

# Traffic Spatial-Temporal Transformer for Traffic Prediction

Zhiyang He

School of Cyberspace Security  
Dongguan University of Technology  
Dongguan, China  
hezhiyang668@163.com

\* Ye Ding

School of Cyberspace Security  
Dongguan University of Technology  
Dongguan, China  
dingye@dgut.edu.cn

**Abstract**—Accurate traffic prediction can help administrators better plan and manage urban traffic, alleviating the traffic pressure. Local spatial-temporal dependency is the strongest and most direct dependency within traffic data. However, recent researches in traffic prediction using stacked or coupled fusion methods to combine temporal and spatial learning networks have not fully captured local spatial-temporal dependency in traffic data. This paper introduces a recurrent neural network structure that captures local spatial-temporal dependency by considering the spatial relationship of each time with its current, past, and future time simultaneously. Additionally, a period enhanced attention mechanism is introduced to capture long-term temporal dependency. Finally, the two modules are combined to construct a Traffic Spatial-Temporal Transformer for traffic prediction. Experimental results demonstrate that the proposed transformer outperforms baselines in terms of prediction accuracy.

**Keywords**- traffic prediction; spatial-temporal dependency

## I. INTRODUCTION

With the development of intelligent sensors and urban computing, Intelligent Transportation Systems (ITS) can collect and analyze voluminous traffic data generated by the busy traffic to extract valuable information. As in [1], ITS can help administrators make informed decisions for rational urban traffic planning and effective management to alleviate traffic pressure. Traffic prediction is essential in ITS, forecasting future traffic conditions from historical data, which involves temporal and non-Euclidean spatial dependencies.

Taking the two dependencies into consideration, the traffic state of a road in the road network at a time is related to its own and the connected roads' state at the current, past, and future

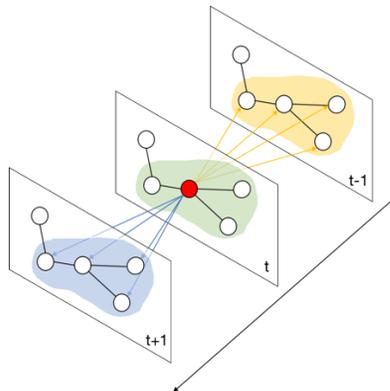


Figure 1. The local spatial-temporal dependency of traffic data.

time, which is called local spatial-temporal dependency, as shown in Fig. 1. Local spatial-temporal dependency is the most direct and strongest dependency present in traffic data, and it is also the most critical dependency that needs to be captured for traffic prediction.

Recent researches on traffic prediction mainly focused on the network structures that fuse temporal learning network (TLN) and spatial learning network (SLN). According to [2], the fusion methods in these network structures can be divided into two categories: stacked fusion and coupled fusion. Stacked fusion connects TLN and SLN serially or in parallel, capturing temporal and spatial dependencies separately, such as TGCN [3], STGCN [4], ASTGCN [5] and GWN [6]. This fusion method leads to each network neglecting other dependency while capturing their respective dependency and fails to capture the local spatial-temporal dependency. The coupled fusion method is usually embedding GNN-based SLN into RNN-based TLN, capturing spatial dependency in each iteration, such as DCRNN [7] and DGCRN [8]. This method can capture the spatial-temporal dependency between the present and past time, but without the future time. BiLSTM [9] can address this issue by processing bidirectionally but with high costs. S-LSTM [10], a recurrent network structure with a message-passing mechanism that treats the entire time series as a hidden state with global state, provides lower cost and more accurate temporal feature extraction compared to BiLSTM.

To achieve more accurate traffic prediction, we enhanced S-LSTM with a dynamic graph generation process during each iteration and utilized it for graph convolution during hidden state updates, which allows us to capture local spatial-temporal dependency in traffic data effectively. Additionally, we designed an enhanced attention mechanism based on period to capture long-term temporal dependency. This mechanism decomposes time series into the frequency domain to obtain periods and strengthens the attention between each time point and its related periodic time points based on frequency. Our contributions are summarized as follows:

- We propose a recurrent network structure named LSTN to simultaneously capture the spatial and temporal dependencies of each road at each time along with the neighbor roads at current, past and future time.
- We propose an enhanced attention named PEA, which enhanced the attention of each time to its periodic correlated time.

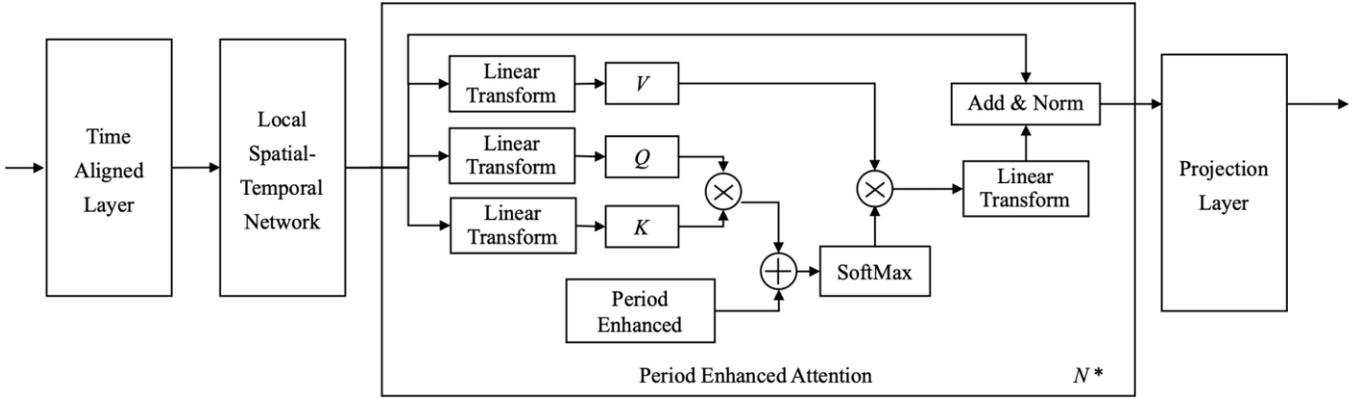


Figure 2. The overall architecture of TTST

- We propose a transformer called TSTT that combines LSTN and PEA for traffic prediction. We conducted experiments based on two real-world traffic datasets, and the results show that TTST outperforms baselines.

## II. RELATED WORK

Early researches treated traffic prediction as a time series forecasting problem, only focusing on capturing the temporal dependency in traffic data. In past decade, due to the development of Convolution Neural Network (CNN) in area of computer vision, some researches such as ConvLSTM [11] and PredRNN [12], have applied CNN to catch the Euclidean spatial dependency. In recent years, with the development of Graph Neural Network, especially GCN [13], many researches have increasingly integrated TLN with GCN to effectively capture the spatial-temporal dependency in traffic data.

TGCN [1] stacks GCN and GRU to separately capture time dependency and spatial dependency. STGCN [2] combines one-dimensional convolution and gate mechanism for time dependency, and stacks it with GCN to capture spatial dependency. ASTGCN [3] extracts adjacent, daily, and weekly time segments, computes temporal and spatial attention for each of these three groups of data, and then combines GCN with one-dimensional convolution to capture the spatial-temporal dependency. GWN [4] proposes an adaptive adjacency matrix learning approach that combines GCN to capture spatial dependency while utilizing gated TCN to capture temporal dependency. DCRNN [5] introduces a bidirectional random walk graph diffusion convolutional recurrent network and employs seq2seq architecture for modeling spatial-temporal dependency. DGCRN [6] combines GCN with RNN and uses a hyper-network to generate an dynamic graph before each recurrent of RNN to capture dynamic changes of traffic data. Based on the generated dynamic graph and the original adjacent matrix, GCN can capture more spatial information.

## III. PRELIMINARIES

### A. Traffic data

We define the traffic data as  $D = (V, A, X_t)$ . Here,  $V$  is a set with  $N = |V|$  vertices, where each  $v \in V$  represents a traffic sensor or road segment in the traffic network;  $A$  is the adjacent matrix, where each  $a_{ij} \in A$  represents a connection from  $v_i$  to  $v_j$

in the traffic network when it equals 1;  $X_t = (x_t^1, x_t^2, \dots, x_t^N)$  is the traffic state, where  $x_t^i$  represents the state collecting from  $v_i$  at time  $t$ .

### B. Problem formulation

The traffic prediction problem can be formulated as follows, given the historical data  $X = (X_{t-h}, X_{t-h-1}, \dots, X_t)$  of the past  $h$  time steps, predict the traffic state  $Y = (Y_{t+1}, Y_{t+2}, \dots, Y_{t+f})$  of the  $f$  future time steps.

## IV. METHODS

### A. Overview

Fig. 2 shows the overall architecture of Traffic Spatial-Temporal Transformer (TSTT) including two main parts: Local Spatial-Temporal Network (LSTN) and Period Enhanced Attention (PEA). TSTT first applies a linear transformation layer to align the time dimension of the input from  $h$  to  $h+f$ , which can better learn the temporal variation of cross past and future. Then the aligned input is fed into LSTN to capture local spatial-temporal dependency. Next, the output of LSTN is further processed through PEA to capture long-term time dependency. Finally, PEA's output is projected from the time dimension of  $h+f$  to  $f$  using a linear transformation layer, and the projected output serves as the prediction result of TTST.

### B. LSTN

We propose LSTN by incorporating dynamic graph generation and GCN into S-LSTM. Fig. 3 shows the structure of LSTN. The hidden state  $H^k$  in  $k$ -th recurrence include the sub-hidden time state  $h_i^k$  at each time  $i$  and the global spatial state  $z^k$ , as in (1). Following S-LSTM, we initialize  $h_i^0$  and  $z^0$  in  $H^0$  with the same parameter  $h^0$ . The transition from  $H^{k-1}$  to  $H^k$  involves the generation of dynamic graph  $g^k$  and the updates of  $h_i^k$  and  $z^k$ .

$$H^k = \{h_i^k, \dots, h_i^k, z^k\} \quad (1)$$

Building upon the S-LSTM's process of extracting local temporal dependency based on each time state and its temporal context, LSTN integrates graph convolution into this process to capture local spatial-temporal dependency. GCN is typically based on the adjacency matrix in traffic data, which only records direct connections between roads. However, there are rich relationships between indirectly connected roads. To comprehensively consider the relationships between different



in a time series. To further capture the long-term time dependency in traffic data and account for its periodicity, we propose PEA by enhancing the attention scores for related periodic time points within Self-Attention.

PEA consists of two parts: the period enhanced matrix generation and the self-attention. Fig. 4 shows the process of period enhanced matrix generation part, it starts by using Fast Fourier Transform (FFT) to obtain frequency intensity of the

data, computing the average frequency intensity for each channel and selecting the top- $d$  highest frequency intensity, as shown in (6), where  $f_d$  is the  $d$ -th highest frequency intensity.

$$\{f_1, \dots, f_d\} = \text{Top}_d(\text{Avg}(\text{FFT}(X))) \quad (6)$$

Then the periods are computed based on the historical data length  $T$ , as in (7), and the enhanced coefficients are obtained by applying LeakyReLU operation to frequency intensity, as in

TABLE I. THE RESULTS OF EXPERIMENT

Data	Model	3 steps (15 mins)			6 steps (30 mins)			12 steps (60 mins)		
		RMSE	MAE	MAPE	RMSE	MAE	MAPE	RMSE	MAE	MAPE
PeMS - Bay	STGCN	2.962	1.352	2.94%	3.942	1.728	3.94%	4.905	2.11	5.12%
	ASTGCN	3.217	1.396	3.08%	4.286	1.814	4.52%	4.984	2.213	5.15%
	GWN	2.812	1.349	2.81%	3.795	1.692	3.86%	4.618	2.025	4.79%
	DCRNN	2.854	1.372	2.87%	3.892	1.765	4.12%	4.811	2.143	5.09%
	DGCRN	<u>2.69</u>	<u>1.28</u>	<u>2.66%</u>	<u>3.63</u>	<u>1.65</u>	<u>3.55%</u>	<u>4.42</u>	<u>1.89</u>	<u>4.43%</u>
	TSTT	<b>2.651</b>	<b>1.275</b>	<b>2.64%</b>	<b>3.592</b>	<b>1.623</b>	<b>3.53%</b>	<b>4.218</b>	<b>1.756</b>	<b>4.39%</b>
METR - LA	STGCN	5.236	3.052	7.12%	6.212	3.476	8.67%	8.142	4.014	10.52%
	ASTGCN	5.512	3.127	7.21%	6.275	3.631	8.63%	8.256	4.132	10.66%
	GWN	5.374	2.935	7.02%	6.174	3.205	8.39%	7.933	3.768	10.34%
	DCRNN	5.059	2.962	7.08%	6.193	3.355	8.48%	8.122	3.855	10.39%
	DGCRN	<u>5.01</u>	<u>2.62</u>	<u>6.63%</u>	<u>6.05</u>	<u>2.99</u>	<u>8.02%</u>	<u>7.19</u>	<u>3.44</u>	<u>9.73%</u>
	TSTT	<b>4.965</b>	<b>2.596</b>	<b>6.57%</b>	<b>5.832</b>	<b>2.895</b>	<b>7.93%</b>	<b>7.094</b>	<b>3.245</b>	<b>9.45%</b>

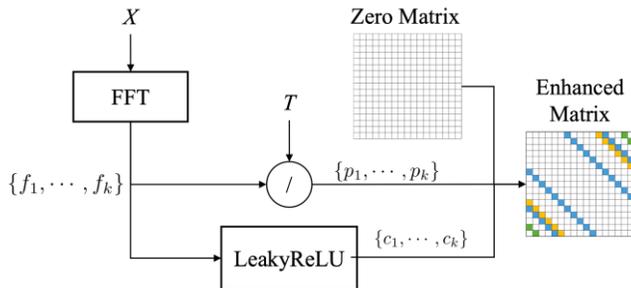


Figure 4. The process of period enhanced matrix generation.

(8), where  $p_d$  and  $c_d$  represents the period and enhanced coefficient for  $f_d$ .

$$\{p_1, \dots, p_d\} = \{T/f_1, \dots, T/f_d\} \quad (7)$$

$$\{c_1, \dots, c_d\} = \text{LeakyReLU}(\{f_1, \dots, f_d\}) \quad (8)$$

The enhanced matrix is generated by updating a zero matrix  $M$  based on the period and frequency intensity. The update process for each  $m_{ij} \in M$  is as shown in (9). The detail of self-attention is as shown in (10) and (11), where  $W_q$ ,  $W_k$  and  $W_v$  are the parameters of linear transformations,  $d_k$  is the dimension of key.

$$m_{ij} = \{c_k : (i - j) \bmod p_k = 0\} \quad (9)$$

$$Q, K, V = XW_q, XW_k, XW_v \quad (10)$$

$$\text{Attention}(Q, K, V) = \text{softmax}((QK^T + M) / (d_k)^{-1/2}) \quad (11)$$

## V. Experiment

### A. Setup

We employed two real-world datasets for conducting the experiment: METR-LA[7] and PeMS-Bay[15], and adopt the five following baselines to compare with our model: STGCN[4], ASTGCN[5], GWN[6], DCRNN[7], and DGCRN[8]. The following three metrics are used to evaluate the models:

- Mean Absolute Error (MAE):

$$\text{MAE} = (\sum_{i=1}^m |\hat{y}_i - y_i|) / m$$

- Root Mean Squared Error (RMSE):

$$\text{RMSE} = ((\sum_{i=1}^m (\hat{y}_i - y_i)^2) / m)^{-1/2}$$

- Mean Absolute Percent Error (MAPE):

$$\text{MAPE} = 100\% * (\sum_{i=1}^m |\hat{y}_i - y_i| / y_i) / m$$

where  $m$  is the number of samples in the test dataset,  $\hat{y}_i$  and  $y_i$  are the predicted and true values of sample  $i$ .

### B. Result

Table 1 shows the experimental results. The bold results indicate the best performance, and the underlined results represent the second-best. Overall, TTST performs better than other models in the three evaluated metrics. To be more specific, the advantage of TTST is more evident in long-term

predictions (12 steps) compared to short-term predictions (3 steps and 6 steps).

## VI. CONCLUSIONS

In this paper, we fully consider the spatial-temporal dependency in traffic data and propose a transformer called TTST to achieve traffic prediction. TTST includes a recurrent neural network combined with GCN to capture the local spatial-temporal dependency, and a period enhanced attention to capture the long-range temporal dependency. We have tested TTST on two real-world datasets for traffic prediction task and the results showed TTST outperforms baselines with a greater advantage in long-term prediction compared to short-term prediction.

## ACKNOWLEDGMENT

This work is supported by National Natural Science Foundation of China (U19A2067, 61976051).

## REFERENCES

- [1] M. Shaygan, C. Meese, W. Li, X.G. Zhao, and M. Nejad, "Traffic prediction using artificial intelligence: review of recent advances and emerging opportunities," *Transportation Research Part C: Emerging Technologies*, vol. 145, p. 103921, 2022.
- [2] G. Jin, Y. Liang, Y. Fang, J. Huang, J. Zhang, and Y. Zheng, "Spatio-temporal graph neural networks for predictive learning in urban computing: A survey," arXiv preprint arXiv:2303.14483, 2023.
- [3] L. Zhao, Y. Song, C. Zhang, Y. Liu, P. Wang, T. Lin, M. Deng, and H. Li, "T-GCN: A Temporal Graph Convolutional Network for Traffic Prediction," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 9, pp. 3848-3858, 2019.
- [4] B. Yu, H. Yin, and Z. Zhu, "Spatio-temporal graph convolutional networks: A deep learning framework for traffic forecasting," arXiv preprint arXiv:1709.04875, 2017.
- [5] S. Guo, Y. Lin, N. Feng, C. Song, and H. Wan, "Attention based spatial-temporal graph convolutional networks for traffic flow forecasting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 922-929, 2019.
- [6] Z. Wu, S. Pan, G. Long, J. Jiang, and C. Zhang, "Graph Wavenet for Deep Spatial-Temporal Graph Modeling," arXiv preprint arXiv:1906.00121, 2019.
- [7] Y. Li, R. Yu, C. Shahabi, and Y. Liu, "Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting," arXiv preprint arXiv:1707.01926, 2017.
- [8] F. Li, J. Feng, H. Yan, G. Jin, F. Yang, F. Sun, D. Jin, and Y. Li, "Dynamic Graph Convolutional Recurrent Network for Traffic Prediction: Benchmark and Solution," *ACM Transactions on Knowledge Discovery from Data*, vol. 17, no. 1, pp. 1-21, 2023.
- [9] M. Schuster and K.K. Paliwal, "Bidirectional Recurrent Neural Networks," *IEEE Transactions on Signal Processing*, vol. 45, no. 11, pp. 2673-2681, 1997.
- [10] Y. Zhang, Q. Liu, and L. Song, "Sentence-State LSTM for Text Representation," arXiv preprint arXiv:1805.02474, 2018.
- [11] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo, "Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting," in *Advances in Neural Information Processing Systems (NIPS)*, vol. 28, 2015.
- [12] Y. Wang, M. Long, J. Wang, Z. Gao, and P. S. Yu, "PredRNN: Recurrent Neural Networks for Predictive Learning Using Spatiotemporal LSTMs," in *Proceedings of the Neural Information Processing Systems (NIPS)*, vol. 30, 2017.
- [13] T.N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," arXiv preprint arXiv:1609.02907, 2016.
- [14] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7132-7141, 2018.
- [15] R. Jiang, D. Yin, Z. Wang, Y. Wang, J. Deng, H. Liu, Z. Cai, J. Deng, X. Song, and R. Shibasaki, "DL-Traff: Survey and Benchmark of Deep Learning Models for Urban Traffic Prediction," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 4515-4525, October 2021.