

# Transformers in Time Series: A Survey

Qingsong Wen<sup>1</sup>, Tian Zhou<sup>2</sup>, Chaoli Zhang<sup>2</sup>, Weiqi Chen<sup>2</sup>, Ziqing Ma<sup>2</sup>, Junchi Yan<sup>3</sup>, Liang Sun<sup>1</sup>

<sup>1</sup>DAMO Academy, Alibaba Group, Bellevue, USA

<sup>2</sup>DAMO Academy, Alibaba Group, Hangzhou, China

<sup>3</sup>Department of CSE, MoE Key Lab of Artificial Intelligence, Shanghai Jiao Tong University  
{qingsong.wen, tian.zt, chaoli.zcl, jarvus.cwq, maziqing.mzq, liang.sun}@alibaba-inc.com,  
yanjunchi@sjtu.edu.cn

## Abstract

Transformers have achieved superior performances in many tasks in natural language processing and computer vision, which also triggered great interests in the time series community. Among multiple advantages of transformers, the ability to capture long-range dependencies and interactions is especially attractive for time series modeling, leading to exciting progress in various time series applications. In this paper, we systematically review transformer schemes for time series modeling by highlighting their strengths as well as limitations. In particular, we examine the development of time series transformers in two perspectives. From the perspective of network structure, we summarize the adaptations and modification that have been made to transformer in order to accommodate the challenges in time series analysis. From the perspective of applications, we categorize time series transformers based on common tasks including forecasting, anomaly detection, and classification. Empirically, we perform robust analysis, model size analysis, and seasonal-trend decomposition analysis to study how transformers perform in time series. Finally, we discuss and suggest future directions to provide useful research guidance. A corresponding resource list which will be continuously updated can be found in the GitHub repository<sup>1</sup>.

## 1 Introduction

The innovation of Transformer in deep learning [Vaswani *et al.*, 2017] has brought great interests recently due to its excellent performances in natural language processing (NLP) [Kenton and Toutanova, 2019], computer vision (CV) [Dosovitskiy *et al.*, 2021], speech processing [Dong *et al.*, 2018], and other disciplines [Chen *et al.*, 2021b]. Over the past few years, numerous Transformer variants have been proposed to advance the state-of-the-art performances of various tasks significantly. There are quite a few literature reviews from different aspects, such as in NLP applications [Qiu *et al.*, 2020; Han *et al.*, 2021], CV applica-

tions [Han *et al.*, 2020; Khan *et al.*, 2021; Selva *et al.*, 2022], efficient Transformers [Tay *et al.*, 2020], and attention models [Chaudhari *et al.*, 2021; Galassi *et al.*, 2020].

Transformers have shown great modeling ability for long-range dependencies and interactions in sequential data and thus are appealing to time series modeling. Many variants of Transformer have been proposed to address special challenges in time series modeling and have been successfully applied to various time series tasks, such as forecasting [Li *et al.*, 2019; Zhou *et al.*, 2021; Zhou *et al.*, 2022], anomaly detection [Xu *et al.*, 2022; Tuli *et al.*, 2022], classification [Zerveas *et al.*, 2021; Yang *et al.*, 2021], and so on. For example, seasonality or periodicity is an important feature of time series [Wen *et al.*, 2021a]. How to effectively model long-range and short-range temporal dependency and capture seasonality simultaneously remains a challenge [Wu *et al.*, 2021; Zhou *et al.*, 2022]. As Transformer for time series is an emerging subject in deep learning, a systematic and comprehensive survey on Transformers in time series would greatly benefit the time series community. We note that there exist several surveys related to deep learning for time series, include forecasting [Lim and Zohren, 2021; Benidis *et al.*, 2020; Torres *et al.*, 2021], classification [Islam Fawaz *et al.*, 2019], anomaly detection [Choi *et al.*, 2021; Blázquez-García *et al.*, 2021], and data augmentation [Wen *et al.*, 2021b], but little was given to Transformers for time series.

In this paper, we aim to fill the gaps by summarizing the main developments of time series Transformers. We first give a brief introduction about vanilla Transformer, and then propose a new taxonomy from perspectives of both network modifications and application domains for time series Transformers. For network modifications, we discuss the improvements made on both low-level (i.e. module) and high-level (i.e. architecture) of Transformers, with the aim to optimize the performance of time series modeling. For applications, we analyze and summarize Transformers for popular time series tasks including forecasting, anomaly detection, and classification. For each time series Transformer, we analyze its insights, strengths and limitations. To provide practical guidelines on how to effectively use Transformers for time series modeling, we conduct extensive empirical studies that examine multiple aspects of time series modeling, including robustness analysis, model size analysis, and seasonal-trend

<sup>1</sup> <https://github.com/qingsongedu/time-series-transformers-review>

decomposition analysis. We conclude this work by discussing possible future directions for time series Transformer, including inductive biases for time series Transformers, Transformers and GNN for time series, pre-trained Transformers for time series, and Transformers with NAS for time series. To the best of our knowledge, this is the first work to comprehensively and systematically review the key developments of Transformers for modeling time series data. We hope this survey will ignite further research interests in time series Transformers.

## 2 Preliminaries of the Transformer

### 2.1 Vanilla Transformer

The vanilla Transformer [Vaswani *et al.*, 2017] follows most competitive neural sequence models with an encoder-decoder structure. Both encoder and decoder are composed of multiple identical blocks. Each encoder block consists of a multi-head self-attention module and a position-wise feed-forward network (FFN) while each decoder block inserts cross-attention models between the multi-head self-attention module and the position-wise feed-forward network (FFN).

### 2.2 Input Encoding and Positional Encoding

Unlike LSTM or RNN, Transformer has no recurrence and no convolution. Instead, it utilizes the positional encoding added in the input embeddings, to model the sequence information. We summarize some positional encodings below.

#### Absolute Positional Encoding

In vanilla Transformer, for each position index  $t$ , encoding vector is given by

$$PE(t)_i = \begin{cases} \sin(\omega_i t) & i \% 2 = 1 \\ \cos(\omega_i t) & i \% 2 = 0 \end{cases} \quad (1)$$

where  $\omega_i$  is the hand-crafted frequency for each dimension. Another way is to learn a set of positional embeddings for each position which is more flexible [Kenton and Toutanova, 2019; Gehring *et al.*, 2017].

#### Relative Positional Encoding

Following the intuition that pairwise positional relationships between input elements is more beneficial than positions of elements, relative positional encoding methods have been proposed. For example, one of such methods is to add a learnable relative positional embedding to keys of attention mechanism [Shaw *et al.*, 2018].

Besides the absolute and relative positional encodings, there are methods using hybrid positional encodings that combine them together [Ke *et al.*, 2020]. Generally, the positional encoding is added to the token embedding and fed to Transformer.

### 2.3 Multi-head Attention

With Query-Key-Value (QKV) model, the scaled dot-product attention used by Transformer is given by

$$Attention(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = softmax(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{D_k}})\mathbf{V} \quad (2)$$

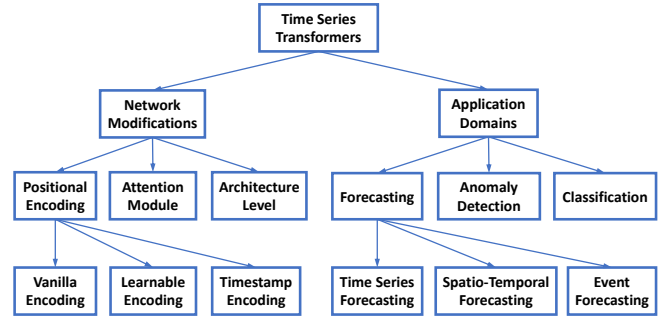


Figure 1: Taxonomy of Transformers for time series modeling from the perspectives of network modifications and application domains.

where queries  $\mathbf{Q} \in \mathcal{R}^{N \times D_k}$ , keys  $\mathbf{K} \in \mathcal{R}^{M \times D_k}$ , values  $\mathbf{V} \in \mathcal{R}^{M \times D_v}$  and  $N, M$  denote the lengths of queries and keys (or values),  $D_k, D_v$  denote the dimensions of keys (or queries) and values. Transformer uses multi-head attention with  $H$  different sets of learned projections instead of a single attention function as

$$MultiHeadAttn(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = Concat(head_1, \dots, head_H)\mathbf{W}^O,$$

where  $head_i = Attention(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V)$ .

### 2.4 Feed-forward and Residual Network

The point-wise feed-forward network is a fully connected module as

$$FFN(\mathbf{H}') = ReLU(\mathbf{H}'\mathbf{W}^1 + \mathbf{b}^1)\mathbf{W}^2 + \mathbf{b}^2, \quad (3)$$

where  $\mathbf{H}'$  is outputs of previous layer,  $\mathbf{W}^1 \in \mathcal{R}^{D_m \times D_f}$ ,  $\mathbf{W}^2 \in \mathcal{R}^{D_f \times D_m}$ ,  $\mathbf{b}^1 \in \mathcal{R}^{D_f}$ ,  $\mathbf{b}^2 \in \mathcal{R}^{D_m}$  are trainable parameters. In a deeper module, a residual connection module followed by Layer Normalization Module is inserted around each module. That is,

$$\mathbf{H}' = LayerNorm(SelfAttn(\mathbf{X}) + \mathbf{X}), \quad (4)$$

$$\mathbf{H} = LayerNorm(FFN(\mathbf{H}') + \mathbf{H}'), \quad (5)$$

where  $SelfAttn(\cdot)$  denotes self attention module and  $LayerNorm(\cdot)$  denotes the layer normal operation.

## 3 Taxonomy of Transformers in Time Series

To summarize the existing time series Transformers, we propose a taxonomy from perspectives of network modifications and application domains as illustrated in Fig. 1. Based on the taxonomy, we review the existing time series Transformers systematically. From the perspective of network modifications, we summarize the changes made on both module level and architecture level of Transformer in order to accommodate special challenges in time series modeling. From the perspective of applications, we classify time series Transformers based on its application tasks including forecasting, anomaly detection, classification and clustering. In the following two sections, we would delve into the existing time series Transformers from these two perspectives.

## 4 Network Modifications for Time Series

### 4.1 Positional Encoding

As Transformers are permutation equivalent and the ordering of time series matters, it is of great importance to encode the positions of input time series into Transformers. A common design is to first encode positional information as vectors and then inject them into the model as an additional input together with the input time series. How to obtain these vectors when modeling time series with Transformers can be divided into three main categories.

**Vanilla Positional Encoding.** A few works [Li *et al.*, 2019] simply introduce vanilla positional encoding (Section 2.2) used in [Vaswani *et al.*, 2017], which is then added to the input time series embeddings and fed to Transformer. Although this plain application can extract some positional information from time series, they were unable to fully exploit the important features of time series data.

**Learnable Positional Encoding.** As the vanilla positional encoding is hand-crafted and less expressive and adaptive, several studies found that learning appropriate positional embeddings from time series data can be much more effective. Compared to the fixed vanilla positional encoding, learned embeddings are more flexible and can adapt itself to specific tasks. [Zerveas *et al.*, 2021] introduce an embedding layer in Transformer that learn embedding vectors for each position index jointly with other model parameters. [Lim *et al.*, 2019] used an LSTM network to encode positional embeddings, with the aim to better exploit sequential ordering information in time series.

**Timestamp Encoding.** When modeling time series in real-world scenarios, the timestamp information is commonly accessible, including calendar timestamps (e.g., second, minute, hour, week, month and year) and special timestamps (e.g., holidays and events). These timestamps are quite informative in real applications but hardly leveraged in vanilla Transformers. To mitigate the issue, Informer [Zhou *et al.*, 2021] proposed to encode timestamps as additional positional encoding by using learnable embedding layers. A similar timestamp encoding scheme was used in Autoformer [Wu *et al.*, 2021] and FEDformer [Zhou *et al.*, 2022].

### 4.2 Attention Module

Central to Transformer is the self-attention module. It can be viewed as a fully connected layer with the weights that are dynamically generated based on the pairwise similarity of input patterns. As a result, it shares the same maximum path length as fully connected layers, but with much less number of parameters, making it suitable for modeling long-term dependencies.

As we show in the previous section the self-attention module in vanilla Transformer has a time and memory complexity of  $O(L^2)$  ( $L$  is the input time series length), which becomes the computational bottleneck when dealing with long sequences. Many efficient Transformers were proposed to reduce the quadratic complexity that can be classified into two main categories: (1) explicitly introducing a sparsity bias into the attention mechanism like LogTrans [Li *et al.*, 2019] and Pyraformer [Liu *et al.*, 2022]; (2) exploring the low-rank

Table 1: Complexity comparisons of popular time series transformers with different attention modules.

Methods	Training		Testing
	Time	Memory	Steps
Transformer [Vaswani <i>et al.</i> , 2017]	$O(L^2)$	$O(L^2)$	$L$
LogTrans [Li <i>et al.</i> , 2019]	$O(L \log L)$	$O(L^2)$	1
Informer [Zhou <i>et al.</i> , 2021]	$O(L \log L)$	$O(L \log L)$	1
Pyraformer [Liu <i>et al.</i> , 2022]	$O(L)$	$O(L)$	1
FEDformer [Zhou <i>et al.</i> , 2022]	$O(L)$	$O(L)$	1

property of self-attention matrix to speed up the computation, e.g. Informer [Zhou *et al.*, 2021] and FEDformer [Zhou *et al.*, 2022]. In Table 1, we summarize both the time and memory complexity of popular Transformers applied to time series modeling.

### 4.3 Architecture-Level Innovation

In addition to accommodate individual modules in Transformers for modeling time series, a number of works [Zhou *et al.*, 2021; Liu *et al.*, 2022] seek to renovate Transformers on the architecture level. Recent works introduce hierarchical architecture into Transformer to take into account multi-resolution aspect of time series. Informer [Zhou *et al.*, 2021] inserts max-pooling layers with stride 2 between attention blocks, which down-sample series into its half slice. Pyraformer [Liu *et al.*, 2022] designs a  $C$ -ary tree based attention mechanism, in which nodes at the finest scale correspond to the original time series, while nodes in the coarser scales represent series at lower resolutions. Pyraformer developed both intra-scale and inter-scale attentions in order to better capture temporal dependencies across different resolutions. Besides the ability of integrating information at different multi-resolutions, a hierarchical architecture also enjoys the benefits of efficient computation, particularly for long time series.

## 5 Applications of Time Series Transformers

In this section, we will review the application of Transformer to important time series tasks, including forecasting, anomaly detection, and classification.

### 5.1 Transformers in Forecasting

We examine three types of forecasting tasks here, i.e. time series forecasting, spatial-temporal forecasting, and event forecasting.

#### Time Series Forecasting

Forecasting is the most common and important application of time series. LogTrans [Li *et al.*, 2019] proposed convolutional self-attention by employing causal convolutions to generate queries and keys in the self-attention layer. It introduces sparse bias, a Logsparse mask, in self-attention model that reduces computational complexity from  $O(L^2)$  to  $O(L \log L)$ .

Instead of explicitly introducing sparse bias, Informer [Zhou *et al.*, 2021] selects  $O(\log L)$  dominant queries based on queries and key similarities, thus achieving similar improvements as LogTrans in computational complexity. It also designs a generative style decoder to produce long term forecasting directly and thus avoids accumulative

error in using one forward step prediction for long term forecasting.

AST [Wu *et al.*, 2020] uses a generative adversarial encoder-decoder framework to train a sparse Transformer model for time series forecasting. It shows that adversarial training can improve the time series forecasting by directly shaping the output distribution of network to avoid the error accumulation through one-step ahead inference.

Autoformer [Wu *et al.*, 2021] devises a simple seasonal-trend decomposition architecture with an auto-correlation mechanism working as an attention module. The auto-correlation block is not a traditional attention block. It measures the time-delay similarity between inputs signal and aggregate the top-k similar sub-series to produce the output with a reduced complexity of  $O(L \log L)$ .

FEDformer [Zhou *et al.*, 2022] applies attention operation in the frequency domain with Fourier transform and wavelet transform. It achieves a linear complexity by randomly selecting a fixed size subset of frequency. It is worth noting that due to the success of Autoformer and FEDformer, it has attracts more attention in the community to explore self-attention mechanism in frequency domain for time series modeling.

TFT [Lim *et al.*, 2021] designs a multi-horizon forecasting model with static covariate encoders, gating feature selection and temporal self-attention decoder. It encodes and selects useful information from various covariates to perform forecasting. It also preserves interpretability incorporating global, temporal dependency and event.

SSDNet [Lin *et al.*, 2021] and ProTran [Tang and Matteson, 2021] combine Transformer with state space models to provide probabilistic forecasts. SSDNet first uses Transformer to learn the temporal pattern and estimate the parameters of SSM, and then applies SSM to perform the seasonal-trend decomposition and maintain the interpretable ability. ProTran designs a generative modeling and inference procedures based on variational inference.

Pyraformer [Liu *et al.*, 2022] designs a hierarchical pyramidal attention module with binary tree following path, to capture temporal dependencies of different ranges with linear time and memory complexity.

Aliformer [Qi *et al.*, 2021] makes sequential forecasting for time series data by using a Knowledge-guided attention with a branch to revise and denoise the the attention map.

### Spatio-Temporal Forecasting

In spatio-temporal forecasting, we need to take into account both temporal and spatio-temporal dependency for accurate forecasting. Traffic Transformer [Cai *et al.*, 2020] designs an encoder-decoder structure using a self attention module to capture temporal-temporal dependencies and a Graph neural network module to capture spatial dependencies. Spatial-temporal Transformer [Xu *et al.*, 2020] for traffic flow forecasting takes a step further. Besides introducing a temporal Transformer block to capture temporal dependencies, it also design a spatial Transformer block, together with a graph convolution network, to better capture spatial-spatial dependencies. Spatio-temporal graph Transformer [Yu *et al.*, 2020] designs a attention-based graph convolution mechanism that is

able to learn a complicated temporal-spatial attentions pattern to improve pedestrian trajectory prediction.

### Event Forecasting

Event sequence data with irregular and asynchronous timestamps are naturally observed in many real-life applications, which is in contrast to regular time series data with equal sampling intervals. Event forecasting or prediction aims to predict the times and marks of future events given the history of past events, and it is often modeled by temporal point processes (TPP) [Shchur *et al.*, 2021].

Recently, several neural TPP models start to incorporate Transformers in order to improve the performance of event prediction. Self-attentive Hawkes process (SAHP) [Zhang *et al.*, 2020] and Transformer Hawkes process (THP) [Zuo *et al.*, 2020] adopt Transformer encoder architecture to summarize the influence of history events and compute the intensity function for event prediction. They modify the positional encoding by translating time intervals into sinusoidal function such that the intervals between events can be utilized. Later, a more flexible named attentive neural Datalog through time (A-NDTT) [Mei *et al.*, 2022] is proposed to extend SAHP/THP schemes by embedding all possible events and times with attention as well. Experiments show that it can better capture sophisticated event dependencies than existing methods.

### 5.2 Transformers in Anomaly Detection

Deep learning also triggers novel developments for anomaly detection [Ruff *et al.*, 2021]. As deep learning is a kind of representation learning, reconstruction model plays a significant role in anomaly detection tasks. Reconstruction model aims to learn a neural network that maps vectors from a simple predefined source distribution  $Q$  to the actual input distribution  $P^+$ .  $Q$  is usually a Gaussian or uniform distribution. Anomaly score is defined by reconstruction error. Intuitively, the higher reconstruction error which means less likely to be from the input distribution, the higher anomaly score. A threshold is set to discriminate anomaly from normality.

Recently, [Meng *et al.*, 2019] revealed the advantage of using Transformer for anomaly detection over other traditional models for temporal dependency (e.g. LSTM ). Besides a higher detection quality (measured by F1), transformer based anomaly detection is significantly more efficient than LSTM based methods, mostly due to parallel computing in Transformer architecture. In multiple studies, including TranAD [Tuli *et al.*, 2022], MT-RVAE [Wang *et al.*, 2022], and TransAnomaly [Zhang *et al.*, 2021], researchers proposed to combine Transformer with neural generative models, such as VAEs [Kingma and Welling, 2013] and GANs [Goodfellow *et al.*, 2014], for better reconstruction models in anomaly detection.

TranAD proposes an adversarial training procedure to amplify reconstruction errors as a simple Transformer-based network tends to miss small deviation of anomaly. GAN style adversarial training procedure is designed by two Transformer encoders and two Transformer decoders to gain stability. Ablation study shows that, if Transformer-based encoder-decoder is replaced, F1 score performance drops nearly 11%,

indicating the significance of Transformer architecture to anomaly detection.

While both MT-RVAE and TransAnomaly combine VAE with Transformer, they share different purposes. TransAnomaly combines VAE with Transformer to allow more parallelization and reduce training cost by nearly 80%. In MT-RVAE, multiscale Transformer is designed to extract and integrate time-series information at different scales. It overcomes the shortcomings of traditional Transformer where only local information is extracted for sequential analysis.

Several time series Transformer are designed for multivariate time series that combine Transformer with graph based learning architecture, such as GTA [Chen *et al.*, 2021d]. Note that, MT-RVAE is also for multivariate time series but with few dimensions or insufficient close relationships among sequences where the graph neural network model does not work well. To deal with such challenge, MT-RVAE modifies positional encoding module and introduces feature-learning module. GTA contains graph convolution structure to model the influence propagation process. Similar to MT-RVAE, GTA also considers “global” information, yet by replacing vanilla multi-head attention with multi-branch attention mechanism, that is, a combination of global-learned attention, vanilla multi-head attention and neighborhood convolution.

AnomalyTrans [Xu *et al.*, 2022] combines Transformer and Gaussian Prior-Association to make rare anomalies more distinguishable. Although sharing similar motivation as TranAD, AnomalyTrans achieves this goal in a very different way. The insight is that it is harder for anomalies to build strong associations with the whole series while easier with adjacent time points compared with normality. In AnomalyTrans, prior-association and series-association are modeled simultaneously. Besides reconstruction loss, anomaly model is optimized by the minimax strategy to constrain the prior- and series- associations for more distinguishable association discrepancy.

### 5.3 Transformers in Classification

Transformer is proved to be effective in various time series classification tasks due to its prominent capability in capturing long-term dependency. Classification Transformers usually employ a simple encoder structure, in which self-attention layers performs representation learning and feed forward layers produce probability of each class.

GTN [Liu *et al.*, 2021] uses a two-tower Transformer with each tower respectively working on time-step-wise attention and channel-wise attention. To merge the feature of the two towers, a learnable weighted concat (also known as ‘gating’) is used. The proposed extension of Transformer achieves state-of-the-art results on 13 multivariate time series classification. [Rußwurm and Körner, 2020] studied the self-attention based Transformer for raw optical satellite time series classification and obtained the best results comparing with recurrent and convolutional neural networks.

Pre-trained Transformers are also investigated in classification tasks. [Yuan and Lin, 2020] studies the Transformer for raw optical satellite image time series classification. The authors use self-supervised pre-trained schema because of lim-

Table 2: The **MSE** comparisons in robustness experiment of forecasting 96 steps for ETTm2 dataset with prolonging input length.

Model	Transformer	Autoformer	Informer	Reformer	LogFormer	
<i>Input Len</i>	96	<b>0.557</b>	<b>0.239</b>	0.428	<b>0.615</b>	<b>0.667</b>
	192	0.710	0.265	<b>0.385</b>	0.686	0.697
	336	1.078	0.375	1.078	1.359	0.937
	720	1.691	0.315	1.057	1.443	2.153
	1440	0.936	0.552	1.898	0.815	0.867

Table 3: The **MSE** comparisons in model size experiment of forecasting 96 steps for ETTm2 dataset with different number of layers.

Model	Transformer	Autoformer	Informer	Reformer	LogFormer	
<i>Layer Num</i>	3	0.557	<b>0.234</b>	0.428	0.597	0.667
	6	<b>0.439</b>	0.282	0.489	<b>0.353</b>	<b>0.387</b>
	12	0.556	0.238	0.779	0.481	0.562
	24	0.580	0.266	0.815	1.109	0.690
	48	0.461	NaN	1.623	OOM	2.992

ited labeled data. [Zerveas *et al.*, 2021] introduced an unsupervised pre-trained framework and the model is pre-trained with proportionally masked data. The pre-trained models are then fine-tuned in downstream tasks such as classification. [Yang *et al.*, 2021] proposes to use large-scale pre-trained speech processing model for downstream time series classification problems and generates 19 competitive results on 30 popular time series classification datasets.

## 6 Experimental Evaluation and Discussion

In this section, we conduct empirical studies to analyze how Transformers work on time series data. Specifically, we test different algorithms with different configurations on a typical benchmark dataset ETTm2 [Zhou *et al.*, 2021].

### Robustness Analysis

A lot of works we describe above carefully design attention module to lower the quadratic calculation and memory complexity, though they practically use a short fixed-size input to achieve the best result in their reported experiments. It makes us question the actual usage of such an efficient design. We perform a robust experiment with prolonging input sequence length to verify their prediction power and robustness when dealing with long-term input sequences.

As shown in Table 2, when we compare the prediction results with prolonging input length, various Transformer-based model deteriorates quickly. This phenomenon makes a lot of carefully designed Transformers impractical in long-term forecasting tasks since they cannot effectively utilize long input information. More works need to be done to fully utilize long sequence input other than simply running it.

### Model Size Analysis

Before being introduced into the field of time series prediction, Transformer has shown dominant performance in NLP and CV community [Vaswani *et al.*, 2017; Kenton and Toutanova, 2019; Qiu *et al.*, 2020; Han *et al.*, 2021; Han *et al.*, 2020; Khan *et al.*, 2021; Selva *et al.*, 2022]. One of the key advantages Transformer holds in these fields is being able to increase prediction power through increasing model size. Usually the model capacity is controlled by Trans-

Table 4: The MSE comparisons in ablation experiments of seasonal-trend decomposition analysis. ‘Ori’ means the original version without the decomposition. ‘Decomp’ means with decomposition. The experiment is performed on ETTm2 dataset with prolonging output length.

Model		FEDformer		Autoformer		Informer		LogTrans		Reformer		Transformer		Promotion
MSE		Ori	Decomp	Ori	Decomp	Ori	Decomp	Ori	Decomp	Ori	Decomp	Ori	Decomp	Relative
Out Len	96	0.457	0.203	0.581	0.255	0.365	0.354	0.768	0.231	0.658	0.218	0.604	0.204	53%
	192	0.841	0.269	1.403	0.281	0.533	0.432	0.989	0.378	1.078	0.336	1.060	0.266	62%
	336	1.451	0.325	2.632	0.339	1.363	0.481	1.334	0.362	1.549	0.366	1.413	0.375	75%
	720	3.282	0.421	3.058	0.422	3.379	0.822	3.048	0.539	2.631	0.502	2.672	0.537	82%

former’s layer number, which is commonly set between 12 to 128 in CV and NLP.

But as shown in our experiments in Table 3, when we compare the prediction result with different Transformer models with various layer numbers, the shallowest Transformer with 3 to 6 layers triumphs. It raises a question about how to design a proper Transformer architecture with deep layers to increase the model’s capacity and achieve better forecasting performance.

### Seasonal-Trend Decomposition Analysis

In the latest studies, researchers [Wu *et al.*, 2021; Zhou *et al.*, 2022; Lin *et al.*, 2021; Liu *et al.*, 2022] begin to realize that the seasonal-trend decomposition is a crucial part for Transformer’s performance in time series forecasting. As a simple experiment shown in table 4, we use a moving average trend decomposition architecture proposed in [Wu *et al.*, 2021] to test various attention modules. A seasonal-trend decomposition model can significantly boost model’s performance by 50 % to 80%. It is a unique block and such performance boosting through decomposition seems a consistent phenomenon in time series forecasting for Transformer’s application, which worth further investigating.

## 7 Future Research Opportunities

Here we highlight a few directions that are potentially promising for the research of Transformers in time series.

### 7.1 Inductive Biases for Time Series Transformers

Vanilla Transformer does not make any assumptions about data patterns and characteristics. Although it is a general and universal network for modeling long-range dependencies, it also come with a price, i.e. lots of data are needed to train Transformer to avoid data overfitting. One of the key features of time series data is its seasonal/periodic and trend patterns [Wen *et al.*, 2019; Cleveland *et al.*, 1990]. Some recent studies have shown that incorporating series periodicity [Wu *et al.*, 2021] or frequency processing [Zhou *et al.*, 2022] into time series Transformer can enhance performance significantly. Thus, one future direction is to consider more effective ways to induce inductive biases into Transformers based on the understanding of time series data and characteristics of specific tasks.

### 7.2 Transformers and GNN for Time Series

Multivariate and spatio-temporal time series are becoming more and more common in applications, calling for additional techniques to handle high dimensionality, especially the ability to capture the underlying relationships among dimensions.

Introducing graph neural networks (GNNs) is a natural way to model spatial dependency or relationships among dimensions. Recently, several studies have demonstrated that the combination of GNN and transformers/attentions could bring not only significant performance improvement like in traffic forecasting [Cai *et al.*, 2020; Xu *et al.*, 2020] and multi-modal forecasting [Li *et al.*, 2021], but also better understanding of the spatio-temporal dynamics and latent causality. It is an important future direction to combine Transformers and GNNs for effectively spatial-temporal modeling in time series.

### 7.3 Pre-trained Transformers for Time Series

Large-scale pre-trained Transformer models have significantly boosted the performance for various tasks in NLP [Kenton and Toutanova, 2019; Brown *et al.*, 2020] and CV [Chen *et al.*, 2021a]. However, there are limited works on pre-trained Transformers for time series, and existing studies mainly focus on time series classification [Zerveas *et al.*, 2021; Yang *et al.*, 2021]. Therefore, how to develop appropriate pre-trained Transformer models for different tasks in time series remains to be examined in the future.

### 7.4 Transformers with NAS for Time Series

Hyper-parameters, such as embedding dimension, number of heads, and number of layers, can largely affect the performance of Transformer. Manual configuring these hyper-parameters is time-consuming and often results in suboptimal performance. Neural architecture search (NAS) [Elsken *et al.*, 2019; Wang *et al.*, 2020] has been a popular technique for discovering effective deep neural architectures, and automating Transformer design using NAS in NLP and CV can be found in recent studies [So *et al.*, 2019; Chen *et al.*, 2021c]. For industry-scale time series data which can be of both high-dimension and long-length, automatically discovering both memory- and computational-efficient Transformer architectures is of practical importance, making it an important future direction for time series Transformer.

## 8 Conclusion

In this paper, we provide a comprehensive survey on time series Transformers in various tasks. We organize the reviewed methods in a new taxonomy consisting of network modifications and applications domains. We summarize representative methods in each category, discuss their strengths and limitations by experimental evaluation, and highlight future research directions.

## References

- [Benidis *et al.*, 2020] Konstantinos Benidis, Syama Sundar Rangapuram, Valentin Flunkert, Bernie Wang, Danielle Maddix, Caner Turkmen, Jan Gasthaus, et al. Neural forecasting: Introduction and literature overview. *arXiv preprint arXiv:2004.10240*, 2020.
- [Blázquez-García *et al.*, 2021] Ane Blázquez-García, Angel Conde, Usue Mori, and Jose A Lozano. A review on outlier/anomaly detection in time series data. *ACM Computing Surveys (CSUR)*, 54(3):1–33, 2021.
- [Brown *et al.*, 2020] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, et al. Language models are few-shot learners. *NeurIPS*, 2020.
- [Cai *et al.*, 2020] Ling Cai, Krzysztof Janowicz, Gengchen Mai, Bo Yan, and Rui Zhu. Traffic transformer: Capturing the continuity and periodicity of time series for traffic forecasting. *Transactions in GIS*, 24(3):736–755, 2020.
- [Chaudhari *et al.*, 2021] Sneha Chaudhari, Varun Mithal, Gungor Polatkan, and Rohan Ramanath. An attentive survey of attention models. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 12(5):1–32, 2021.
- [Chen *et al.*, 2021a] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, et al. Pre-trained image processing transformer. In *CVPR*, 2021.
- [Chen *et al.*, 2021b] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *NeurIPS*, 2021.
- [Chen *et al.*, 2021c] Minghao Chen, Houwen Peng, Jianlong Fu, and Haibin Ling. AutoFormer: Searching transformers for visual recognition. In *CVPR*, 2021.
- [Chen *et al.*, 2021d] Zekai Chen, Dingshuo Chen, Xiao Zhang, Zixuan Yuan, and Xiuzhen Cheng. Learning graph structures with transformer for multivariate time series anomaly detection in IoT. *IEEE Internet of Things Journal*, 2021.
- [Choi *et al.*, 2021] Kukjin Choi, Jihun Yi, Changhwa Park, and Sungroh Yoon. Deep learning for anomaly detection in time-series data: Review, analysis, and guidelines. *IEEE Access*, 2021.
- [Cleveland *et al.*, 1990] Robert B Cleveland, William S Cleveland, Jean E McRae, and Irma Terpenning. STL: A seasonal-trend decomposition procedure based on loess. *Journal of Official Statistics*, 6(1):3–73, 1990.
- [Dong *et al.*, 2018] Linhao Dong, Shuang Xu, and Bo Xu. Speech-transformer: a no-recurrence sequence-to-sequence model for speech recognition. In *ICASSP*, 2018.
- [Dosovitskiy *et al.*, 2021] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.
- [Elsken *et al.*, 2019] Elsken, Thomas, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *Journal of Machine Learning Research*, 2019.
- [Galassi *et al.*, 2020] Andrea Galassi, Marco Lippi, and Paolo Torroni. Attention in natural language processing. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [Gehring *et al.*, 2017] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional sequence to sequence learning. In *ICML*, 2017.
- [Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 2014.
- [Han *et al.*, 2020] Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on visual transformer. *arXiv preprint arXiv:2012.12556*, 2020.
- [Han *et al.*, 2021] Xu Han, Zhengyan Zhang, Ning Ding, Yuxian Gu, Xiao Liu, Yuqi Huo, Jiezhong Qiu, Liang Zhang, Wentao Han, Minlie Huang, et al. Pre-trained models: Past, present and future. *AI Open*, 2021.
- [Ismail Fawaz *et al.*, 2019] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. *Data mining and knowledge discovery*, 2019.
- [Ke *et al.*, 2020] Guolin Ke, Di He, and Tie-Yan Liu. Rethinking positional encoding in language pre-training. *arXiv preprint arXiv:2006.15595*, 2020.
- [Kenton and Toutanova, 2019] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- [Khan *et al.*, 2021] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *arXiv preprint arXiv:2101.01169*, 2021.
- [Kingma and Welling, 2013] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [Li *et al.*, 2019] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyu Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. In *NeurIPS*, 2019.
- [Li *et al.*, 2021] Longyuan Li, Jian Yao, Li Wenliang, Tong He, Tianjun Xiao, Junchi Yan, David Wipf, and Zheng Zhang. Grin: Generative relation and intention network for multi-agent trajectory prediction. In *NeurIPS*, 2021.
- [Lim and Zohren, 2021] Bryan Lim and Stefan Zohren. Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society*, 2021.
- [Lim *et al.*, 2019] Bryan Lim, Sercan O Arik, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *arXiv preprint arXiv:1912.09363*, 2019.
- [Lim *et al.*, 2021] Bryan Lim, Sercan Ö Arik, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4):1748–1764, 2021.
- [Lin *et al.*, 2021] Yang Lin, Irena Koprinska, and Mashud Rana. SSDNet: State space decomposition neural network for time series forecasting. In *ICDM*, 2021.
- [Liu *et al.*, 2021] Minghao Liu, Shengqi Ren, Siyuan Ma, Jiahui Jiao, Yizhou Chen, Zhiguang Wang, and Wei Song. Gated transformer networks for multivariate time series classification. *arXiv preprint arXiv:2103.14438*, 2021.
- [Liu *et al.*, 2022] Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Shahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *ICLR*, 2022.



- [Mei et al., 2022] Hongyuan Mei, Chenghao Yang, and Jason Eisner. Transformer embeddings of irregularly spaced events and their participants. In *ICLR*, 2022.
- [Meng et al., 2019] Hengyu Meng, Yuxuan Zhang, Yuanxiang Li, and Honghua Zhao. Spacecraft anomaly detection via transformer reconstruction error. In *ICASSE*, 2019.
- [Qi et al., 2021] Xinyuan Qi, Kai Hou, Tong Liu, Zhongzhong Yu, Sihao Hu, and Wenwu Ou. From known to unknown: Knowledge-guided transformer for time-series sales forecasting in Alibaba. *arXiv preprint arXiv:2109.08381*, 2021.
- [Qiu et al., 2020] Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. Pre-trained models for natural language processing: A survey. *Science China Technological Sciences*, pages 1–26, 2020.
- [Ruff et al., 2021] Lukas Ruff, Jacob R Kauffmann, Robert A Van dermeulen, Grégoire Montavon, Wojciech Samek, Marius Kloft, Thomas G Dietterich, and Klaus-Robert Müller. A unifying review of deep and shallow anomaly detection. *Proceedings of the IEEE*, 2021.
- [Rußwurm and Körner, 2020] Marc Rußwurm and Marco Körner. Self-attention for raw optical satellite time series classification. *ISPRS Journal of Photogrammetry and Remote Sensing*, 169:421–435, 11 2020.
- [Selva et al., 2022] Javier Selva, Anders S Johansen, Sergio Escalera, Kamal Nasrollahi, Thomas B Moeslund, and Albert Clapés. Video transformers: A survey. *arXiv preprint arXiv:2201.05991*, 2022.
- [Shaw et al., 2018] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.
- [Shchur et al., 2021] Oleksandr Shchur, Ali Caner Türkmen, Tim Januschowski, and Stephan Günnemann. Neural temporal point processes: A review. In *IJCAI*, 2021.
- [So et al., 2019] David So, Quoc Le, and Chen Liang. The evolved transformer. In *ICML*, 2019.
- [Tang and Matteson, 2021] Binh Tang and David Matteson. Probabilistic transformer for time series analysis. In *NeurIPS*, 2021.
- [Tay et al., 2020] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. *arXiv preprint arXiv:2009.06732*, 2020.
- [Torres et al., 2021] José F. Torres, Dalil Hadjout, Abderrazak Sebaa, Francisco Martínez-Álvarez, and Alicia Troncoso. Deep learning for time series forecasting: a survey. *Big Data*, 2021.
- [Tuli et al., 2022] Shreshth Tuli, Giuliano Casale, and Nicholas R Jennings. TranAD: Deep transformer networks for anomaly detection in multivariate time series data. In *VLDB*, 2022.
- [Vaswani et al., 2017] 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*, 2017.
- [Wang et al., 2020] Xiaoxing Wang, Chao Xue, Junchi Yan, Xiaokang Yang, Yonggang Hu, and Kewei Sun. MergeNAS: Merge operations into one for differentiable architecture search. In *IJCAI*, 2020.
- [Wang et al., 2022] Xixuan Wang, Dechang Pi, Xiangyan Zhang, Hao Liu, and Chang Guo. Variational transformer-based anomaly detection approach for multivariate time series. *Measurement*, page 110791, 2022.
- [Wen et al., 2019] Qingsong Wen, Jingkun Gao, Xiaomin Song, Liang Sun, Huan Xu, and Shenghuo Zhu. RobustSTL: A robust seasonal-trend decomposition algorithm for long time series. In *AAAI*, 2019.
- [Wen et al., 2021a] Qingsong Wen, Kai He, Liang Sun, Yingying Zhang, Min Ke, and Huan Xu. RobustPeriod: Time-frequency mining for robust multiple periodicities detection. In *SIGMOD*, 2021.
- [Wen et al., 2021b] Qingsong Wen, Liang Sun, Fan Yang, Xiaomin Song, Jingkun Gao, Xue Wang, and Huan Xu. Time series data augmentation for deep learning: A survey. In *IJCAI*, 2021.
- [Wu et al., 2020] Sifan Wu, Xi Xiao, Qianggang Ding, Peilin Zhao, Ying Wei, and Junzhou Huang. Adversarial sparse transformer for time series forecasting. In *NeurIPS*, 2020.
- [Wu et al., 2021] Haixu Wu, Jiehui Xu, Jianmin Wang, and Ming-sheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In *NeurIPS*, 2021.
- [Xu et al., 2020] 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.
- [Xu et al., 2022] Jiehui Xu, Haixu Wu, Jianmin Wang, and Ming-sheng Long. Anomaly Transformer: Time series anomaly detection with association discrepancy. In *ICLR*, 2022.
- [Yang et al., 2021] Chao-Han Huck Yang, Yun-Yun Tsai, and Pin-Yu Chen. Voice2series: Reprogramming acoustic models for time series classification. In *ICML*, 2021.
- [Yu et al., 2020] Cunjun Yu, Xiao Ma, Jiawei Ren, Haiyu Zhao, and Shuai Yi. Spatio-temporal graph transformer networks for pedestrian trajectory prediction. In *ECCV*, 2020.
- [Yuan and Lin, 2020] Yuan Yuan and Lei Lin. Self-supervised pre-training of transformers for satellite image time series classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:474–487, 2020.
- [Zerveas et al., 2021] George Zerveas, Srideepika Jayaraman, Dhaval Patel, Anuradha Bhamidipaty, and Carsten Eickhoff. A transformer-based framework for multivariate time series representation learning. In *KDD*, 2021.
- [Zhang et al., 2020] Qiang Zhang, Aldo Lipani, Omer Kirnap, and Emine Yilmaz. Self-attentive Hawkes process. In *ICML*, 2020.
- [Zhang et al., 2021] Hongwei Zhang, Yuanqing Xia, Tijin Yan, and Guiyang Liu. Unsupervised anomaly detection in multivariate time series through transformer-based variational autoencoder. In *CCDC*, 2021.
- [Zhou et al., 2021] 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*, 2021.
- [Zhou et al., 2022] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. *arXiv preprint arXiv:2201.12740*, 2022.
- [Zuo et al., 2020] Simiao Zuo, Haoming Jiang, Zichong Li, Tuo Zhao, and Hongyuan Zha. Transformer Hawkes process. In *ICML*, 2020.