

Supplemental Material to

LIGHTCTS: A Lightweight Framework for Correlated Time Series Forecasting

1 TIME AND SPACE COMPLEXITY ANALYSIS OF S/T-OPERATORS

1.1 Key Notation

First, we recap the key notation related to *correlated time series* (CTS) in Table 1.

Table 1. Notation

Notation	Description
X	An indexed set of correlated time series (CTS)
H	Embedded CTS feature maps
N	Number of time series in X
T	Number of time steps in X
F	Number of features in X
D	Embedding size of S/T-operators
P	Number of historical time steps used in CTS forecasting

Notably, we suppose that the input CTS data $X \in \mathbb{R}^{N \times P \times F}$ has been embedded into $H \in \mathbb{R}^{N \times P \times D}$ by the input and embedding module. Consequently, the input shape of S/T-operators is $(N \times P \times D)$.

The detailed complexity analyses for T-operators and S-operators are presented in Sections 1.2 and 1.3, respectively.

1.2 T-operator

1.2.1 CNN family. In this section, we present the detailed compositions of T-operators from the CNN family, and specify their time and space complexities. We take the TCN (Temporal Convolutional Network) as the representative [2]. The main difference between the TCN and the standard CNN is the convolution order such that the TCN should ensure the convolution output is based on the information from previous time steps. Nevertheless, the convolution order does not influence the time and space complexities.

With the raw input CTS $X \in \mathbb{R}^{N \times P \times F}$ being embedded into the latent representation $H \in \mathbb{R}^{N \times P \times D}$ by the embedding module, a TCN layer cast on H is formalized as follows.

$$\text{TCN}(H \mid \delta, K) = H', \text{ where}$$

$$H'[i; t; d] = \sum_{k=0}^{K-1} (H[i; t - \delta \times k; :] \cdot W^d[k; :])$$

is the d -th ($d \in [0, D)$) output feature map at time step t ($t \in [0, P)$) of the i -th ($i \in [0, N)$) time series, $W^d \in \mathbb{R}^{D \times K}$ is the d -th convolutional kernel, and K (often as small as 2 or 3) and δ are the kernel size and dilation rate, respectively.

Specifically, we adopt D kernels, each of which is in size of $(K \times D)$. Notably, the kernel size K is usually a small value (e.g., 2 or 3) that can be regarded as a constant.

Specifically, the input is of shape $N \times P \times D$ and the kernel size is set to K , thus each kernel weight matrix is of shape $1 \times K \times D$. In addition, the shape of the output is also $N \times P \times D$, which means there are D kernels. Since all parameters of TCNs come from convolution kernels, the **space complexity** is $O(KD^2)$ (i.e., D kernels each taking $1 \times K \times D$ parameters). Notably, the kernel size K is usually a small value (e.g., 2 or 3) that can be regarded as a constant, the space complexity can be simplified as $O(D^2)$. It should be noticed that, since features from different time steps and nodes typically share one set of convolution kernels [13, 14, 16], the space complexity is not related to N and P .

Besides, for one kernel in one convolution calculation, its time complexity is $O(D)$ since the number of its input channels (i.e., embedding size) is D^1 . With P convolution operations along the temporal dimension for N CTS nodes using D kernels, the overall **time complexity** is $O(D^2 \cdot N \cdot P)$.

1.2.2 RNN family. We proceed to detail the compositions of RNN-based T-operators and derive their time and space complexities. There are mainly two RNN-based T-operators, namely LSTMs (Long Short-Term Memory) [8] and GRUs (Gate Recurrent Units) [5]. Their main difference is the number of gates, which is a constant that does not affect the complexities. In the following, we take the LSTM as the example.

Equation 1 presents the workflow of an LSTM operator in a particular time step t :

$$\begin{aligned} i_t &= \sigma(W_{ii}H(i)_t + W_{hi}h_{t-1}), \\ f_t &= \sigma(W_{if}H(i)_t + W_{hf}h_{t-1}), \\ g_t &= \tanh(W_{ig}H(i)_t + W_{hg}h_{t-1}), \\ o_t &= \sigma(W_{io}H(i)_t + W_{ho}h_{t-1}), \\ c_t &= f_t \odot c_{t-1} + i_t \odot g_t, \\ h_t &= o_t \odot \tanh(c_t), \end{aligned} \quad (1)$$

where $H(i)_t$ is the i -th time series features at the t -th time step of H , h_t is the hidden state at the t -th time step of LSTM, i_t , f_t , g_t , and o_t denote the outputs of the input gate, the forget gate, the cell update, and the output gate, respectively. The computation of the four components i_t , f_t , g_t , and o_t are the same in terms of time complexity. Thus, we take the input gate i_t as an example for time complexity analysis:

$$i_t = \sigma(W_{ii}X(i)_t + W_{hi}h_{t-1}) \quad (2)$$

where $W_{ii} \in \mathbb{R}^{H_h \times D}$ and $W_{hi} \in \mathbb{R}^{H_h \times H_h}$ are the weight matrices of the input gate, which determine the value of the input, H_h is the size of the hidden state.

Since the LSTM is made up of these gates which consist of two weight matrices with shapes of $H_h \times D$ and $H_h \times H_h$, the space complexity of an LSTM T-operator is $O(H_h D + H_h^2)$. For simplification, we substitute H_h with D as they are normally similar in values for practical implementations. As a result, the **space complexity** is simplified as $O(D^2)$.

Similarly, since features from different time steps and nodes typically share one set of gates [1, 3, 4, 9], the space complexity is not related to N and P .

On the other hand, Equation 2 consists of the computations of $W_{ii}H(i)_t$ and $W_{hi}h_{t-1}$. Thus, the time complexity of Equation 2 is $O(H_h(H_h + D))$, simplified as $O(D^2)$. Considering that we have N time series with the length of P , the overall **time complexity** of an LSTM operator is $O(D^2 NP)$.

1.2.3 Transformer family. A standard Transformer-based T-operator consists of an MHA layer and an FFN layer [11].

In particular, an **MHA layer** is a concatenation of a number h of repeated self-attention modules (heads) in parallel:

$$\text{MHA}(H^{\text{PE}}) = \text{concat}(\{\text{head}_i(H^{\text{PE}})\}_{i=1}^h) \quad (3)$$

$$\text{head}_i(H^{\text{PE}}) = \text{softmax}(H^{\text{I}}) \cdot V_i \quad (4)$$

$$H^{\text{I}} = (Q_i \cdot K_i^{\text{T}}) / \sqrt{D/h}, \quad (5)$$

$$Q_i = H^{\text{PE}} \cdot W_i^Q, K_i = H^{\text{PE}} \cdot W_i^K, V_i = H^{\text{PE}} \cdot W_i^V. \quad (6)$$

¹As explained above, we also disregard the kernel size K in the time complexity analysis.

where all learnable matrices $W_i^* \in \mathbb{R}^{D \times \frac{D}{h}}$ are used to convert the embedding size from D to D/h . In Equation 4, each head $\text{head}_i(\cdot)$ is a self-attention module that produces attention scores for its input $H^{\text{PE}} \in \mathbb{R}^{N \times D}$ by making the input interact with itself. It should be noticed that, H^{PE} is the position embedding results of the input CTS feature H , of which computation only involves matrix addition calculations, not contributing to the time complexity.

An **FFN layer** merges the output $H^{\text{MHA}} \in \mathbb{R}^{P \times D}$ of Equation 3 and provides non-linear activations via two fully-connected layers:

$$\text{FFN}(H^{\text{MHA}}) = \text{FFN}^{(1)}(H^{\text{MHA}}) \cdot W_2 + b_2, \quad (7)$$

$$\text{FFN}^{(1)}(H^{\text{MHA}}) = \text{ReLU}(H^{\text{MHA}} \cdot W_1 + b_1), \quad (8)$$

where $W_1 \in \mathbb{R}^{D \times D'}$, $W_2 \in \mathbb{R}^{D' \times D}$, $b_1 \in \mathbb{R}^{D'}$, and $b_2 \in \mathbb{R}^D$ are learnable matrices and biases. The two layers first enlarge the embedding size from D to D' , and then scale it back to D .

For one Transformer-based T-operator, it has $3 \times h$ linear projection matrices $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{D \times D/h}$ in the MHA, and two FFN layers with the weight matrices of $\mathbb{R}^{D \times D_F}$ and $\mathbb{R}^{D_F \times D}$, respectively, D_F is the embedding size of the first fully-connected layer of the FFN layer. Therefore, the overall **space complexity** is positively correlated with $\mathcal{O}(3 \times h \times D \times D/h + 2 \times D \times D_F)$, which equals to $\mathcal{O}(D^2)$ (D_F is typically a quadruple of D [11], thus we substitute D_F with D). It should be also noticed that, as the most studies do [6, 10, 12, 15, 17], the features from all time steps and nodes share one set of Transformer parameters, thus, the parameter number is unrelated to N and P .

As for the time complexity, considering $3 \times h$ linear projection matrices with the shape of $\mathbb{R}^{D \times D/h}$ in the MHA, and we need to exploit the temporal features of N nodes, the time complexity of this part is positively correlated with $(3 \times N \times h \times P \times D \times D/h)$, which can be simplified as $\mathcal{O}(D^2NP)$. Since we conduct self-attention operations on the basis of time steps, according to the computation processes of Equation 4 and Equation 5, the time complexity of this part is $\mathcal{O}(N \times h \times P \times P \times D/h + N \times h \times P \times P \times D/h)$, which can be simplified as $\mathcal{O}(P^2DN)$. For the FFN layer, since it consists of two linear layers with D_F and D kernels, the time complexity of this part is $\mathcal{O}(2 \times N \times P \times D \times D_F)$, which is $\mathcal{O}(D^2NP)$. In summary, the total **time complexity** is $\mathcal{O}(D^2NP + P^2DN)$, which equals to $\mathcal{O}(NDP(P + D))$.

1.3 S-operator

1.3.1 GCN family. Each GCN S-operator performs the following two steps: neighborhood information aggregation and nonlinear transformation. The neighborhood information aggregation is performed by the matrix multiplication between the spatial adjacency matrix and the input node features, followed by a nonlinear transformation performed by convolution operation over the aggregated feature. In particular, a GCN S-operator conducted one GCN operation on one time step slice $H_G \in \mathbb{R}^{N \times D}$ is shown as follows:

$$\text{GCN}(H_G) = \sigma(A \cdot H_G \cdot W) \quad (9)$$

where $A \in \mathbb{R}^{N \times N}$ is a normalized adjacency matrix representing the spatial correlations, $W \in \mathbb{R}^{D \times D}$ is the GCN weight matrix, and $\sigma(\cdot)$ is the nonlinear activation function.

The **space complexity** derives from the GCN weight matrix $W \in \mathbb{R}^{D \times D}$, which is $\mathcal{O}(D^2)$, and is also unrelated to N and P . This is because the features from all time steps and nodes also share the same GCN weight matrix [7, 9, 12, 14].

Besides, from Equation 9, the time complexity for one time step slice is calculated as $\mathcal{O}(N^2D + ND^2)$. Considering P time steps, the **time complexity** can be expressed as $\mathcal{O}(NDP(N + D))$.

1.3.2 Transformer family. Similar to Transformer-based T-operators, a standard Transformer-based S-operator also consists of an MHA layer and an FFN layer and adopts the same computation mechanism, except for the slicing method of the input features. For S-operators, since we focus

on exploiting the correlations among the information of each node within P time steps, the input features for the self-attention are partitioned into N slices of $P \times D$ features, which are P slices of $N \times D$ features in Transformer-based T-operators.

For one Transformer-based S-operator, it also has $3 \times h$ linear projection matrices $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{D \times D/h}$ in the MHA, and two FFN layers with the weight matrices of $\mathbb{R}^{D \times D_F}$ and $\mathbb{R}^{D_F \times D}$, respectively. Therefore, the overall **space complexity** is also $O(3 \times h \times D \times D/h + 2 \times D \times D_F)$, which equals to $O(D^2)$. Similarly, as most studies do [6, 10, 12, 15, 17], the features from all time steps and nodes share one set of Transformer parameters, thus, the parameter number is unrelated to N and P .

As for the time complexity, considering $3 \times h$ linear projection matrices with the shape of $\mathbb{R}^{D \times D/h}$ in the MHA, and we need to exploit the spatial correlations of P time steps, the time complexity of this part is positively correlated with $(3 \times P \times h \times N \times D \times D/h)$, which can be simplified as $O(D^2NP)$. Since we conduct self-attention operations on the basis of nodes, like the computation processes of Equation 4 and Equation 5, the time complexity of this part is $O(P \times h \times N \times N \times D/h + P \times h \times N \times N \times D/h)$, which can be simplified as $O(N^2DP)$. For the FFN layer, since it consists of two linear layers with D_F and D kernels, the time complexity of this part is positively correlated with $(2 \times N \times P \times D \times D_F)$, which is $O(D^2NP)$.

Overall, the total **time complexity** is $O(D^2NP + N^2DP)$, which equals to $O(NDP(N + D))$.

2 PARAMETER STUDIES ON THE OTHER DATASETS

The experiments in this section are the supplement to those in Section 5.3.1 to Section 5.3.3 (on Page 9 to 11) in the submission.

2.1 Multi-step Forecasting

We perform parameter studies on the other three multi-step datasets, i.e. PEMS04, METR-LA, and PEMS-BAY. The results are shown in Fig. 1, Fig. 2, and Fig. 3.

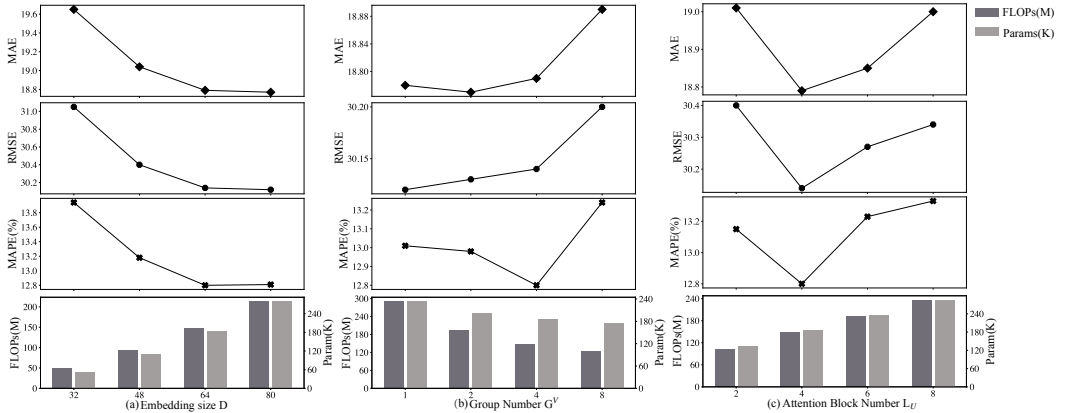


Fig. 1. Impact of (a) embedding size D , (b) group number G^T , and (c) attention block number L_S in GL-Former for multi-step forecasting on the PEMS04 dataset.

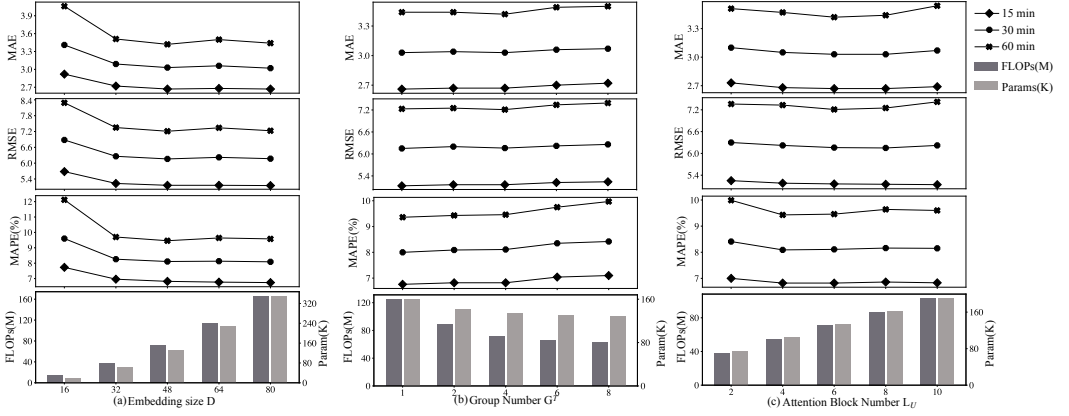


Fig. 2. Impact of (a) embedding size D , (b) group number G^T , and (c) attention block number L_U in GL-Former for multi-step forecasting on the METR-LA dataset.

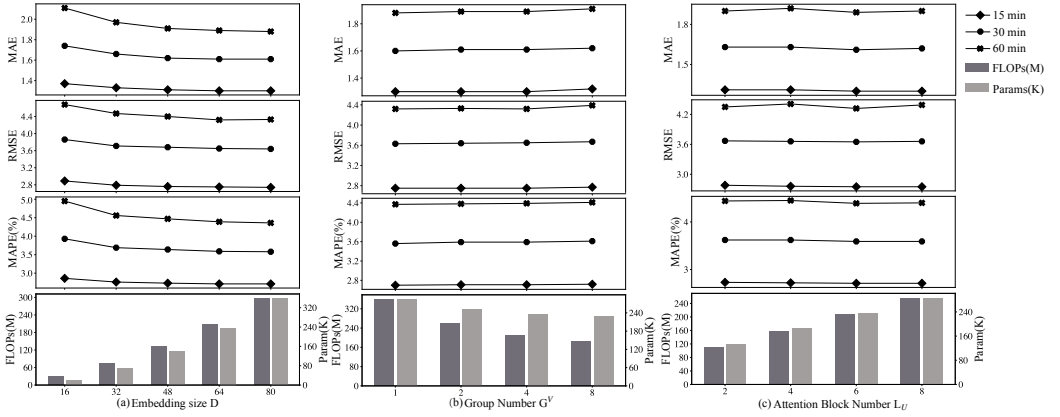


Fig. 3. Impact of (a) embedding size D , (b) group number G^T , and (c) attention block number L_U in GL-Former for multi-step forecasting on the PEMS-BAY dataset.

2.2 Single-step Forecasting

We also perform parameter studies on the other single-step dataset, i.e. Electricity. The result is shown in Fig. 4.

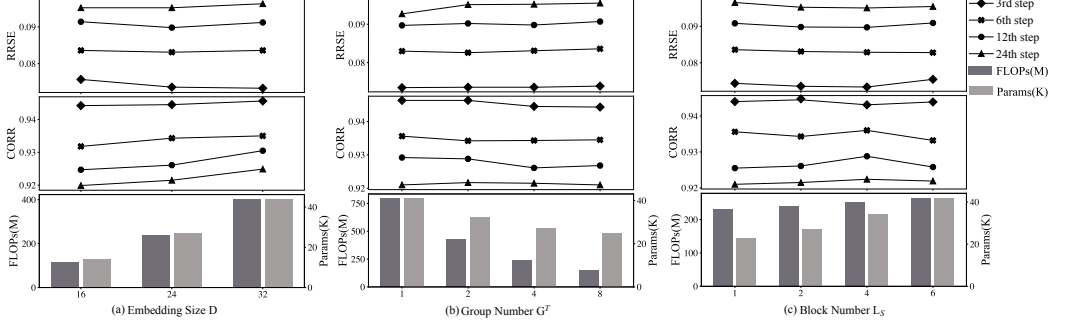


Fig. 4. Impact of (a) embedding size D, (b) group number G^T , and (c) attention block number L_S in GL-Former for single-step forecasting on the Electricity dataset.

3 EXPERIMENTS OF THE IMPACTS OF G^M AND G^F

The experiments in this section are the supplement to those in Section 5.3.2 (on Page 10 to 11) in the submission.

3.1 Multi-step Forecasting

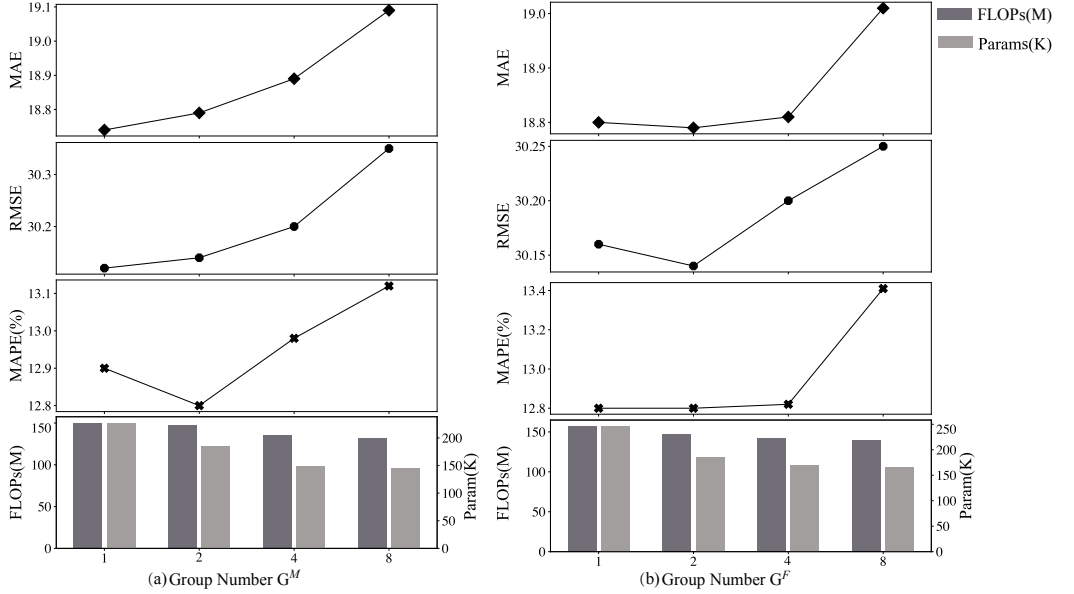


Fig. 5. Impact of (a) the group number of MHA layers G^M and (b) the group number of FFN layers G^F for multi-step forecasting on the PEMS04 dataset.

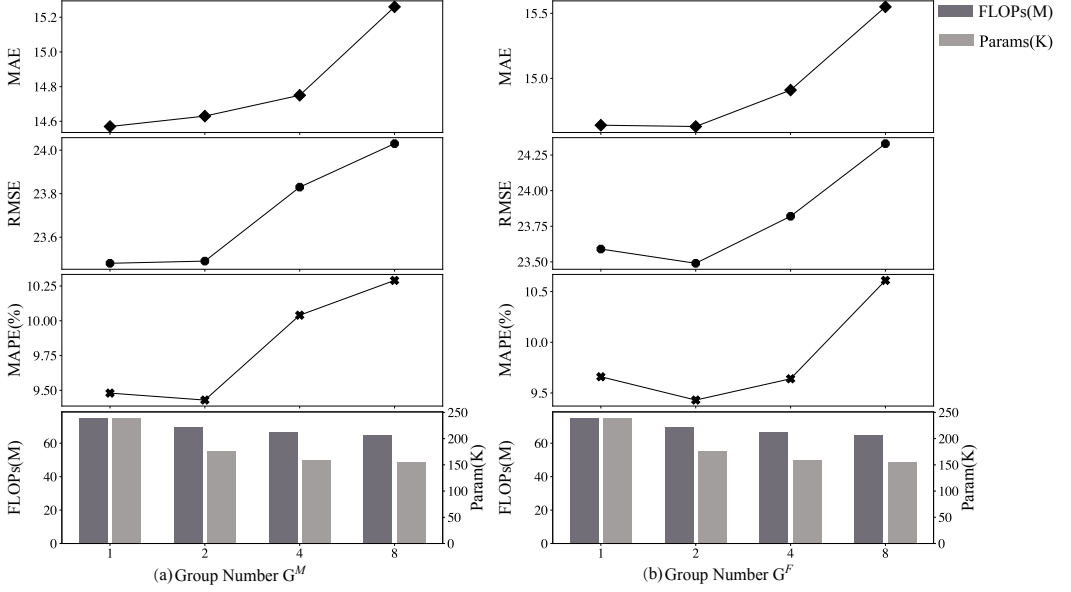


Fig. 6. Impact of (a) the group number of MHA layers G^M and (b) the group number of FFN layers G^F for multi-step forecasting on the PEMS08 dataset.

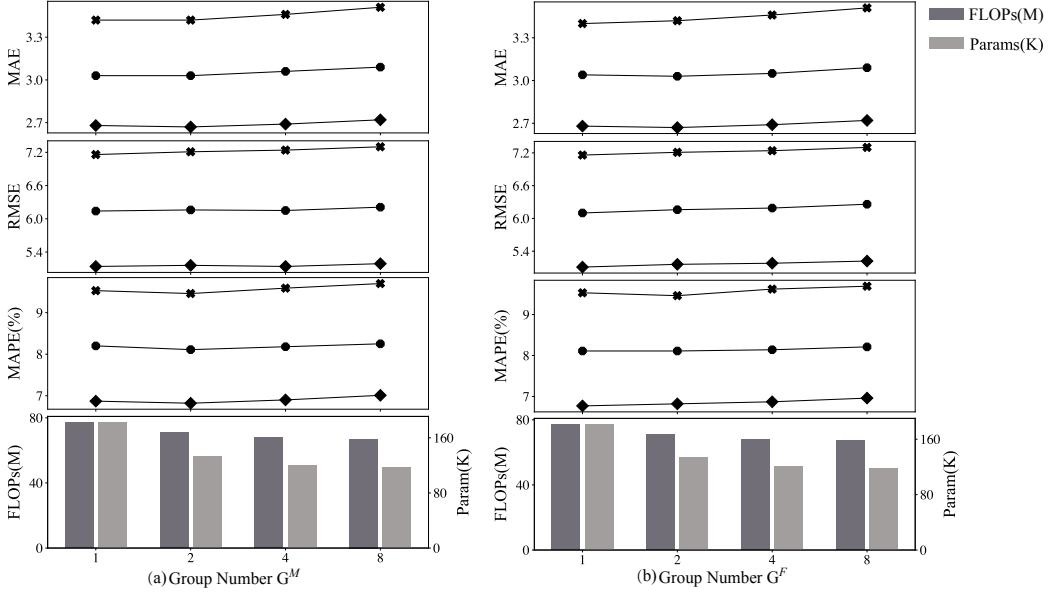


Fig. 7. Impact of (a) the group number of MHA layers G^M and (b) the group number of FFN layers G^F for multi-step forecasting on the METR-LA dataset.

3.2 Single-step Forecasting

4 ABLATION STUDIES ON THE OTHER DATASETS

The experiments in this section are the supplement to those in Section 5.4 (on Page 11 to 12) in the submission.

REFERENCES

- [1] Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. Adaptive graph convolutional recurrent network for traffic forecasting. In *NeurIPS*, pages 17804–17815, 2020.
- [2] Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv:1803.01271*, 2018.
- [3] Yen-Yu Chang, Fan-Yun Sun, Yueh-Hua Wu, and Shou-De Lin. A memory-network based solution for multivariate time-series forecasting. *arXiv preprint arXiv:1809.02105*, 2018.
- [4] Weiqi Chen, Ling Chen, Yu Xie, Wei Cao, Yusong Gao, and Xiaojie Feng. Multi-range attentive bicomponent graph convolutional network for traffic forecasting. In *AAAI*, pages 3529–3536, 2020.
- [5] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [6] Jake Grigsby, Zhe Wang, and Yanjun Qi. Long-range transformers for dynamic spatiotemporal forecasting. *arXiv preprint arXiv:2109.12218*, 2021.
- [7] Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. In *AAAI*, pages 922–929, 2019.
- [8] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [9] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *ICLR*, pages 1–10, 2018.
- [10] Cheonbok Park, Chunggi Lee, Hyojin Bahng, Yunwon Tae, Seungmin Jin, Kihwan Kim, Sungahn Ko, and Jaegul Choo. ST-GRAT: A novel spatio-temporal graph attention networks for accurately forecasting dynamically changing road speed. In *CIKM*, pages 1215–1224, 2020.
- [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, volume 30, 2017.
- [12] Xinle Wu, Dalin Zhang, Chenjuan Guo, Chaoyang He, Bin Yang, and Christian S Jensen. AutoCTS: Automated correlated time series forecasting. *PVLDB*, 15(4):971–983, 2021.
- [13] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang. Connecting the dots: Multivariate time series forecasting with graph neural networks. In *KDD*, pages 753–763, 2020.
- [14] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, and Chengqi Zhang. Graph wavenet for deep spatial-temporal graph modeling. In *IJCAI*, pages 1907–1913, 2019.
- [15] Mingxing Xu, Wenrui Dai, Chunmiao Liu, Xing Gao, Weiyao Lin, Guo-Jun Qi, and Hongkai Xiong. Spatial-temporal transformer networks for traffic flow forecasting. *arXiv preprint arXiv:2001.02908*, 2020.
- [16] Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: a deep learning framework for traffic forecasting. In *IJCAI*, pages 3634–3640, 2018.
- [17] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *AAAI*, pages 11106–11115, 2021.