

---

# Preformer: Predictive Transformer with Multi-Scale Segment-wise Correlations for Long-Term Time Series Forecasting

---

Dazhao Du<sup>1,2</sup> Bing Su<sup>3</sup> Zhewei Wei<sup>3</sup>

## Abstract

Transformer-based methods have shown great potential in long-term time series forecasting. However, most of these methods adopt the standard point-wise self-attention mechanism, which not only becomes intractable for long-term forecasting since its complexity increases quadratically with the length of time series, but also cannot explicitly capture the predictive dependencies from contexts since the corresponding key and value are transformed from the same point. This paper proposes a predictive Transformer-based model called *Preformer*. Preformer introduces a novel efficient *Multi-Scale Segment-Correlation* mechanism that divides time series into segments and utilizes segment-wise correlation-based attention for encoding time series. A multi-scale structure is developed to aggregate dependencies at different temporal scales and facilitate the selection of segment length. Preformer further designs a predictive paradigm for decoding, where the key and value come from two successive segments rather than the same segment. In this way, if a key segment has a high correlation score with the query segment, its successive segment contributes more to the prediction of the query segment. Extensive experiments demonstrate that our Preformer outperforms other Transformer-based methods.

(Borovykh et al., 2017; Sen et al., 2019) and Transformer-based models (Vaswani et al., 2017; Wu et al., 2020). RNN-based methods suffer from the problem of gradient vanishing, gradient exploding, and lack of parallelism. TCN-based methods need deeper layers to achieve larger local receptive fields. For both categories of methods, signals must pass through a long path between two far-away temporal locations, hence the number of operations required to associate two elements increases with their temporal distance. Differently, transformer-based methods directly model the relationships between any element pairs and can better capture long-term dependencies, which is crucial for long-term forecasting.

On the other hand, the standard self-attention mechanism of Transformer calculates all similarities between any element pairs, where the computational and space complexities increase quadratically with the length of the time series. Recent works (Li et al., 2019; Kitaev et al., 2019; Zhou et al., 2021) explore different sparse attention mechanisms to suppress the contribution of irrelevant time steps and ease the computational pressure. These models still perform dot-product attention to time steps individually and utilize the point-wise connections to capture temporal dependencies. However, a single point may have limited influence on predicting the future. Autoformer (Wu et al., 2021) conducts the series-wise dependencies discovery by performing Auto-Correlation of the time series to the top-k time delayed series. The aggregation operation acts on the whole delayed series and complicated Fourier transforms are required for Auto-Correlation.

These methods perform the correlation either at the point level or at the overall series level, which not only require high computational redundancy to intensively tackle point pairs or perform time-frequency domain transformations, but also do not directly reflect the true dependencies within the time series. For instance, in traffic flow forecasting, the flows of a single previous time point and the shifted whole time series may have limited contribution on the flow of the future period. There exist stronger correlations at the segment level, e.g., the flow at evening today should be more related to the flow at evening yesterday than the flow of other time periods any day and the flow of a period that

## 1. Introduction

Long-term time series forecasting has a wide range of real-world applications such as financial investment, traffic management, and electricity management. Existing deep learning-based methods can be divided into three categories, i.e., RNN-based models (Qin et al., 2017; Lai et al., 2018; Song et al., 2018; Salinas et al., 2020), TCN-based models

---

<sup>1</sup>Institute of Software Chinese Academy of Sciences, Beijing, China <sup>2</sup>University of Chinese Academy of Sciences, Beijing, China <sup>3</sup>Renmin University of China, Beijing, China. Correspondence to: Bing Su <subingats@gmail.com>.

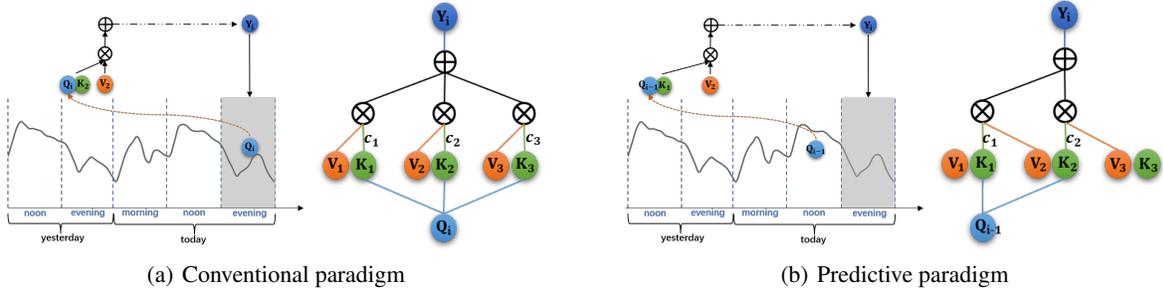


Figure 1. (a) Segment-Correlation with the conventional paradigm. (b) Segment-Correlation with the predictive paradigm. The blue, green, and orange nodes represent the queries, keys and values generated by different segments (different time periods), respectively. The gray part in the time series is the segment to be predicted. Multiplication and addition represent weighted aggregation operations.  $c_j$  is the correlation score between  $Q_i$  and  $K_j$ . In the predictive paradigm, to obtain the output  $Y_i$  at the current segment  $i$ , the similarity between  $Q_{i-1}$  at the previous segment  $i - 1$  and the key  $K_j$  is used as the weight of the value  $V_{j+1}$  in the aggregation operation.

has passed a long time ago.

To this end, we propose a novel sparse attention mechanism called *Multi-Scale Segment-Correlation (MSSC)*. Segment-wise correlation not only reduces the amount of calculation since the number of segments is much smaller than the number of points, but also better explore the locality of neighboring points. The length of segment is a critical hyperparameter. Long segments ignore fine-grained information while short segments have high computational redundancy. To tackle this issue, MSSC performs correlation calculations and fusion on multiple segment lengths, i.e., multi-scale resolutions, while maintaining low complexity.

Standard attention paradigm can be applied to MSSC for time series encoding, but it is not well suited for forecasting since the unknown prediction segment generates the query by itself to predict itself. For prediction tasks, it is more reasonable to utilize the query of the previous segment to generate the prediction of the unknown segment. In this case, if the given query for prediction is similar to keys of some segments, their next segments rather than these segments themselves should contribute more to the prediction. As shown in Figure 1, when we predict the flow at evening today given the flow at noon, we use the given flow at noon as the query rather than perform correlation for the unknown evening flow. If the flow at noon today is highly correlated with the flow at noon yesterday, then the flow of this evening should largely depend on the flow at evening yesterday rather than the flow at noon yesterday.

Motivated by this, we further propose a *Predictive Multi-Scale Segment-Correlation (PreMSSC)*, where current segment output  $Y_i$  can be obtained via using the previous segment  $Q_{i-1}$  to query all segments  $\{K_1, K_2, \dots\}$  and weighting the values of their next segments  $\{V_2, V_3, \dots\}$  by the calculated correlations. We derive our predictive model, namely *Predictive Transformer (Preformer)*, via replacing

the standard self-attention and cross-attention in the original Transformer model with MSSC and PreMSSC, respectively.

The main contributions of this paper are as follows:

- We propose a novel MSSC mechanism to replace the canonical point-wise attention mechanism, which can increase the efficiency, extract more relevant information from the time series, and avoid the selection of segment length.
- We design a PreMSSC paradigm for forecasting by introducing a one-segment delay between the keys and their corresponding values. We develop a long-term forecasting model namely Preformer based on MSSC and PreMSSC.
- Extensive experiments show that our Preformer model achieves better forecasting performance than other Transformer-based prediction models.

## 2. Related Work

**Time Series Forecasting.** Early works on the TSF problem are based on classical mathematical models such as vector autoregression (VAR) (Lütkepohl, 2005) and auto regressive intergrated moving average (ARIMA) (Box et al., 2015). Support vector regression (SVR) (Cao & Tay, 2003) introduces a traditional machine learning method to regress the future. Gaussian Process (Roberts et al., 2013) predicts the distribution of future values without assuming any certain form of the prediction function. However, all these classical models can not handle complicated data distributions or high-dimensional data.

With the development of deep learning, neural networks have shown stronger modeling ability than classical models. Recurrent Neural Network (RNN) (Connor et al., 1992) and

Temporal Convolution Network (TCN) (Oord et al., 2016; Bai et al., 2018) are two common types of deep models for modeling sequence data. LSTNet (Lai et al., 2018) combines convolutional layers and recurrent layers to capture both long-term and short-term dependencies. There are also some works (Qin et al., 2017; Song et al., 2018) that introduce additional attention mechanism to RNN to achieve better performance in forecasting. However, RNN-based models suffer from the gradient vanishing and gradient exploding problem. Popular variants of RNN such as LSTM (Hochreiter & Schmidhuber, 1997) and GRU (Chung et al., 2014) can not solve this problem fundamentally. The lack of parallelizability is another main limitation of RNN-based models. Benefiting from the good parallelism of convolution operations, TCN-based models (Borovykh et al., 2017; Sen et al., 2019) have also achieved good results in time series tasks. Both RNN-based and TCN-based models do not explicitly model the dependencies between two far-away temporal locations, and the information exchange between them must go through a long path.

**Transformer-based models.** Transformer (Vaswani et al., 2017) was originally proposed as a sequence-to-sequence model in natural language processing to deal with machine translation. Due to its powerful modeling capabilities, it has even been widely applied in processing non-sequential data such as images (Dosovitskiy et al., 2020; Carion et al., 2020). Self-attention plays an important role for explicitly discovering the dependencies between any element pairs, but both the time and space complexities increase quadratically with the length of the sequence, which limits the application of Transformer in long-term forecasting.

Therefore, various sparse self-attention mechanisms have been proposed in recent years. Logfomer (Li et al., 2019) proposes LogSparse self-attention which selects elements in exponentially increasing intervals to break the memory bottleneck. Informer (Zhou et al., 2021) defines a sparsity measurement for queries and selects dominant queries based on this measurement to obtain ProbSparse self-attention. These works use point-wise dot product to compute attention score, and differ in the way of selecting point pairs.

AutoFomer (Wu et al., 2021) develops an Auto-Correlation mechanism to replace self-attention, which utilizes series-wise correlation instead of point-wise dot product. In this work, we introduce a new Segment-Correlation mechanism to explore the context information within neighboring points and capture the segment-wise correlation in the sequence. Our method differs from the Auto-Correlation mechanism in the way of correlation computation and aggregation. Instead of the complicated Fast Fourier Transforms calculation in Auto-Correlation, we directly segment the time series based on implicit period and compute the correlation between segments. Benefiting from the simplicity and intuitiveness

of Segment-Correlation, we have customized a multi-scale structure and predictive paradigm for long-term forecasting.

### 3. Method

#### 3.1. Problem Definition

Multi-horizon time series forecasting aims to predict values at multiple future time steps. Typically, given the previous time series  $X_{1:t_0} = \{x_1, x_2, \dots, x_{t_0}\}$ , where  $x_t \in \mathbb{R}^{d_x}$  and  $d_x$  is the dimension of the variable, we aim to predict the future values  $Y_{t_0+1:t_0+\tau} = \{y_{t_0+1}, y_{t_0+2}, \dots, y_{t_0+\tau}\}$ , where  $y_t \in \mathbb{R}^{d_y}$  is the prediction at every time step  $t$  and  $d_y$  is the dimension of the output variable. The prediction model  $f$  can be formulated as:

$$\hat{Y}_{t_0+1:t_0+\tau} = f(X_{1:t_0}; \Omega), \tag{1}$$

where  $\hat{Y}_{t_0+1:t_0+\tau}$  is the predicted time series and  $\Omega$  is the learnable parameters of the model. For long-term forecasting, the future time duration  $\tau$  to be predicted, is longer. The problem can be categorized into two types based on whether the dimension of the output variable  $d_y$  is larger than one: univariate forecasting and multivariate forecasting.

#### 3.2. The Preformer Model

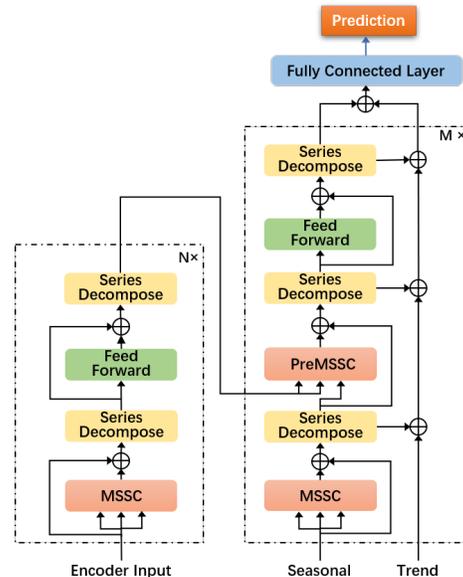


Figure 2. The overall framework of the Preformer model.

As shown in Figure 2, the overall architecture of Preformer is similar to the vanilla Transformer. We replace the dot-product self-attention in Transformer with MSSC and replace the cross-attention with PreMSSC. Time series decomposition methods which deconstruct time series into

several components are widely used in time series analysis. Decoupling the time series into several components has been shown to be helpful for forecasting tasks (Taylor & Letham, 2018; Oreshkin et al., 2019). Here we use the series decomposition module proposed in Autoformer (Wu et al., 2021) to decompose the time series into trend and seasonal components. A fully connected layer is added to convert the decoder output into prediction values, which will be used to compute the MSE loss  $\mathcal{L}_{MSE}$  with the ground truth.

**Model inputs.** The way the original time series is input to the Preformer (encoder and decoder) is consistent with AutoFormer (Wu et al., 2021). Following other methods (Zhou et al., 2021; Wu et al., 2021), we utilize additional time-dependent features (e.g., a set of dummy variables like hour-of-the-day, day-of-the-week, etc) called covariates as parts of inputs. Specifically, we concatenate the values of the past original time series  $X_{1:t_0}$  and covariates  $C_{1:t_0}$  to derive the inputs of the Preformer encoder  $\mathcal{X}_{1:t_0} = [X_{1:t_0}; C_{1:t_0}]$ . Because these covariates can be predetermined (e.g. the hour-of-the-day at any time), we also use  $C_{t_0+1:t_0+\tau}$  as parts of the decoder inputs. Before being fed to the MSSC module, inputs to the encoder and decoder are transformed into the feature dimension through an embedding layer which has been widely used in Transformer-based models.

**Encoder.** The encoder consists of  $N$  identical layers, where each layer consists of a MSSC module and a feed forward network each followed by a series decomposition module with residual connections. The input of the encoder is  $\mathcal{X}_{1:t_0}$ , which includes the past time series values and covariates. The series decomposition module passes the inputs  $\mathcal{X}$  through an average pooling layer to get the trend component  $\mathcal{X}^{trend}$ , and subtracts the trend component from the inputs to get the seasonal component  $\mathcal{X}^{season} = \mathcal{X} - \mathcal{X}^{trend}$ . All decomposition modules in the encoder eliminate the trend component, which makes the encoder focus on seasonal pattern modeling.

**Decoder.** The decoder consists of  $M$  identical layers. Different from the encoder, there is an additional PreMSSC module where key and value matrices are transformed from the outputs of the encoder in each decoder layer. The inputs of the decoder include two components: seasonal and trend components. Each component is composed of two parts: information from the latter half part of the original time series and placeholders filled by scalars. Specifically, the latter half part of the original time series  $X_{half}$  (i.e.,  $X_{t_0/2:t_0}$ ) are decomposed into  $X_{half}^{season}$  and  $X_{half}^{trend}$ , which will be concatenated with placeholders to get seasonal inputs  $[X_{half}^{season}, X_0]$  and trend inputs  $[X_{half}^{trend}, X_{mean}]$ , where  $X_0, X_{mean} \in \mathbb{R}^{\tau \times d_x}$  denote the placeholders filled with zero and the mean of  $X_{half}$  respectively. The decomposition modules in the decoder extract the trend part from hidden variables progressively, which is finally added to the

seasonal part to derive the output.

An additional fully connected layer takes the output of the decoder as input, and it generates final prediction values  $\hat{Y}_{t_0+1:t_0+\tau} \in \mathbb{R}^{\tau \times d_y}$ .

### 3.3. Segment-Correlation

Segment-Correlation is the key module in Preformer, which performs segment-wise attention instead of point-wise attention. We denote the input of each Segment-Correlation module as  $H \in \mathbb{R}^{L \times d}$ , where  $L$  and  $d$  are the length and dimension of the input respectively. Formally, for the single head situation, the input series  $H$  will be projected by three projection matrices to obtain the query, key and value, i.e.,  $Q = HW_q, K = HW_k, V = HW_v$ . Then all the  $Q, K$  and  $V$  are segmented into several segments having the same length  $L_{seg}$ :  $\{Q_1, Q_2, \dots, Q_m\}, \{K_1, K_2, \dots, K_n\}, \{V_1, V_2, \dots, V_n\}$ , where  $Q_i \in \mathbb{R}^{L_{seg} \times d}, K_i \in \mathbb{R}^{L_{seg} \times d}, V_i \in \mathbb{R}^{L_{seg} \times d}, m, n$  denote the number of segments, and  $L_{seg}$  is a hyperparameter that determines the computational complexity via controlling the length of segments.

The correlation measurement  $c_{ij}$  between any pair of query segment  $Q_i$  and key segment  $K_j$  can be computed by the function:

$$c_{ij} = \text{Correlation}(Q_i, K_j) = \frac{1}{d \times L_{seg}} Q_i \odot K_j, \quad (2)$$

where  $\odot$  is the dot product operator between two matrices of the same size. For each query segment  $Q_i$ , whose correlation measurements with all the key segments will be normalized by the Softmax function to obtain the aggregation weight:

$$\hat{c}_{i1}, \hat{c}_{i2}, \dots, \hat{c}_{in} = \text{Softmax}(c_{i1}, c_{i2}, \dots, c_{in}). \quad (3)$$

The output at the position of the  $i$ -th segment  $Y_i$  is the weighted sum of all the value segments  $\{V_j \mid j = 1, \dots, n\}$ :

$$Y_i = \sum_{j=1}^n \hat{c}_{ij} V_j. \quad (4)$$

Lastly, output of the Segment-Correlation module can be obtained by concatenating all the  $Y_i$  along length dimension:

$$\text{SC}(H; L_{seg}) = \text{Concat}(Y_1, \dots, Y_m), \quad (5)$$

where SC is the abbreviation of Segment-Correlation. The output of multi-head Segment-Correlation can be computed by concatenating and projecting the outputs of all heads described above. We omit the formulation since it is similar to the canonical multi-head attention (Vaswani et al., 2017).

**Multi-Scale Segment-Correlation.** Since  $L_{seg}$  determines the resolution of the smallest unit involved in the Segment-Correlation calculation, a large  $L_{seg}$  means

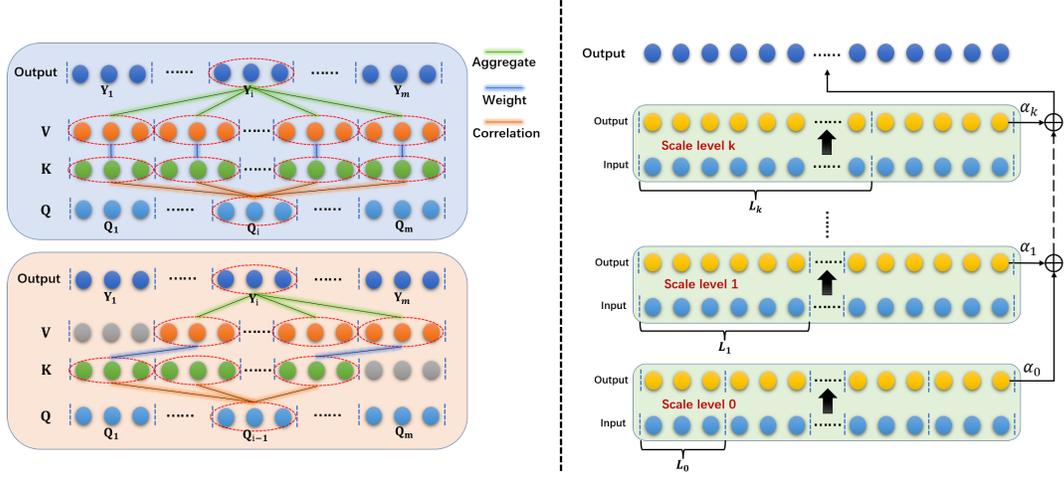


Figure 3. Segment-Correlation (upper left), predictive paradigm (bottom Left) and the multi-scale architecture (right). In Segment-Correlation, the query, key and value are segmented evenly. Then the correlation measurements between all segments are calculated, and the derived weights are used to aggregate the value segments to get the output. Only the query of the  $i$ -th segment  $Q_i$  and the output  $Y_i$  are illustrated here for simplicity. In the multi-scale structure, each scale level is a Segment-Correlation module with or without predictive paradigm, representing MSSC and PreMSSC respectively.

that the coarse-grained temporal dependencies in the time series can be captured, while Segment-Correlation with a small  $L_{seg}$  can capture fine-grained dependencies. To attenuate the impact of hyperparameter  $L_{seg}$  selection on performance, we propose Multi-Scale Segment-Correlation (MSSC) by fusing the output of multiple Segment-Correlation with different  $L_{seg}$ . Specifically, as the scale level increases, we start with a small initial segment length  $L_0$  and increase the segment length exponentially, i.e.,  $L_l = 2^l L_0$ , where  $l \in \{0, 1, \dots, l_{max}\}$  denotes the scale level and  $l_{max} = \lfloor \log_2(\frac{L}{L_0}) \rfloor$ . The input  $H$  is passed through these Segment-Correlation layers of different scale levels to obtain the outputs of the corresponding scale levels. We aggregate the outputs of all scale levels to get the output of the entire MSSC module, where the weight of the  $l$ -th level  $\alpha_l$  is set to decrease exponentially as the scale level increases. Therefore, MSSC can be formulated as:

$$\text{MSSC}(H) = \sum_{l=0}^{l_{max}} \alpha_l \cdot \text{SC}(H; 2^l L_0), \quad (6)$$

$$\alpha_l = \frac{2^l}{\sum_{l=0}^{l_{max}} 2^l}.$$

**Predictive paradigm.** In the decoding phase, the query for the period to be predicted is its preceding segment rather than itself. Therefore, if some segments with respect to keys are highly relevant to the query, their future segments rather than themselves should contribute more to the prediction of the query. Inspired by this intuition, we propose a predictive paradigm for cross-attention by introducing a segment delay between the keys and their corresponding values. The

basic idea of the predictive paradigm is shown in Figure 1, Section 1. For Segment-Correlation without the predictive paradigm, the output at the position of the  $i$ -th segment  $Y_i$  is obtained according to Equation (4). If we reformulate  $\hat{c}_{ij}$  as  $\hat{c}_{(Q_i, K_j)}$ , then we have:

$$Y_i = \sum_{j=1}^n \hat{c}_{(Q_i, K_j)} V_j. \quad (7)$$

There are two differences between the predictive paradigm and the non-predictive paradigm. Firstly, to get the output of the current segment  $Y_i$ , the query of the previous segment  $Q_{i-1}$  is used to calculate correlations with the keys of all segments  $\{K_1, \dots, K_{n-1}\}$ . Secondly, the correlation  $\hat{c}_{(Q_{i-1}, K_j)}$  corresponding to the key of a certain segment  $K_j$  is regarded as the weight of the next segment  $V_{j+1}$  to aggregate values. That is, the segment of value is delayed by one segment relative to the segment of the corresponding key. Therefore, we can obtain  $Y_i$  by the following equation:

$$Y_1 = \sum_{j=1}^{n-1} \hat{c}_{(Q_m, K_j)} V_{j+1}, \quad (8)$$

$$Y_i = \sum_{j=1}^{n-1} \hat{c}_{(Q_{i-1}, K_j)} V_{j+1}, \quad i > 1.$$

### 3.4. Complexity Analysis

For Single-Scale Segment-Correlation, if the segment length  $L_{seg} = L_0$ , the computational complexity is  $\mathcal{O}(L^2/L_0)$ . For Multi-Scale Segment-Correlation, the computational

Table 1. Multivariate time-series forecasting results on six datasets

Models	Preformer		Autoformer		Informer		LogTrans		LSTNet		LSTM		TCN		
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
ETTh2	96	<b>0.213</b>	<b>0.295</b>	<u>0.255</u>	<u>0.339</u>	0.365	0.453	0.768	0.642	3.142	1.365	2.041	1.073	3.041	1.330
	192	<b>0.269</b>	<b>0.329</b>	<u>0.281</u>	<u>0.340</u>	0.533	0.563	0.989	0.757	3.154	1.369	2.249	1.112	3.072	1.339
	336	<b>0.324</b>	<b>0.363</b>	<u>0.339</u>	<u>0.372</u>	1.363	0.887	1.334	0.872	3.160	1.369	2.568	1.238	3.105	1.348
	720	<b>0.418</b>	<b>0.416</b>	<u>0.422</u>	<u>0.419</u>	3.379	1.388	3.048	1.328	3.171	1.368	2.720	1.287	3.135	1.354
Electricity	96	<b>0.180</b>	<b>0.297</b>	<u>0.201</u>	<u>0.317</u>	0.274	0.368	0.258	0.357	0.680	0.645	0.375	0.437	0.985	0.813
	192	<b>0.189</b>	<b>0.302</b>	<u>0.222</u>	<u>0.334</u>	0.296	0.386	0.266	0.368	0.725	0.676	0.442	0.473	0.996	0.821
	336	<b>0.201</b>	<b>0.319</b>	<u>0.231</u>	<u>0.338</u>	0.300	0.394	0.280	0.380	0.828	0.727	0.439	0.473	1.000	0.824
	720	<b>0.232</b>	<b>0.342</b>	<u>0.254</u>	<u>0.361</u>	0.373	0.439	0.283	0.376	0.957	0.811	0.980	0.814	1.438	0.784
Exchange	96	<b>0.148</b>	<b>0.282</b>	<u>0.197</u>	<u>0.323</u>	0.847	0.752	0.968	0.812	1.551	1.058	1.453	1.049	3.004	1.432
	192	<b>0.268</b>	<b>0.378</b>	<u>0.300</u>	<b>0.369</b>	1.204	0.895	1.040	0.851	1.477	1.028	1.846	1.179	3.048	1.444
	336	<b>0.447</b>	<b>0.499</b>	<u>0.509</u>	<u>0.524</u>	1.672	1.036	1.659	1.081	1.507	1.031	2.136	1.231	3.113	1.459
	720	<b>1.092</b>	<b>0.812</b>	<u>1.447</u>	<u>0.941</u>	2.478	1.310	1.941	1.127	2.285	1.243	2.984	1.427	3.150	1.458
Traffic	96	<b>0.560</b>	<b>0.349</b>	<u>0.613</u>	<u>0.388</u>	0.719	0.391	0.684	0.384	1.107	0.685	0.843	0.453	1.438	0.784
	192	<b>0.565</b>	<b>0.349</b>	<u>0.616</u>	<u>0.382</u>	0.696	0.379	0.685	0.390	1.157	0.706	0.847	0.453	1.463	0.794
	336	<b>0.577</b>	<u>0.351</u>	<u>0.622</u>	<b>0.337</b>	0.777	0.420	0.733	0.408	1.216	0.730	0.853	0.455	1.479	0.799
	720	<b>0.597</b>	<b>0.358</b>	<u>0.660</u>	<u>0.408</u>	0.864	0.472	0.717	0.396	1.481	0.805	1.500	0.805	1.499	0.804
Weather	96	<b>0.227</b>	<b>0.292</b>	<u>0.266</u>	<u>0.336</u>	0.300	0.384	0.458	0.490	0.594	0.587	0.369	0.406	0.615	0.589
	192	<b>0.275</b>	<b>0.322</b>	<u>0.307</u>	<u>0.367</u>	0.598	0.544	0.658	0.589	0.560	0.565	0.416	0.435	0.629	0.600
	336	<b>0.324</b>	<b>0.352</b>	<u>0.359</u>	<u>0.395</u>	0.578	0.523	0.797	0.652	0.597	0.587	0.455	0.454	0.639	0.608
	720	<b>0.394</b>	<b>0.393</b>	<u>0.419</u>	<u>0.428</u>	1.059	0.741	0.869	0.675	0.618	0.599	0.535	0.520	0.639	0.610
ILI	24	<b>3.143</b>	<b>1.185</b>	<u>3.483</u>	<u>1.287</u>	5.764	1.677	4.480	1.444	6.026	1.770	5.914	1.734	6.624	1.830
	36	<b>2.793</b>	<b>1.054</b>	<u>3.103</u>	<u>1.148</u>	4.755	1.467	4.799	1.467	5.340	1.668	6.631	1.845	6.858	1.879
	48	<b>2.845</b>	<u>1.090</u>	<b>2.669</b>	<b>1.085</b>	4.763	1.469	4.800	1.468	6.080	1.787	6.736	1.857	6.968	1.892
	60	<u>2.957</u>	<b>1.124</b>	<b>2.770</b>	<u>1.125</u>	5.264	1.564	5.278	1.560	5.548	1.720	6.870	1.879	7.127	1.918

<sup>1</sup> Reported metrics except Preformer all come from the Autoformer paper (Wu et al., 2021).

complexity is the sum of all scales, which is still at the same computational level as  $\mathcal{O}(L^2/L_0)$  benefiting from the exponential increase in segment length. We demonstrate this in Appendix B due to space limitations.

Most other sparse mechanisms introduce some extra operations, such as selections of dominant queries in Informer and fast Fourier transform in Autoformer. Although their theoretical complexity is  $\mathcal{O}(L \log L)$ , the actual operation efficiency is not even as good as our Multi-Scale Segment-Correlation. Please see Section 4.4 for more details.

## 4. Experiments

### 4.1. Experimental Setup

**Datasets** We conduct experiments on the following six datasets as in Wu et al. (2021). (1) *ETT*<sup>1</sup> contains data related to electricity which is collected from two Chinese stations in two years. In order to explore the model’s performance on data with different granularities, we use different sampling frequencies to get hourly data  $\{ETTh1, ETTh2\}$  and 15-minutes data  $\{ETTm1, ETTm2\}$ . (2) *Electricity*<sup>2</sup> contains the hourly electricity consumption of 321 clients in 2 years. (3) *Exchange* (Lai et al., 2018) collects the daily exchange rates of eight countries ranging from 1990

to 2016. (4) *Traffic*<sup>3</sup> collects the hourly road occupancy rates from the California Department of Transportation in two years. (5) *Weather*<sup>4</sup> records climatological data of the Max-Planck-Institute every 10 minutes in 2020 year, which contains 21 climate features including air pressure and temperature etc. (6) *ILI*<sup>5</sup> collects the weekly outpatient visits for influenza-like illness (ILI) from 2002 to 2021. We split the datasets following Wu et al. (2021). The train/val/test contains 12/4/4 months of data for the ETT datasets, while we split other datasets into train/val/test by the ratio of 7:1:2.

**Implementation details** We use the ADAM (Kingma & Ba, 2014) optimizer with an initial learning rate 1e-4 and the learning rate decay to train our model. Early stop training strategy is utilized to avoid overfitting. The number of training epochs is set to 10 and the batch size is set to 32. All experiments are implemented with PyTorch (Paszke et al., 2019) and conducted on single NVIDIA TITAN RTX 24GB GPU. For all experimental settings, Preformer contains 2 encoder layers and 1 decoder layer.

**Baselines** Several models are selected to compare with Preformer in multivariate time-series forecasting, including three transformer-based models: Autoformer (Wu et al., 2021), Informer (Zhou et al., 2021), LogTrans (Li et al.,

<sup>3</sup><http://pems.dot.ca.gov/>.

<sup>4</sup><https://www.bgc-jena.mpg.de/wetter/>.

<sup>5</sup><https://gis.cdc.gov/grasp/fluview/fluportaldashboard.html>.

<sup>1</sup><https://github.com/zhouhaoyi/ETDataset>.

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/ElectricityLoadDiagrams20112014>.

Table 2. Univariate time-series forecasting results on two typical datasets

Models	Preformer		Autoformer		N-BEATS		Informer		LogTrans		DeepAR		Prophet		ARIMA		
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
ETTM2	96	0.072	0.205	<b>0.065</b>	<b>0.189</b>	0.082	0.219	0.088	0.225	0.082	0.217	0.099	0.237	0.287	0.456	0.211	0.362
	192	<b>0.109</b>	<b>0.255</b>	0.118	0.256	0.120	0.268	0.132	0.283	0.133	0.284	0.154	0.310	0.312	0.483	0.261	0.406
	336	<b>0.139</b>	<b>0.295</b>	0.154	0.305	0.226	0.370	0.180	0.336	0.201	0.361	0.277	0.428	0.331	0.474	0.317	0.448
	720	<b>0.167</b>	<b>0.327</b>	0.182	0.335	0.188	0.338	0.300	0.435	0.268	0.407	0.332	0.468	0.534	0.593	0.366	0.487
Exchange	96	0.141	0.293	0.241	0.387	0.156	0.299	0.591	0.615	0.279	0.441	0.417	0.515	0.828	0.762	<b>0.112</b>	<b>0.245</b>
	192	<b>0.252</b>	<b>0.389</b>	0.273	0.403	0.669	0.665	1.183	0.912	1.950	1.048	0.813	0.735	0.909	0.974	0.304	0.404
	336	<b>0.426</b>	<b>0.512</b>	0.508	0.539	0.611	0.605	1.367	0.984	2.438	1.262	1.331	0.962	1.304	0.988	0.736	0.598
	720	<b>0.790</b>	<b>0.712</b>	0.991	0.768	1.111	0.860	1.872	1.072	2.010	1.247	1.894	1.181	3.238	1.566	1.871	0.935

<sup>1</sup> Reported metrics except Preformer all come from the Autoformer paper (Wu et al., 2021).

2019), two RNN-based models: LSTM (Hochreiter & Schmidhuber, 1997), LSTNet (Lai et al., 2018) and one TCN-based model: TCN (Bai et al., 2018). Since there are some models designed specifically for univariate forecasting, we compare Preformer with more competitive baselines, including two deep learning models: N-BEATS (Oreshkin et al., 2019) and DeepAR (Salinas et al., 2020), two traditional machine learning methods: Prophet (Taylor & Letham, 2018) and ARIMA (Box et al., 2015).

#### 4.2. Main Results

To comprehensively compare Preformer and baselines, we conduct thorough multivariate and univariate forecasting experiments on various datasets under multiple settings. For the ILI dataset, the input length is fixed to 24 and the prediction length includes {24, 36, 48, 60}. For other datasets, the input length is fixed to 96 and the prediction length is chosen from {96, 192, 336, 720}.

**Multivariate Time-series Forecasting** From Table 1, we find that Preformer is better than other models except Autoformer in all cases and performs slightly worse than Autoformer in only a few settings. For example, under the input-96-predict-192 setting, compared to previous state-of-the-art results, Preformer has achieved **4.3%** (0.281 → 0.269) relative MSE improvement in ETTm2, **14.9%** (0.222 → 0.189) in Electricity, **10.7%** (0.300 → 0.268) in Exchange, **8.3%** (0.616 → 0.565) in Traffic, **10.4%** (0.307 → 0.275) in Weather, and **10%** (3.103 → 2.793) in ILI. Moreover, Preformer shows stable performance in cases of long prediction horizons, which shows that it is suitable for long-term forecasting. Besides, the overall performance of Transformer-based models is better than RNN-based models and TCN-based models, proving the potential of the Transformer-based models in time-series forecasting.

**Univariate Time-series Forecasting** We show the results on two typical datasets, i.e., the ETTm2 dataset with obvious periodicity and the Exchange dataset without obvious periodicity in Table 2. We observe that our Preformer achieves the best results in all long-term forecasting cases

whose prediction length is not less than 192. For example, under the input-96-predict-336 setting, compared to previous state-of-the-art results, Preformer has achieved **9.7%** (0.154 → 0.139) relative improvement on MSE in ETTm2 and **16%** (0.508 → 0.426) in Exchange. Also, Autoformer performs best in the input-96-predict-96 setting of the ETTm2 dataset while ARIMA performs best in the same setting of the Exchange dataset, which illustrates their advantages in shorter-term prediction.

#### 4.3. Ablation Study

##### Performance of the Segment-Correlation mechanism

We conduct experiments on the ETTh1 dataset to compare several different sparse attention mechanisms. Segment-correlation in Table 3 refers to the Preformer model with the multi-scale structure and predictive paradigm, and the initial segment length  $L_0$  is set to 4. For a fair comparison, we replace MSSC and PreMSSC with other sparse attention mechanisms and set hyperparameters such as the hidden dimension and the head numbers to be consistent. From Table 3, we observe that compared with other sparse mechanisms, Segment-Correlation achieves the best performance. Full attention and LogSparse attention fail in long-term forecasting due to high memory complexity.

Table 3. Ablations of the Segment-Correlation mechanism

Input Length	96			336			
	Prediction Length	336	720	1440	336	720	1440
Segment-Correlation	MSE	<b>0.628</b>	<b>0.632</b>	<b>0.769</b>	<b>0.560</b>	<b>0.581</b>	<b>0.755</b>
	MAE	<b>0.551</b>	<b>0.571</b>	<b>0.650</b>	<b>0.531</b>	<b>0.558</b>	<b>0.657</b>
Auto-Correlation	MSE	0.630	0.687	0.821	0.597	0.619	0.859
	MAE	0.552	0.593	0.671	0.539	0.578	0.692
Full Attention	MSE	0.673	0.679	-	0.661	0.678	-
	MAE	0.564	0.585	-	0.566	0.592	-
LogSparse Attention	MSE	0.679	0.688	-	0.646	0.668	-
	MAE	0.567	0.591	-	0.565	0.592	-
ProbSparse Attention	MSE	0.702	0.705	0.831	0.689	0.701	0.861
	MAE	0.577	0.599	0.673	0.579	0.604	0.694

<sup>1</sup> The "-" indicates the out-of-memory.

**Impact of the multi-scale structure** The multi-scale structure can effectively extract dependencies at different temporal resolutions, which is important for time series forecasting. To illustrate this, we remove the multi-scale structure from all PreMSSC and MSSC modules in Preformer to get the model without multi-scale structure. As shown in Table 4, the prediction performance of Preformer with the multi-scale structure (*MS+Predictive*) is better than that without the multi-scale structure (*Only Predictive*) in all cases. Especially on the ETTh1 and Exchange datasets, the multi-scale structure can lead to considerable performance improvements.

**Impact of the predictive paradigm** To explore whether the predictive paradigm is really effective for prediction tasks, we conduct ablation studies of the predictive paradigm on three datasets with different settings. As shown in Table 4, *MS+predictive* represents the standard Preformer model, while *Only MS* means that replacing the PreMSSC module in Preformer’s decoder with the MSSC module. The experimental results show that Preformer with the predictive paradigm achieves better performance in almost all cases, which proves that the proposed predictive paradigm is helpful for the prediction tasks.

Table 4. Ablations of predictive paradigm and multi-scale structure

Models	<i>MS+Predictive</i> (ours)		<i>Only Predictive</i> (without MS)		<i>Only MS</i> (without Predictive)		
	MSE	MAE	MSE	MAE	MSE	MAE	
ETTh1	96	<b>0.414</b>	0.439	0.480	0.472	0.416	<b>0.436</b>
	192	<b>0.445</b>	<b>0.455</b>	0.493	0.484	0.472	0.472
	336	<b>0.466</b>	<b>0.468</b>	0.511	0.497	0.471	0.475
	720	<b>0.471</b>	<b>0.487</b>	0.605	0.553	0.484	0.499
Exchange	96	<b>0.148</b>	<b>0.282</b>	0.186	0.305	0.187	0.305
	192	<b>0.268</b>	<b>0.378</b>	0.359	0.457	0.280	0.386
	336	<b>0.447</b>	<b>0.499</b>	0.704	0.659	0.520	0.542
	720	<b>1.092</b>	<b>0.812</b>	1.400	0.931	1.232	0.883
Traffic	96	<b>0.560</b>	<b>0.349</b>	0.561	0.352	0.567	0.351
	192	<b>0.565</b>	<b>0.349</b>	0.573	0.356	0.583	0.358
	336	<b>0.577</b>	<b>0.351</b>	<b>0.577</b>	0.353	0.581	0.354
	720	<b>0.597</b>	<b>0.358</b>	0.599	0.363	0.598	0.362

**The sensitivity analysis of the segment length** The segment length  $L_{seg}$  is a critical hyperparameter in Segment-Correlation. The multi-scale structure can facilitate the selection of segment length. We conduct multivariate forecasting experiments (input-96 to predict-336) on the ETTh1 dataset to explore the sensitivity of segment length using our Preformer with or without the multi-scale structure. As shown in Figure 4, Preformer without the multi-scale structure is very sensitive to the choice of segment length, while the performance of Preformer with the multi-scale structure does not change much with the segment length. More detailed results can be found in Table 5. These results all demonstrate the effectiveness of the multi-scale structure.

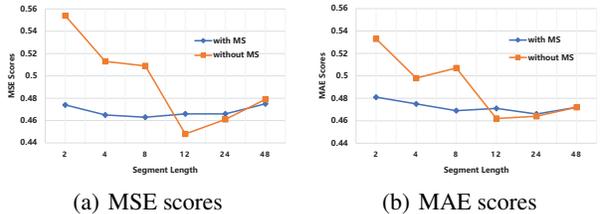


Figure 4. Effects of the multi-scale structure on sensitivity of segment length hyperparameter.

#### 4.4. Efficiency analysis

As shown in Figure 5, we compare the running memory and time of models with different sparse attention mechanisms during the training phase. Segment-Correlation in the figure means MSSC with the multi-scale structure and  $L_0$  is set to 4. For memory efficiency analysis, we set the batch size to 16 in order to avoid out-of-memory. For time efficiency analysis, we run each attention mechanism 1000 times to get the average training time. It is worth noting that our Multi-Scale Segment-Correlation is more efficient in time and memory than other sparse mechanisms whose theoretical complexity  $\mathcal{O}(L \log L)$  in practical applications.

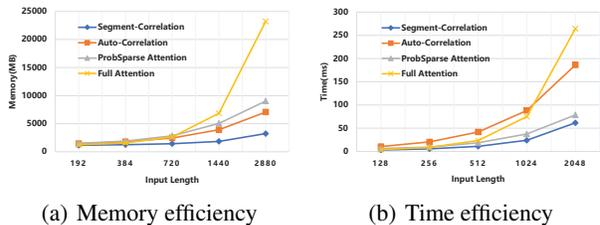


Figure 5. Memory and time consumption during the training phase.

## 5. Conclusion

In this paper, we propose a Transformer-based model called Preformer for long-term time series forecasting. In Preformer, we introduce a sparse and efficient attention mechanism called Multi-Scale Segment-Correlation (MSSC) which utilizes correlations between segment pairs to discover dependencies and aggregate information in time series. Further, we design a predictive paradigm and combine it with MSSC to get Predictive Multi-Scale Segment-Correlation (PreMSSC), which can discover predictive dependencies from contexts. Under various experimental settings in different datasets, Preformer with MSSC and PreMSSC can yield state-of-the-art prediction performance, which demonstrates the effectiveness of our Preformer.

## References

- Bai, S., Kolter, J. Z., and Koltun, V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- Borovykh, A., Bohte, S., and Oosterlee, C. W. Conditional time series forecasting with convolutional neural networks. *arXiv preprint arXiv:1703.04691*, 2017.
- Box, G. E., Jenkins, G. M., Reinsel, G. C., and Ljung, G. M. *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.
- Cao, L.-J. and Tay, F. E. H. Support vector machine with adaptive parameters in financial time series forecasting. *IEEE Transactions on neural networks*, 14(6):1506–1518, 2003.
- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pp. 213–229, 2020.
- Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- Connor, J., Atlas, L. E., and Martin, D. R. Recurrent networks and narma modeling. In *In Advances in Neural Information Processing Systems*, pp. 301–308, 1992.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9:1735–1780, 1997.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2014.
- Kitaev, N., Kaiser, L., and Levskaya, A. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2019.
- Lai, G., Chang, W.-C., Yang, Y., and Liu, H. Modeling long- and short-term temporal patterns with deep neural networks. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pp. 95–104, 2018.
- Li, S., Jin, X., Xuan, Y., Zhou, X., Chen, W., Wang, Y.-X., and Yan, X. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Lütkepohl, H. *New Introduction to Multiple Time Series Analysis*. Springer-Verlag Berlin Heidelberg, 2005.
- Oord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. Wavenet: A generative model for raw audio. In *SSW*, 2016.
- Oreshkin, B. N., Carпов, D., Chapados, N., and Bengio, Y. N-beats: Neural basis expansion analysis for interpretable time series forecasting. *arXiv preprint arXiv:1905.10437*, 2019.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, volume 32, pp. 8026–8037, 2019.
- Qin, Y., Song, D., Chen, H., Cheng, W., Jiang, G., and Cottrell, G. A dual-stage attention-based recurrent neural network for time series prediction. In *Twenty-Sixth International Joint Conference on Artificial Intelligence*, 2017.
- Roberts, S., Osborne, M., Ebdon, M., Reece, S., Gibson, N., and Aigrain, S. Gaussian processes for time-series modelling. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 371 (1984):20110550, 2013.
- Salinas, D., Flunkert, V., Gasthaus, J., and Januschowski, T. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.
- Sen, R., Yu, H.-F., and Dhillon, I. Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2019.
- Song, H., Rajan, D., Thiagarajan, J. J., and Spanias, A. Attend and diagnose: Clinical time series analysis using attention models. In *Thirty-second AAAI conference on artificial intelligence*, 2018.
- Taylor, S. J. and Letham, B. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.

Wu, H., Xu, J., Wang, J., and Long, M. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.

Wu, N., Green, B., Ben, X., and O'Banion, S. Deep transformer models for time series forecasting: The influenza prevalence case. *arXiv preprint arXiv:2001.08317*, 2020.

Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *The Thirty-Fifth AAAI Conference on Artificial Intelligence*, volume 35, pp. 11106–11115, 2021.

## A. Evaluation Metrics

On each prediction window, we use two metrics which are commonly used in prediction task to evaluate performances, i.e., Mean Square Error (MSE) and Mean Absolute Error (MAE). They are defined as follows:

$$\begin{aligned} \text{MSE} &= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \\ \text{MAE} &= \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|, \end{aligned} \tag{9}$$

where  $y_i$  is ground truth and  $\hat{y}_i$  is prediction result. For multivariate prediction, the metrics can be calculated and averaged on each single variable.

## B. The Computational Complexity of Multi-Scale Segment-Correlation

The overall computational complexity of the MSSC module is the sum of all scales, which can be formulated as:

$$\begin{aligned} \mathcal{O}(\text{MSSC}) &= \mathcal{O}\left(\frac{L^2}{2^0 L_0}\right) + \mathcal{O}\left(\frac{L^2}{2^1 L_0}\right) + \dots + \mathcal{O}\left(\frac{L^2}{2^{\lfloor \log_2(\frac{L}{L_0}) \rfloor} L_0}\right) \\ &= \left[2 - \left(\frac{1}{2}\right)^{\lfloor \log_2(\frac{L}{L_0}) \rfloor}\right] \mathcal{O}\left(\frac{L^2}{L_0}\right) \\ &\leq \left[2 - \left(\frac{1}{2}\right)^{\log_2(\frac{L}{L_0})}\right] \mathcal{O}\left(\frac{L^2}{L_0}\right) \\ &= \left(2 - \frac{L_0}{L}\right) \mathcal{O}\left(\frac{L^2}{L_0}\right). \end{aligned} \tag{10}$$

We find that the computational complexity of MSSC does not exceed twice the computational complexity of the Segment-Correlation without the multi-scale structure (i.e., Single-Scale Segment-Correlation).

## C. Sensitivity of the Segment Length

We explore the impact of the multi-scale structure by comparing the effect of choice at different segment lengths on MSE scores with or without the multi-scale structure. We use the standard deviation of the MSE scores to evaluate the stability of forecasting performance. As shown in Table 5, models with multi-scale structure have more stable performance on all datasets, i.e., the MSE scores vary less with different choices of segment length.

Table 5. In the cases of the Preformer model with or without multi-scale structure, the MSE scores vary with the segment length on four datasets

Segment length		2	4	8	12	24	48	std ( $\times 1000$ )
ETTm2	with MS	0.326	0.327	0.326	0.327	0.327	0.327	<b>0.47</b>
	without MS	0.328	0.327	0.326	0.325	0.327	0.324	1.34
Electricity	with MS	0.201	0.195	0.199	0.196	0.211	0.206	<b>5.62</b>
	without MS	0.207	0.200	0.198	0.200	0.216	0.207	6.16
Traffic	with MS	0.577	0.578	0.580	0.583	0.585	0.597	<b>6.7</b>
	without MS	0.596	0.583	0.584	0.581	0.590	0.600	7.02
Weather	with MS	0.325	0.325	0.326	0.327	0.332	0.327	<b>2.38</b>
	without MS	0.322	0.326	0.342	0.344	0.337	0.326	8.53



## F. Visualization of Multivariate Time-Series Forecasting

### F.1. Comparison of Transformer-based Models

As shown in Figures 6 to 9, we plot the multivariate forecasting results of several Transformer-based models on the Electricity dataset. Blue lines are the ground truth and orange lines are the prediction results. Our Preformer can accurately predict the periodicity, trend and even some small fluctuations. Though under a very long prediction horizon, i.e., under the input-96-predict-720 setting, our Preformer can also perform well.

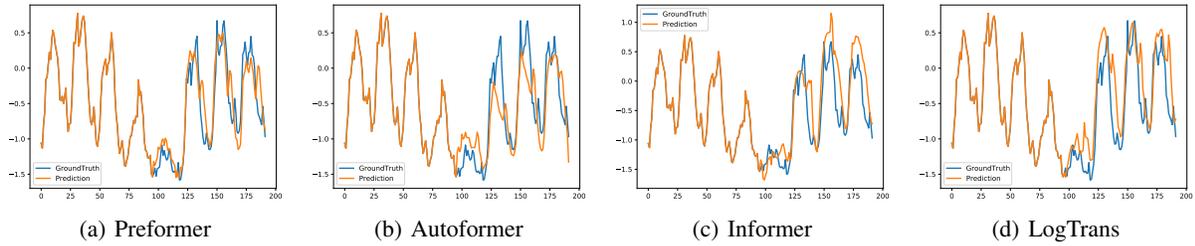


Figure 6. The prediction results on the Electricity dataset under the input-96-predict-96 setting.

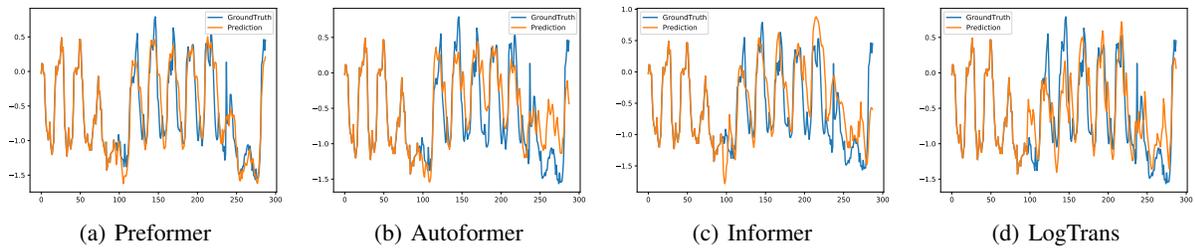


Figure 7. The prediction results on the Electricity dataset under the input-96-predict-192 setting.

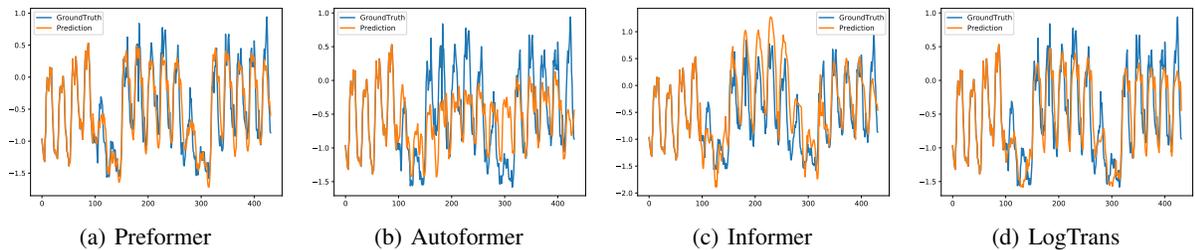


Figure 8. The prediction results on the Electricity dataset under the input-96-predict-336 setting.

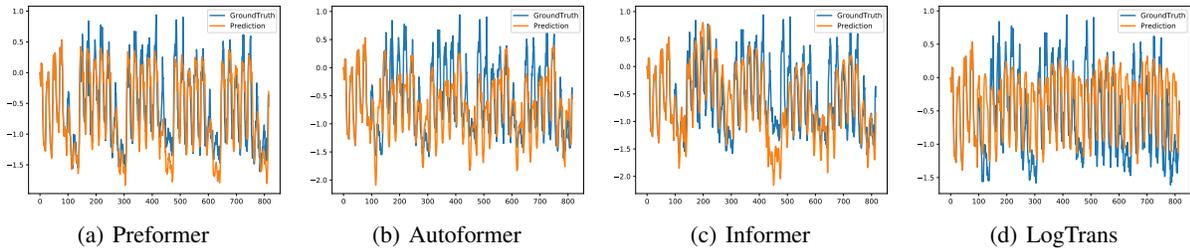


Figure 9. The prediction results on the Electricity dataset under the input-96-predict-720 setting.

### F.2. Results of Preformer on various datasets

We plot forecasting results of our Preformer on more datasets in Figures 10 to 12. We find that Preformer can obtain satisfactory prediction results whether on datasets with strong periodicity or on datasets with strong noise.

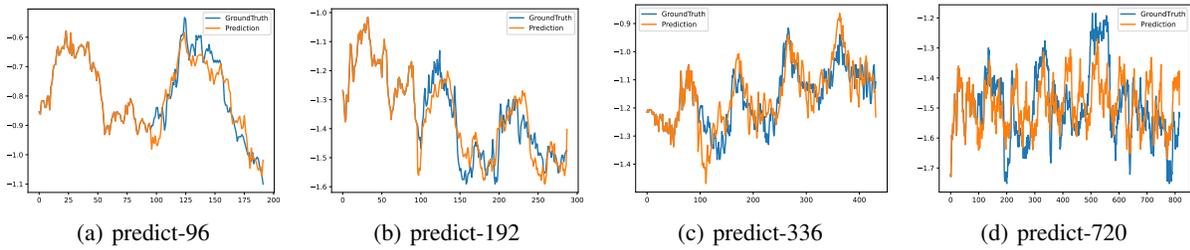


Figure 10. Some prediction results of Preformer on the ETTm1 dataset under different prediction lengths.

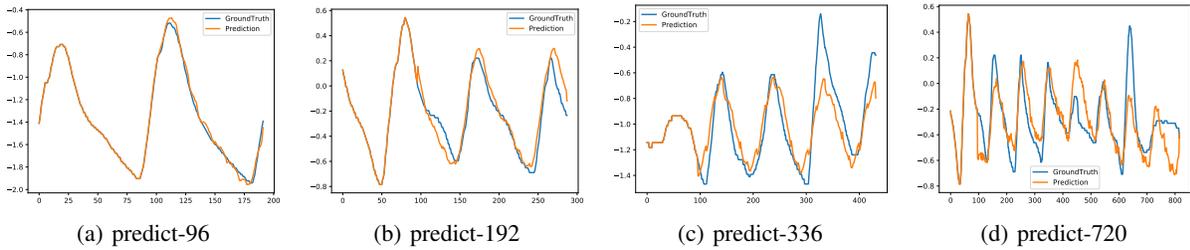


Figure 11. Some prediction results of Preformer on the ETTm2 dataset under different prediction lengths.

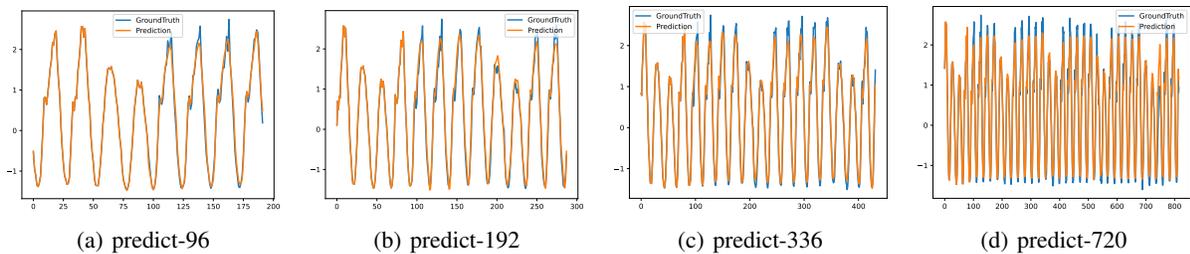


Figure 12. Some prediction results of Preformer on the Traffic dataset under different prediction lengths.

## G. Visualization of Univariate Time-Series Forecasting

### G.1. Comparison of Transformer-based Models

As shown in Figures 13 to 16, we plot the univariate forecasting results of several Transformer-based models on the ETTm2 dataset.

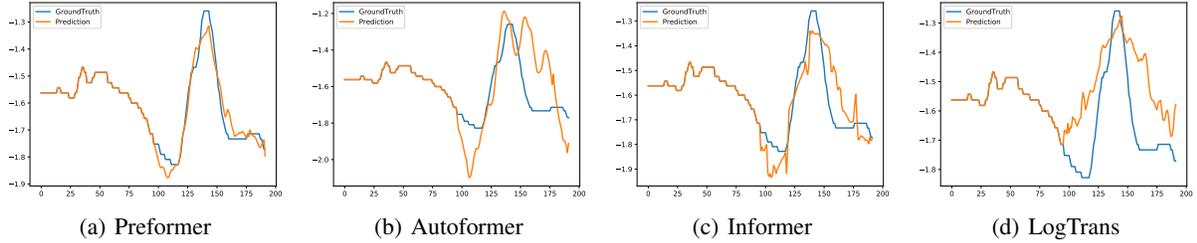


Figure 13. The prediction results on the ETTm2 dataset under the input-96-predict-96 setting.

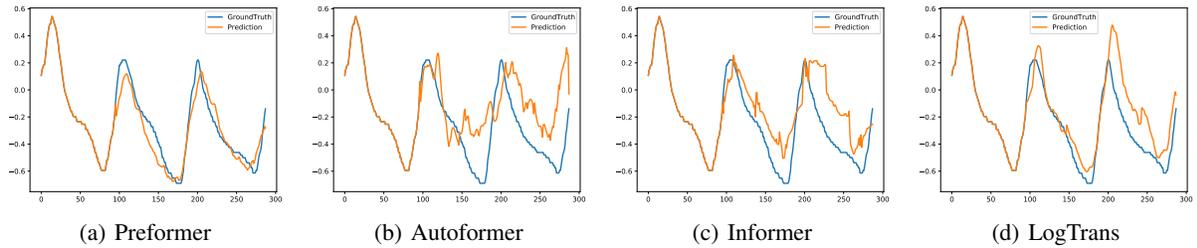


Figure 14. The prediction results on the ETTm2 dataset under the input-96-predict-192 setting.

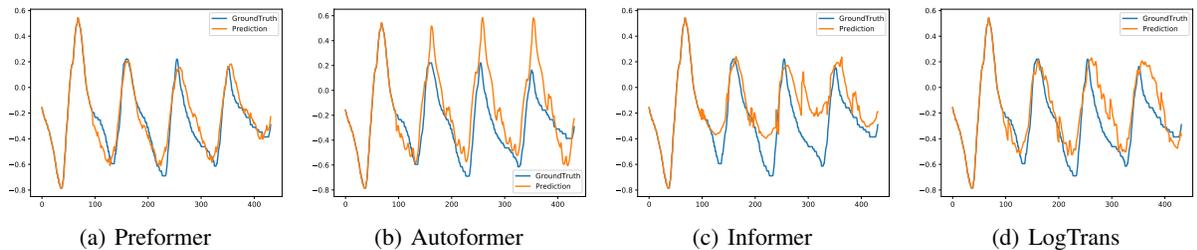


Figure 15. The prediction results on the ETTm2 dataset under the input-96-predict-336 setting.

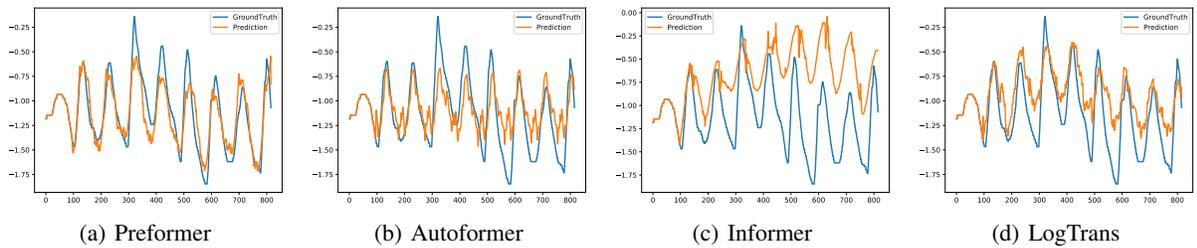


Figure 16. The prediction results on the ETTm2 dataset under the input-96-predict-720 setting.