part C: Emerging Technologies* 19(3):387–399.
- Kingma, D., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, 1097–1105.
- LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *Nature* 521(7553):436–444.
- Li, L.; Su, X.; Wang, Y.; Lin, Y.; Li, Z.; and Li, Y. 2015. Robust causal dependence mining in big data network and its application to traffic flow predictions. *Transportation Research Part C: Emerging Technologies* 58:292–307.
- Lv, Y.; Duan, Y.; Kang, W.; Li, Z.; and Wang, F.-Y. 2015. Traffic flow prediction with big data: a deep learning approach. *IEEE Transactions on Intelligent Transportation Systems* 16(2):865–873.
- Ma, X.; Tao, Z.; Wang, Y.; Yu, H.; and Wang, Y. 2015. Long short-term memory neural network for traffic speed prediction using remote microwave sensor data. *Transportation Research Part C: Emerging Technologies* 54:187–197.
- Min, W., and Wynter, L. 2011. Real-time road traffic prediction with spatio-temporal correlations. *Transportation Research Part C: Emerging Technologies* 19(4):606–616.
- Peris, Á.; Bolaños, M.; Radeva, P.; and Casacuberta, F. 2016. Video description using bidirectional recurrent neural networks. *arXiv preprint arXiv:1604.03390*.
- Polson, N., and Sokolov, V. 2016. Deep learning predictors for traffic flows. *arXiv preprint arXiv:1604.04527*.
- Qi, Y., and Ishak, S. 2014. A hidden markov model for short term prediction of traffic conditions on freeways. *Transportation Research Part C: Emerging Technologies* 43:95–111.
- Shahsavari, B., and Abbeel, P. 2015. Short-term traffic forecasting: Modeling and learning spatio-temporal relations in transportation networks using graph neural networks.
- Springenberg, J. T.; Dosovitskiy, A.; Brox, T.; and Riedmiller, M. 2014. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806*.
- Tan, H.; Feng, G.; Feng, J.; Wang, W.; Zhang, Y.-J.; and Li, F. 2013. A tensor-based method for missing traffic data completion. *Transportation Research Part C: Emerging Technologies* 28:15–27.
- Tan, H.; Wu, Y.; Shen, B.; Jin, P. J.; and Ran, B. 2016a. Short-term traffic prediction based on dynamic tensor completion. *IEEE Transactions on Intelligent Transportation Systems* 17(8):2123–2133.

- Tan, H.; Xuan, X.; Wu, Y.; Zhong, Z.; and Ran, B. 2016b. A comparison of traffic flow prediction methods based on dbn. In *CICTP 2016*. 273–283.
- Vincent, P.; Larochelle, H.; Bengio, Y.; and Manzagol, P.-A. 2008. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, 1096–1103. ACM.
- Vinyals, O.; Toshev, A.; Bengio, S.; and Erhan, D. 2015. Show and tell: A neural image caption generator. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3156–3164.
- Wang, J.; Yu, L.-C.; Lai, K. R.; and Zhang, X. 2016. Dimensional sentiment analysis using a regional cnn-lstm model. In *The 54th Annual Meeting of the Association for Computational Linguistics*, 225.
- Wang, J.; Deng, W.; and Guo, Y. 2014. New bayesian combination method for short-term traffic flow forecasting. *Transportation Research Part C: Emerging Technologies* 43:79–94.
- Wu, Z.; Wang, X.; Jiang, Y.-G.; Ye, H.; and Xue, X. 2015. Modeling spatial-temporal clues in a hybrid deep learning framework for video classification. In *Proceedings of the 23rd ACM international conference on Multimedia*, 461–470. ACM.
- Xingjian, S.; Chen, Z.; Wang, H.; Yeung, D.-Y.; Wong, W.-k.; and Woo, W.-c. 2015. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in Neural Information Processing Systems*, 802–810.
- Yang, S.; Kalpakis, K.; and Biem, A. 2014. Detecting road traffic events by coupling multiple timeseries with a nonparametric bayesian method. *IEEE Transactions on Intelligent Transportation Systems* 15(5):1936–1946.
- Yao, L.; Torabi, A.; Cho, K.; Ballas, N.; Pal, C.; Larochelle, H.; and Courville, A. 2015. Describing videos by exploiting temporal structure. In *Proceedings of the IEEE International Conference on Computer Vision*, 4507–4515.
- Zhang, Y., and Haghani, A. 2015. A gradient boosting method to improve travel time prediction. *Transportation Research Part C: Emerging Technologies* 58:308–324.
- Zheng, X.; Chen, W.; Wang, P.; Shen, D.; Chen, S.; Wang, X.; Zhang, Q.; and Yang, L. 2016. Big data for social transportation. *IEEE Transactions on Intelligent Transportation Systems* 17(3):620–630.