

DeepSTD: Mining Spatio-Temporal Disturbances of Multiple Context Factors for Citywide Traffic Flow Prediction

Chuanpan Zheng^{ID}, Xiaoliang Fan^{ID}, Senior Member, IEEE, Chenglu Wen^{ID}, Senior Member, IEEE, Longbiao Chen, Cheng Wang^{ID}, Senior Member, IEEE, and Jonathan Li, Senior Member, IEEE

Abstract—Deep learning techniques have been widely applied to traffic flow prediction, considering underlying routine patterns, and multiple context factors (e.g., time and weather). However, the complex spatio-temporal dependencies between inherent traffic patterns and multiple disturbances have not been fully addressed. In this paper, we propose a two-phase end-to-end deep learning framework, namely DeepSTD to uncover the spatio-temporal disturbances (STD) to predict the citywide traffic flow. In the *STD Modeling* phase, we propose an STD modeling method to model both the different regional disturbances caused by various region functions and the spatio-temporal propagating effects. In the *Prediction* phase, we eliminate the STD from the historical traffic flow to enhance the leaning of inherent traffic patterns and combine the STD at the prediction time interval to consider the future disturbances. The experimental results on two real-world datasets demonstrate that DeepSTD outperforms the state-of-the-art methods.

Index Terms—Traffic flow prediction, spatio-temporal disturbances, deep learning, intelligent transportation systems.

I. INTRODUCTION

ACCURATE citywide traffic flow prediction is of great importance to the development of modern intelligent transportation systems (ITS) [1]. It has the potential to alleviate traffic congestion, reduce fuel consumption, and thus enhance the overall efficiency of transportation networks [2].

Traffic flow prediction heavily depends on historical traffic data collected from various sensor sources, such as vehicle license plate recognition (VLPR) devices [3], ridesharing platform [4], loop detectors [5], taxi GPS systems [6], and call detail records [7]. Tremendous studies have been conducted

Manuscript received February 18, 2018; revised July 30, 2018, January 25, 2019 and June 8, 2019; accepted July 25, 2019. Date of publication August 9, 2019; date of current version August 28, 2020. This work was supported in part by the National Natural Science Foundation of China under Grant 61872306, Grant 61771413, and Grant 61802325. The Associate Editor for this paper was F.-Y. Wang. (Corresponding authors: Xiaoliang Fan; Jonathan Li.)

C. Zheng, X. Fan, C. Wen, L. Chen, and C. Wang are with the Fujian Key Laboratory of Sensing and Computing for Smart Cities, School of Informatics, Digital Fujian Institute of Urban Traffic Big Data Research, Xiamen University, Xiamen 361005, China (e-mail: zhengchuanpan@stu.xmu.edu.cn; faxniaoliang@xmu.edu.cn; clwen@xmu.edu.cn; longbiaochen@xmu.edu.cn; cwang@xmu.edu.cn).

J. Li is with the Fujian Key Laboratory of Sensing and Computing for Smart Cities, School of Informatics, Digital Fujian Institute of Urban Traffic Big Data Research, Xiamen University, Xiamen 361005, China, and also with the Department of Geography and Environmental Management, Department of Systems Design Engineering, University of Waterloo, Waterloo, ON N2L 3G1, Canada (e-mail: junli@xmu.edu.cn; junli@uwaterloo.ca).

Digital Object Identifier 10.1109/TITS.2019.2932785

on traffic flow prediction in past decades, such as autoregressive integrated moving average (ARIMA) [8], support vector regression (SVR) [9], artificial neural networks (ANN) [10].

Indeed, traffic tends to exhibit inherent patterns [11], [12] and thus shows high potential predictability [13], [14]. Recently, deep neural networks show more superior performance than shallow models in traffic flow prediction [15]–[17], due to the strong capability of capturing spatial or temporal dependencies. Meanwhile, traffic could be affected by multiple factors such as time [18] (e.g., travel peaks during holidays), weather [19] (e.g., a sudden rainstorm in evening peak), and social events [20] (e.g., a pop concert to be held in downtown). We coin the impact of these factors on the traffic flow as disturbance. Some researchers have considered multiple disturbances, separately or collectively, in the prediction models [21]–[24].

However, we believe that the major limitations of state-of-the-art deep learning methods are twofold.

First, the impact of multiple disturbances (e.g., time, weather) in the spatio-temporal view has not been fully addressed. On the one hand, the disturbances vary from region to region due to different region functions. For example, the traffic flow in tourist areas will significantly decrease during rough weather, whereas office areas will be less affected. On the other hand, it is probable that the disturbances could propagate in both spatial [20], [21] (i.e., extend from one region to its neighboring regions), and temporal [18] (i.e., spread over time spans) dimensions.

Second, many works usually directly utilize the historical traffic data to train deep neural networks, which might inevitably introduce noises into the model and thus degrade the prediction accuracy. Indeed, the historical traffic data contain not only regular traffic patterns [11], but also multiple disturbances. Thus, it is necessary to minimize the impact of multiple disturbances from historical traffic data before feeding into prediction models.

In this paper, we propose **DeepSTD**, a two-phase framework to uncover the spatio-temporal disturbances (STD) to predict the citywide traffic flow. Specifically, in the *STD Modeling* phase, we first partition the city into regions and calculate the inherent influence of each point-of-interest (POI) category in every region as Inherent Influence Factor (IIF), which implies the potential region functions. Second, we model the impact of multiple factors on each POI category as Disturbance

Influence Factor (DIF). Third, we fuse IIF and DIF as STD via 3D convolutional neural networks (3D CNN) to model the spatio-temporal propagating effects. In the *Prediction* phase, we first eliminate the STD from the historical traffic flow. Then, the traffic flow without STD is fed into a residual neural network (ResNet) to learn the inherent traffic patterns. Finally, the output of ResNet is combined with the STD of the prediction time interval to generate the final prediction. We evaluate the performance of DeepSTD on two real-world datasets from two major cities of China. Specifically, the Xiamen dataset contains 3.15 billion VLPR records, 83,961 POIs, and the corresponding time and weather data from August 1st, 2015 to August 31st, 2016. The Chengdu dataset includes 1.10 billion GPS records provided by Didi Chuxing,¹ 69,049 POIs, and the corresponding time and weather data in November 2016. The experimental results demonstrate that DeepSTD outperforms state-of-the-art methods.

The contributions of this study are summarized as follows.

- To effectively model the multiple disturbances on the citywide traffic flow in the spatio-temporal view, we design a STD modeling method. On one hand, to model the different regional disturbances caused by various region functions, we calculate the inherent influence factor (IIF) of each POI category in every region, and model the disturbance influence factor (DIF) of multiple factors on every POI category in each time interval. On the other hand, to capture the spatio-temporal propagating effects of STD, we apply a 3D CNN based method to fuse IIF and DIF as STD.
- To further incorporate the STD into the citywide traffic flow prediction, we design a combinational mechanism. We first eliminate the STD from the historical traffic flow to represent the inherent traffic patterns. Then, we feed the traffic flow without STD into a ResNet. Finally, to further consider the future disturbances, the STD of the prediction time interval is combined with the output of ResNet as the final prediction. We conduct extensive experiments to evaluate the effectiveness of both the STD modeling method and the combinational mechanism, respectively.

The rest of this paper is organized as follows. Section II reviews the studies on traffic flow prediction. Section III presents some preliminaries of this study. Section IV details the method of DeepSTD. Section V describes the experimental setup. Section VI compares the predictive performance between DeepSTD and other methods on two real-world datasets. Finally, Section VII concludes this paper.

II. RELATED WORKS

Traffic flow prediction has gained increasing attention with the rapid development and widely deployment of intelligent transportation systems in past decades [25], [26]. Existing approaches can be generally divided into three categories: parametric, non-parametric and hybrid approaches [27].

Parametric approaches are also known as model-driven methods, where the model structure is predetermined

based on certain theoretical assumptions [15], including historical average [28], autoregressive integrated moving average (ARIMA) [8], seasonal ARIMA [29], etc. [30], [31].

Non-parametric methods, such as support vector regression (SVR) [9], [32], artificial neural network (ANN) [10], [22] are designed to capture the non-linear variations in traffic data. Recently, deep neural networks has been successfully applied in various domains, e.g., computer vision, speech recognition, and natural language processing [33]–[36]. This success inspired several attempts to utilize deep learning techniques on traffic prediction problems. Huang *et al.* [37] introduced deep learning approaches into transportation with deep belief networks. Lv *et al.* [25] designed a stacked autoencoder (SAE) model to learn latent traffic features for traffic flow prediction. Ma *et al.* [15] developed long short-term memory (LSTM) networks to capture the long-term temporal dependency for traffic speed prediction. In general, non-parametric techniques provide more accurate prediction results than parametric approaches due to the strong generalization and learning ability [38].

To obtain adaptive models, some works explored hybrid methods. Tan *et al.* [39] applied moving average, exponential smoothing, and ARIMA, to forecast three relevant time series (i.e., weekly similarity, daily similarity, and hourly time series), and then adopted ANN to aggregate the three forecasting values as the final prediction. Zheng *et al.* [40] proposed a neural network model that combines Bayesian and ANN to predict short-term freeway traffic flows. Results show that hybrid models could outperform singular predictors [11].

Recently, many studies focus on the prediction of large-scale transportation networks. Ma *et al.* [16] applied convolutional neural networks (CNNs) to extract the spatio-temporal features from traffic data to predict large-scale, network-wide traffic speed. Zhang *et al.* [23] proposed ST-ResNet, which employed residual neural networks [41] to model the spatio-temporal dependencies for citywide crowd flows prediction. These methods capture the underlying spatio-temporal traffic patterns, but may introduce noises because they directly feed the historical traffic flow into deep neural networks.

Furthermore, transportation systems can be heavily affected by multiple factors, e.g., time, weather, and social events. The impact of weather on the traffic flow was studied in [19]. Researchers in [20] adopted real-time social media data to extract event information to forecast the subway passenger flow under event occurrences. A study presented in [22] demonstrated that the time context such as time-of-day and day-of-week were very effective for predicting traffic flows. However, these methods fail to capture the different regional disturbances caused by various region functions and neglect the spatio-temporal propagating effects.

In summary, many works have been developed for traffic flow prediction. However, the complex dependencies between inherent traffic patterns and multiple disturbances in the spatio-temporal view have not been fully addressed. In this paper, we propose an end-to-end deep-learning-based approach to effectively uncover the spatio-temporal disturbances (STD) and integrate the STD with the inherent traffic patterns in a unified framework to predict the citywide traffic flow.

¹<http://www.didichuxing.com>

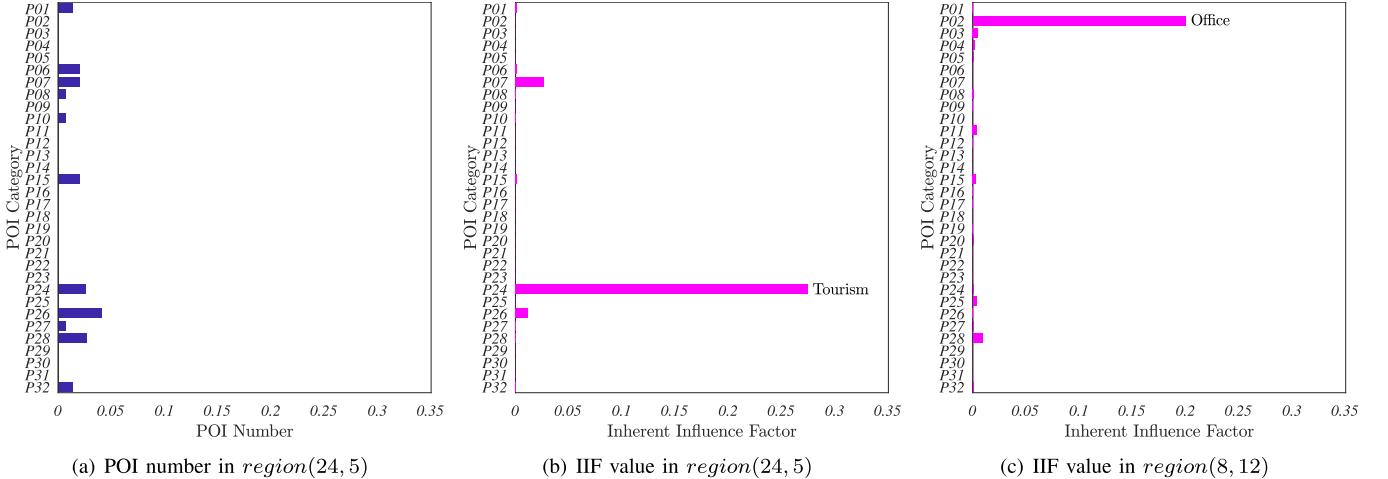


Fig. 4. The number, IIF value of each POI category in *region*(24, 5) and *region*(8, 12). The list of POI categories and the two regions are presented in Table II and Fig. 1, respectively.

c) The *imbalance degree (IBD)* of the POI category distributed in the city. The imbalance distribution of the POI category could strongly influence human mobility over a large scale of the transportation network because it offers the unique service (e.g., airport, railway station) that people can only find in few regions [44]. For example, as shown in Fig. 3(c), the airport would be more influential than shops as the airport could attract visitors from all over the city.

We employ the Shannon entropy to measure the imbalance degree. Shannon entropy is a measurement of uncertainty or randomness of a single random variable [45], in which higher values indicate more uniform distribution. Formally, the Shannon entropy of k^{th} POI category is defined as

$$S(k) = - \sum_i \sum_j (Den(i, j, k) \times \log Den(i, j, k)), \quad (4)$$

where $Den(i, j, k)$ is defined in Eq. 3. When the POI category distributes uniformly in all regions, the Shannon entropy reaches the maximum:

$$S_{max} = \log(I \times J). \quad (5)$$

Thus, we define the IBD of k^{th} POI category as

$$IBD(k) = 1 - \frac{S(k)}{S_{max}}. \quad (6)$$

The $IBD(k)$ denotes the imbalance degree of k^{th} POI category distributed in the city, in which higher value means the POI category is more imbalance and more influential.

Finally, the IIF of k^{th} POI category of $region(i, j)$ is defined as

$$IIF(i, j, k) = Fre(i, j, k) \times Den(i, j, k) \times IBD(k). \quad (7)$$

Therefore, the IIF of all POI categories in all regions is represented as $IIF \in \mathbb{R}^{I \times J \times K}$, where K is the number of POI categories.

The IIF implies the potential region functions in every region. For example, Fig. 4(a) and 4(b) present the number, IIF value (normalized into $[0, 1]$, respectively) of each POI category in *region*(24, 5) of Xiamen. *Region*(24, 5) is a well-known tourist area, including the famous tourist spot,

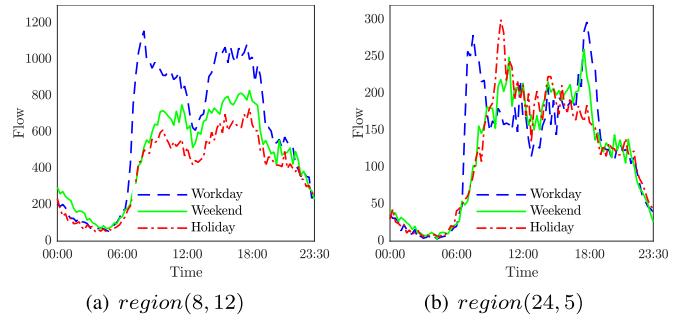


Fig. 5. The impact of time context on *region*(8, 12) and *region*(24, 5). The blue line is October 12nd, 2015 (workday), the green line is October 11th, 2015 (weekend), and the red line is October 5th, 2015 (holiday).

i.e., *South Putuo Temple*. We observe that the proposed IIF value effectively uncovers the region function, as shown in Fig. 4(b).

2) *Disturbance Influence Factor (DIF)*: The disturbances of multiple factors on different regions differ greatly due to different region functions. Fig. 5(a) and 5(b) depict the different influence patterns of time context on *region*(8, 12) and *region*(24, 5). As shown in Fig. 4(c), *region*(8, 12) tends to be an office area, and the traffic flow on weekend and holiday is less than that on workday. Whereas the phenomenon is different in *region*(24, 5) (a tourist area, see Fig. 4(b)). Similarly, as shown in Fig. 6(a) and (b), the impact of the rainstorm is much greater on *region*(24, 5) than that on *region*(8, 12).

To capture the different influence patterns of multiple factors on different regions, we design a simple yet effective component to model the disturbance influence on each POI category. Specifically, we stack two fully-connected layers upon the context feature $\mathcal{C} \in \mathbb{R}^{(N+1) \times D_C}$, where the first layer is viewed as an embedding layer followed by an activation (i.e., *relu* [46]), and the second layer is used to extract the disturbance influence of multiple factors on each POI category with K neurons (K is the number of POI categories). Consequently, we obtain an output represented as $DIF \in \mathbb{R}^{(N+1) \times K}$, indicating the disturbance influence of multiple context factors

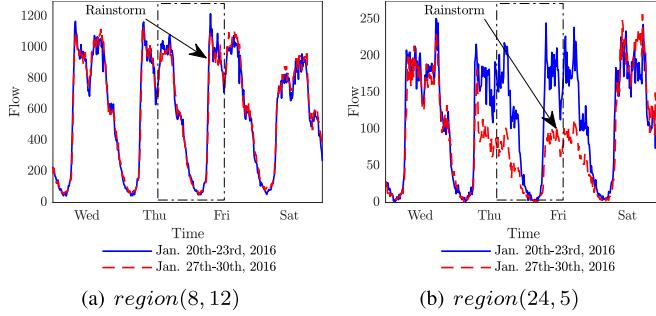


Fig. 6. The impact of weather context on $\text{region}(8, 12)$ and $\text{region}(24, 5)$. The blue line is January 20th-23rd, 2016 and the red line is January 27th-30th, 2016.

on K POI categories in $N + 1$ time intervals (i.e., historical N time intervals and the prediction time interval $t + 1$).

3) *Fusion*: We propose a two-step fusion method to fuse IIF and DIF. In the first step, we combine IIF with DIF as

$$\text{IDIF}(n, i, j, k) = \text{IIF}(i, j, k) \times \text{DIF}(n, k). \quad (8)$$

$\text{IDIF}(n, i, j, k)$ indicates the disturbance of multiple context factors on k^{th} POI category in $\text{region}(i, j)$ at time interval n . The IDIF can be denoted as $\text{IDIF} \in \mathbb{R}^{(N+1) \times I \times J \times K}$.

In the second step, we consider the spatio-temporal propagating effects of the disturbances. On the one hand, the disturbances could propagate along the spatial dimension. The disturbances in one region could lead to a chain of reactions that influence large-scale transportation networks [20], [21]. On the other hand, the disturbances would propagate along the temporal dimension. As shown in Fig. 6(b), we observe that the impact of the rainstorm is not only limited in the rainstorm periods but also before and after the rainstorm. Based on the observations, we employ a 3D CNN upon IDIF to model the spatio-temporal propagating effects. 3D CNN could extract features from both spatial and temporal dimensions by adopting 3D convolutional operations [42], which has been widely applied in various domains [47]–[49].

We adopt six layers of 3D CNN to model the spatio-temporal propagating effects. Specifically, we use $3 \times 3 \times 3$ time-space filters for all convolutional layers. We use 64 filters in first five layers and use one filter in the output layer. To keep size constant, we do not apply pooling layers and utilize zero-padding in all three dimensions in all convolutional operations. Consequently, we obtain an output $\text{STD} \in \mathbb{R}^{(N+1) \times I \times J \times 1}$, indicating the spatio-temporal disturbances in all $I \times J$ regions of $(N + 1)$ time intervals (historical N time intervals and the prediction time interval $t + 1$). Then, we slice the spatio-temporal disturbances of historical N time intervals from STD and reshape it to $\text{STD}_t^{\text{history}} \in \mathbb{R}^{I \times J \times N}$. Likewise, we could obtain the STD of time interval $t + 1$ as $\text{STD}_{t+1} \in \mathbb{R}^{I \times J \times 1}$.

B. Prediction

As the traffic flow contains both regular traffic patterns and multiple disturbances, it is essential to avoid the effect of multiple disturbances when modeling the traffic patterns. To address this issue, in the *Prediction* phase, we propose a three-step approach to incorporate STD with the historical

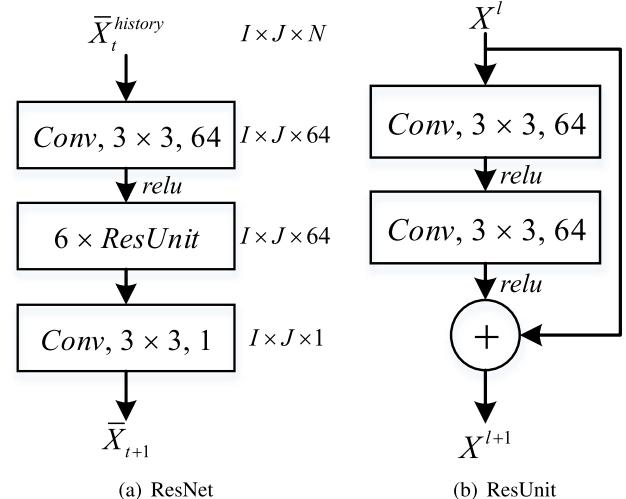


Fig. 7. The structure of the prediction model.

traffic flow to uncover the complex dependencies between inherent traffic patterns and multiple disturbances for predicting the citywide traffic flow. In the first step, we eliminate the STD from the historical traffic flow, as

$$\bar{X}_t^{\text{history}} = X_t^{\text{history}} - \text{STD}_t^{\text{history}}, \quad (9)$$

where $\bar{X}_t^{\text{history}} \in \mathbb{R}^{I \times J \times N}$ denotes the historical traffic flow without disturbances. It would be easier for neural networks to learn the inherent traffic patterns from $\bar{X}_t^{\text{history}}$ rather than directly feed the raw traffic flow X_t^{history} into prediction models.

In the second step, we feed $\bar{X}_t^{\text{history}}$ into a prediction model to capture the underlying traffic patterns for prediction, as

$$\bar{X}_{t+1} = f(\bar{X}_t^{\text{history}}), \quad (10)$$

where f denotes a prediction algorithm. Herein, we employ the residual neural network (ResNet) [41] to construct the prediction model due to its superior learning ability. ResNet allows CNNs to have very deep structures by adopting residual units (ResUnits) [50], as

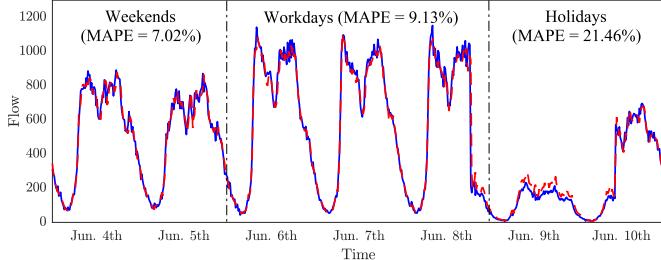
$$X^{l+1} = X^l + \mathbb{F}(X^l), \quad (11)$$

where X^l and X^{l+1} are the input and output of the l^{th} ResUnit, respectively, and $\mathbb{F}(\bullet)$ is the residual function. As shown in Fig. 7, we utilize 14 layers including 6 ResUnits, an input convolution layer, and an output convolution layer to construct the prediction model. In each ResUnit, we stack two convolution layers using 64 filters of 3×3 with zero-padding. We utilize 64 filters in the input layer and one filter in the output layer, respectively. Consequently, we obtain an output as $\hat{X}_{t+1} \in \mathbb{R}^{I \times J \times 1}$, which represents the predicted traffic flow without disturbances at time interval $t + 1$.

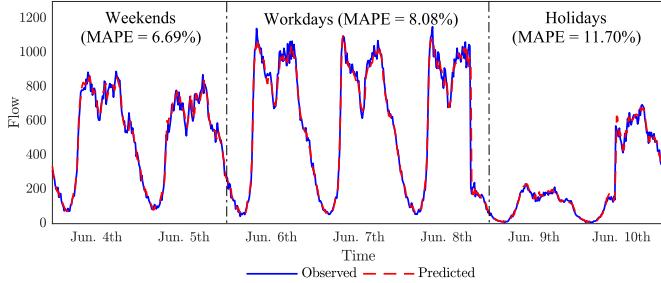
In the third step, we combine the output of ResNet with the STD of time interval $t + 1$ to consider the future disturbances, as

$$\hat{X}_{t+1} = \bar{X}_{t+1} + \text{STD}_{t+1}, \quad (12)$$

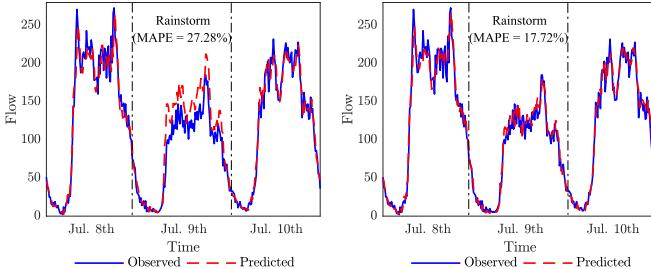
where $\hat{X}_{t+1} \in \mathbb{R}^{I \times J \times 1}$ denotes the final prediction, containing both regular traffic patterns and multiple disturbances at time interval $t + 1$.



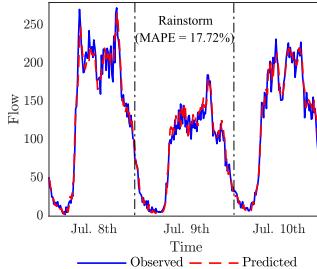
(a) The prediction result of ST-ResNet (MAPE = 12.05%)



(b) The prediction result of DeepSTD (MAPE = 8.72%)

Fig. 8. The prediction results of ST-ResNet and DeepSTD in *region*(8, 12) from June 4th to 10th, 2016.

(a) The prediction result of ST-ResNet

Fig. 9. The prediction results of ST-ResNet and DeepSTD in *region*(24, 5) during the rainstorm (July 9th, 2016).

of DeepSTD and ST-ResNet (generally performs the best among baseline methods) during holidays and rainstorm in Fig. 8 and 9, respectively. As shown in Fig. 8, DeepSTD captures the fluctuation of the sharp decrease of holidays more accurately than ST-ResNet. Likewise, as shown in Fig. 9, it is easy to observe that DeepSTD is more capable of fast responding to the dynamic disturbances of the rainstorm. We believe the poor performance of ST-ResNet comes from the over-reliance on the historical regular patterns. In contrast, our proposed DeepSTD effectively models the spatio-temporal dependencies between inherent traffic patterns and multiple disturbances for prediction, and thus shows high accuracy under unconventional disturbances.

2) Comparison With Variants of DeepSTD:

a) Effect of STD modeling: In the STD modeling, we propose IIF and DIF to model the different regional disturbances caused by various region functions, and 3D CNN based fusion method to capture the spatio-temporal propagating effects. To evaluate the effectiveness of these approaches, we conduct experiments of using partial components. As shown in Table V,

TABLE V
EFFECT OF STD MODELING

Method	RMSE	MAE	MAPE (%)
IIF	52.64	29.21	10.93
DIF	52.37	28.82	10.46
IIF + DIF (without 3D CNN)	51.78	28.58	10.81
POI + DIF (with 3D CNN)	50.49	27.63	10.63
IIF + DIF (with 3D CNN)	49.31	27.32	10.29

TABLE VI
EFFECT OF ELIMINATION AND COMBINATION

Method	RMSE	MAE	MAPE (%)
NoElimination + NoCombination	52.76	29.66	11.63
Elimination + NoCombination	50.16	27.62	11.01
NoElimination + Combination	52.19	28.56	10.63
Elimination + Combination	49.31	27.32	10.29

we observe that: (1) “IIF + DIF (without 3D CNN)” outperforms “IIF” and “DIF”, showing that both IIF and DIF are essential for modeling the STD; (2) “IIF + DIF (with 3D CNN)” performs better than “IIF + DIF (without 3D CNN)”, indicating that the propagating effect is important that should not be neglected; (3) in addition, “IIF + DIF (with 3D CNN)” performs better than “POI + DIF (with 3D CNN)”, illustrating that it is beneficial to calculate the influence of POIs as IIF rather than directly use the number of POIs.

b) Effect of elimination and combination: In the *Prediction* phase, we eliminate STD from the historical traffic flow and combine the STD at the prediction time interval. In order to verify the benefits of these approaches, we conduct experiments of using or without using *Elimination* or *Combination* steps. As presented in Table VI, “NoElimination + NoCombination” means we do not consider STD and directly feed the historical traffic flow into the ResNet to generate the prediction, and it performs the worst. “Elimination + NoCombination” and “NoElimination + Combination” indicate that we adopt either *Elimination* or *Combination* approach, and both of them perform better than “NoElimination + NoCombination”. This demonstrates that it is beneficial to eliminate the STD in historical time intervals and to combine the STD at the prediction time interval. Consequently, “Elimination + Combination” that utilizes both *Elimination* and *Combination* approaches performs the best.

c) Effect of multiple context factors: To evaluate the effect of multiple context factors, we show the performance of various variants that employ different factors in Fig. 10. It can be observed that the variants with combinational factors clearly perform better than those with only one factor, demonstrating the necessity to include various factors. In particular, “T + WN + WC” performs the best, indicating that the integration of multiple context factors is effective for modeling the STD in traffic flow prediction.

d) Impact of Hyperparameters: We further analyze the impact of four hyperparameters in DeepSTD, including the number of historical time intervals, 3D CNN layers, residual units, and filters. As introduced in Section V-A, these four hyperparameters are set as 4, 6, 6, 64, respectively. In the

