# **MoFo: Empowering Long-term Time Series Forecasting with Periodic Pattern Modeling**

Jiaming Ma<sup>1</sup>, Binwu Wang<sup>1,2,\*</sup>, Qihe Huang<sup>1</sup>, Guanjun Wang<sup>1</sup> Pengkun Wang<sup>1,2</sup>, Zhengyang Zhou<sup>1,2</sup>, Yang Wang<sup>1,2,\*</sup>

<sup>1</sup>University of Science and Technology of China (USTC), Hefei, Anhui, China <sup>2</sup>Suzhou Institute for Advanced Research, USTC, Suzhou, Jiangsu, China {JiamingMa, hqh, always}@mail.ustc.edu.cn {wbw2024, pengkun, zzy0929, angyan}@ustc.edu.cn

### **Abstract**

The stable periodic patterns present in the time series data serve as the foundation for long-term forecasting. However, existing models suffer from limitations such as continuous and chaotic input partitioning, as well as weak inductive biases, which restrict their ability to capture such recurring structures. In this paper, we propose MoFo, which interprets periodicity as both the correlation of periodaligned time steps and the trend of period-offset time steps. We first design periodstructured patches—2D tensors generated through discrete sampling—where each row contains only period-aligned time steps, enabling direct modeling of periodic correlations. Period-offset time steps within a period are aligned in columns. To capture trends across these offset time steps, we introduce a period-aware modulator. This modulator introduces an adaptive strong inductive bias through a regulated relaxation function, encouraging the model to generate attention coefficients that align with periodic trends. This function is end-to-end trainable, enabling the model to adaptively capture the distinct periodic patterns across diverse datasets. Extensive empirical results on widely used benchmark datasets demonstrate that MoFo achieves competitive performance while maintaining high memory efficiency and fast training speed. Our code is available at official repository.

# 1 Introduction

Long-term time series forecasting (LTSF) has found widespread applications across various domains [10, 12, 13, 36, 62, 78, 79], with its core challenge lying in understanding and modeling the inherent periodic patterns present within data [20]. To address this, various cutting-edge models have been proposed, among which Transformer [43] has emerged as the de facto backbone for capturing long-range dependencies in LTSF tasks [6, 24, 62, 72]. Despite the promising results achieved, we identify two underexplored potentials that remain to be fully harnessed.

• Continuous but Chaotic Time Steps of Input. Existing models often input consecutive time steps, and popular patch-based methods are no exception. Specifically, the patching method partitions the input sequence in a continuous manner using either convolutional down sampling [22] or sliding window strategies [5, 32, 75], which we refer to as the continuous patch. Taking the Electricity dataset as an example, as shown in the red box in Figure 1(a), a consecutive patch may contain both time steps that are phase-aligned across each period (referred to as period-aligned) and those that are misaligned within each period (period-offset). We further visualize the pairwise correlations between

<sup>\*</sup>Binwu Wang and Yang Wang are corresponding authors.

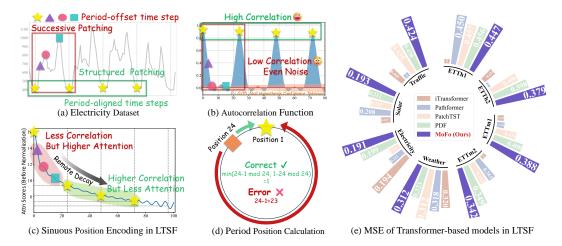


Figure 1: Case visualization and analysis on real-world datasets. (a) Continuous patching strategy. (b) Correlation coefficients within patches. The orange region represents the 95% confidence interval for the null hypothesis from Bartlett's Test [2]. Correlation coefficients that fall within this interval are not statistically significant and cannot be rejected as noise [4]. (c) Attention weights exhibit remote attenuation due to sinuous positional encoding. (d) Absolute and relative position distance in the period perspective. (e) Performance of Transformer-based models for LTSF with 720 horizon.

the first time step in the patch and subsequent ones in Figure 1(b). It is evident that period-aligned time steps, such as 8 AM every day, exhibit strong correlations. In contrast, period-offset time steps rapidly decaying correlations and may even introduce noise. Therefore, the internal correlations within a continuous patch are chaotic, which hampers the model's ability to learn periodic patterns.

Weak Inductive Bias for Periodicity. Due to the permutation invariance of the self-attention mechanism, sinusoidal positional encodings [16, 32] or timestamp position embeddings [62, 77, 79] are commonly used to inject temporal position information. However, as shown in Figure 1(c), sinusoidal encodings may cause attention weights to decay gradually with increasing position distance [50], which conflicts with the repetitive nature of periodic signals. In contrast, timestamp position embeddings encode absolute time intervals, leading to large distances between strongly correlated and phase-aligned time steps. For example, in a daily-period electricity dataset, 8 AM on Monday and 7 AM on Tuesday may exhibit high correlation due to their similar daily phases. However, as shown in Figure 1(d), their absolute time difference is encoded as 23, whereas a more reasonable alternative — their relative cyclical distance — is only 1. Large temporal distances may mislead the model into assigning diminished attention weights. Such weak inductive bias toward periodicity in Transformer further hampers the ability to capture periodic patterns from consecutive and chaotic input.

To address these limitations, we propose MoFo, a novel TransFormer architecture with Modulator that explicitly models periodic patterns by capturing both the correlations among period-aligned time steps and the trends across period-offset ones. Our approach introduces two key components: Period-structured Patching and the Period-Aware Modulator. The former discretely samples time steps to arrange the input into a 2D tensor where each row represents a patch comprising only period-aligned time steps. Period-offset time steps within the period are realigned across columns (i.e., across patches). By modeling the features within each patch, the model can directly learn the underlying periodic correlations. To capture cross-patch periodic trends, we introduce the Period-Aware Modulator, which generates an attention modulation term through a regulated regulated relaxation function. This function assigns attention coefficients based on the periodic relative distances between time steps, encouraging the resulting attention scores to align with the underlying periodic patterns. As a result, the attention mechanism is infused with strong inductive biases that favor periodic dependencies. Crucially, this function is end-to-end trainable, enabling it to dynamically adapt the modulation behavior to the specific periodic characteristics of the input data.

Our contributions are fourfold: • We propose a novel perspective for periodic pattern modeling, termed MoFo, which integrates two innovative strategies for explicit periodicity-aware time series modeling. • We design a Period-Structured Patching strategy that separately manages period-

# 2 Related Work

In recent years, deep learning methods have achieved remarkable success across a broad spectrum of time series tasks, such as forecasting [17, 26, 31, 54, 64, 65, 67, 71], imputation [52], classification [25, 58], and anomaly detection [35, 66, 70]. Among these, Long-Term Series Forecasting (LTSF) has attracted especially intense interest owing to its foundational role in both academic research and practical applications [30, 51, 55]. Existing deep learning approaches for LTSF can be broadly grouped into four categories according to how they handle the temporal dimension: RNN-based, TCN-based, MLP-based (discussed in Appendix B), and Transformer-based models. The Transformer architecture has emerged as the prevailing choice for LTSF, primarily because of its strong ability to capture long-range temporal dependencies. Recent studies have largely concentrated on enhancing both its computational efficiency and modeling power [28, 29], with significant advances driven by two core elements: the self-attention mechanism and patch-based modeling strategies.

Variants of Self-Attention Mechanism for Time Series Forecasting. One of the most interesting aspects of Transformer in LTSF is its self-attention mechanism. LogTrans [16] introducing a convolutional LogSparse attention that ensured long-distance interactions while reducing the number of interactions, is an early and influential attempt to apply Transformer. Informer [77] proposes ProbSparse attention while combining with a distillation mechanism to select the most representative query vectors to compute the attention scores. Autoformer [62] utilizes a decomposition architecture to discover dependencies for building sequence-level connections based on their aggregation of similar subsequences. FEDformer [79] leverages the attention mechanism in the frequency domain, providing the capture of the underlying oscillatory modes and their intensities. Despite the efficacy of these methodologies in reducing computational expense, they often encounter information bottlenecks in long sequence input [33].

Patch-based Time Series Forecasting Method. Patch-based approaches are widely adopted for efficient time series representation, where the input sequence is divided into consecutive segments (patches). This strategy not only enhances the computational efficiency of the Transformer backbone but also promotes more effective modeling of localized temporal dynamics within each patch [11, 60]. PatchTST [32] explicitly reduces the length of input sequence and preserves the local semantic information by dividing the time series into smaller subsequences. Pyraformer [22] employs a pyramidal attention mechanism, wherein the input sequence is downsampled into patches by multilayer convolution, allowing attention to be applied at coarse scales. Crossformer [75] proposes two-stage attention with a dynamic routing mechanism that performs the patch operation from the sequence dimensionality and the channel dimensionality. Pathformer [5] selects multiple patch volumes at the same time and uses a multi-scale router to determine the interactions between different patches to weaken the impact of the selection of patch volume. All existing patch strategies are sequential along the temporal dimensionality, which does not facilitate effective capture and learning of periodic patterns [11, 60].

# 3 MoFo

Given  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T]^{\top} \in \mathbb{R}^T$  with the look-back window T where  $\mathbf{x}_t \in \mathbb{R}$  represents the time series data on time step t, the objective of LTSF is to forecast the next L values  $\mathbf{Y} = [\mathbf{x}_{T+1}, \mathbf{x}_{T+2}, \dots, \mathbf{x}_{T+L}]^{\top} \in \mathbb{R}^L$ . Effective long-term forecasting hinges on the ability to model the intrinsic periodic patterns underlying the data. In this work, as shown in Figure 2, we propose MoFo which integrates two key components: Period-structured Patch and the Period-Aware Modulator, which jointly enhance the model's capacity to capture and exploit periodic patterns.

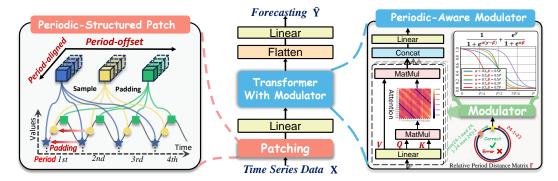


Figure 2: The details of MoFo which contains two core contributions for accurate modeling of periodic patterns: Period-structured Patch and Period-Aware Modulator.

### 3.1 Period-structured Patch in MoFo

Patch-based input representation method has become a popular approach to improve the efficiency of long-time-series modeling of Transformer. However, as discussed in the introduction section, the pairwise correlations between time steps within a continuously partitioned patch can be weak or even negative, which may undermine the model's ability to capture periodic patterns. To address this challenge, we propose a Period-structured Patching strategy. This approach discretely samples period-aligned time steps from the input sequence to form patches. Specifically, given a period length P, which is typically much smaller than the input sequence length T (i.e.,  $P \ll T$ ) in LTSF settings, we compute  $N_P = \lceil T/P \rceil$  as the number of period-length segments required to cover the input time series  $\mathbf{X}$ . Since the length T of the input time series may not be an integer multiple of the period length P, it becomes necessary to apply padding so that all time segments have uniform length without discarding the trailing portion of the input.

Input Padding. Our proposed padding strategy fills incomplete periodic segments with data from adjacent periods. Specifically, we start from the current time step and move backward to delineate periods of length P, and if necessary, we prepend the input time series with the first  $P-(T \bmod P)$  time steps of the first complete period, as shown in Figure 2. This ensures a seamless continuation of the sequence while retaining its underlying periodic structure. As a result, the input series is extended to  $\mathbf{X}_{pad} \in \mathbb{R}^{T'}$  with T'=P\*[T/P], and the padded series can be formally expressed as:

$$\mathbf{X}_{pad} = \begin{cases} \operatorname{Concat}\left(\mathbf{X}_{(T \bmod P):P}, \mathbf{X}\right), & \text{if } T \bmod P > 0, \\ \mathbf{X}, & \text{if } T \bmod P = 0. \end{cases}$$
 (1)

Sampling Patch and Unflatten. We sample time steps at periodic intervals (i.e., period-aligned time steps) from  $\mathbf{X}_{pad}$  and group them into the same patch. For example, for i-th patch, it can be denoted as  $\overline{\mathbf{X}}^i = [x_i, x_{i+P}, \cdots, x_{i+P*\lceil T/P \rceil}] \in \mathbb{R}^{\lceil T/P \rceil}$ , where  $x_{i+P}$  means the data point at the time step (i+P) in  $\mathbf{X}_{pad}$ . Then we unflatten  $\overline{\mathbf{X}}$  to generate the patch-structure input  $\mathbf{X}_{in}$ ,

$$\mathbf{X}_{in} = \text{Unflatten}\left(\overline{\mathbf{X}}\right) \in \mathbb{R}^{P \times \lceil T/P \rceil}.$$
 (2)

where P is the number of patches (also equal to the period length). As illustrated in Figure 2, the structured input  $\mathbf{X}_{in}$  possesses two key properties:  $\mathbf{0}$  Each row (within a patch) contains  $\lceil T/P \rceil$  time steps, with all steps period-aligned within their respective periods. This alignment establishes structured temporal dependencies across periods, enabling the explicit modeling of long-range periodic dependencies.  $\mathbf{0}$  Each column (across patches) contains P time steps from a complete period. By capturing correlations among these patches, the model can effectively learn the underlying periodic trends across period-offset time steps.

**Periodic Dependency Modeling.** Due to this desirable property, even a simple neural architecture—such as a single-layer MLP—can suffice for capturing the underlying periodic dependency. The design is formally described as follows:

$$\mathbf{Z} = \mathbf{X}_{in} \mathbf{W}_{in} + \mathbf{b}_{in} \in \mathbb{R}^{P \times d}, \tag{3}$$

where  $\mathbf{W}_{in} \in \mathbb{R}^{\lceil T/P \rceil \times d}$  and  $\mathbf{b}_{in} \in \mathbb{R}^d$  are learnable parameters. And  $\mathbf{Z}$  is the corresponding output.

### 3.2 Period-Aware Modulator for Periodic Trend Modeling

MoFo further incorporates an enhanced Transformer to model the periodic trends across period-offset time steps within a periods. To further strengthen the model's representational capability, we first modify the standard Transformer architecture (details in Appendix C). Moreover, we introduce a Period-Aware Modulator, which integrates strong inductive biases to address the permutation invariance of self-attention, thereby enhancing the model's ability to capture period trends.

### 3.2.1 Period-Aware Modulator

When computing the attention scores, we explicitly introduce a strong inductive bias: the attention generated by the Transformer is encouraged to align with periodic trends. We achieve this by incorporating a modulation term into the attention computation. First, we design periodic positional encodings that effectively capture the relative distances within a period — a key distinction from existing approaches such as sinusoidal or timestamp-based positional encodings. Specifically, the relative distance between the i-th period and the j-th period is computed as follows,

$$\gamma_{ij} = \min\{(i-j) \bmod P, (j-i) \bmod P\} \in [0, \lfloor P/2 \rfloor], \tag{4}$$

this strategy ensures that time steps that are periodically close remain proximal. For example, in a traffic dataset with an hourly sampling frequency and a daily period (i.e., 24 time steps), the relative period distance between 8:00 AM on Monday and 7:00 AM on Tuesday is 1, while their absolute distance is 23. This allows the model to better perceive and learn periodic trends. and we can get a relative periodic distance matrix, which is denoted as  $\Gamma = \{\gamma_{ij}\}_{i,j=1}^P \in \mathbb{R}^{P \times P}$ .

We proceed by introducing a modulation term within the attention computation, designed to encourage the model to prioritize temporal dependencies that align with the underlying cyclical trend of the data. Specifically, we initially define a modulation matrix  $\mathbf{M} \in \{0,1\}^{P \times P}$ , where each entry  $\mathbf{M}_{ij}$  is an indicator that reveals whether time steps i and j are close in terms of cyclical distance. A natural way to construct this matrix is by using the Heaviside Step function  $\mathcal{H} : \mathbb{R} \to \{0,1\} : \mathbf{M}_{ij} = \mathcal{H}(\beta - \gamma_{ij})$ , and its logarithmic value can be defined as:

$$\mathbf{M}_{ij} = \begin{cases} 1, & \text{if } \gamma_{ij} \le \beta, \\ 0, & \text{if } \gamma_{ij} > \beta, \end{cases} \iff \log \mathbf{M}_{ij} = \begin{cases} 0, & \text{if } \gamma_{ij} \le \beta, \\ -\infty, & \text{if } \gamma_{ij} > \beta. \end{cases}$$
 (5)

where  $\gamma_{ij}$  is the periodic distance between *i*-th period and *j*-th period calculated above. And  $\beta \geq 0$  is a distance penalty threshold controlling the penalization distance. Hence, the attention coefficient between the *i*-th patch and the *j*-th patch is satisfying,

$$\exp\left(\frac{\mathbf{Q}_{i}\mathbf{K}_{j}^{\top}}{\sqrt{d_{h}}} + \log \mathbf{M}_{ij}\right) = \begin{cases} \exp\left(\frac{\mathbf{Q}_{i}\mathbf{K}_{j}^{\top}}{\sqrt{d_{h}}}\right), & \text{if } \gamma_{ij} \leq \beta, \\ 0, & \text{if } \gamma_{ij} > \beta. \end{cases}$$
(6)

The attention score is considered valid only if the cyclical distance between time step j and i is no more than  $\beta$ . Otherwise, the attention coefficient is penalized due to the distance and set to zero. Finally, our attention mechanism can be written as,

Attention 
$$(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_h}} + \log \mathbf{M})\mathbf{V},$$
 (7)

where 
$$\mathbf{Q} = \mathbf{Z}\mathbf{W}_Q^j \in \mathbb{R}^{P \times d_h}, \mathbf{K} = \mathbf{Z}\mathbf{W}_K^j \in \mathbb{R}^{P \times d_h}, \mathbf{V} = \mathbf{Z}\mathbf{W}_V^j \in \mathbb{R}^{P \times d_h},$$
 (8)

where  $\mathbf{W}_Q^j, \mathbf{W}_K^j,$  and  $\mathbf{W}_V^j \in \mathbb{R}^{d imes d_h}$  are learnable parameters.

### 3.2.2 Regulated Relaxation Function for Smooth Approximation

There are still two points with the potential for improvement: ① Due to the non-differentiability of the discrete truncation operation in the step function,  $\beta$  cannot be optimized through backpropagation, which limits the flexibility and adaptability of the method. ② Although the modulated attention mechanism implicitly incorporates temporal positional information by preserving attention coefficients only between time steps with small periodic distances, it imposes uniform inductive biases on nearby steps. As a result, the model still faces challenges related to permutation invariance [15].

To address these limitations, we proposed a regulated relaxation function to approximate the Heaviside Step function to generate the modulator. In contrast to the sharp Heaviside Step function, our function S exhibits a smooth attenuation trend, as shown in Figure 2, which is defined as follows:

# Theorem 1. Regulated Relaxation Function

Define a continuous differentiable function  $S(\cdot; \alpha, \beta) : \mathbb{R}^+ \cup \{0\} \to [0, 1]$  as follows,

$$S(\gamma; \alpha, \beta) = \frac{1}{1 + \exp(\alpha(\gamma - \beta))} + \frac{\exp(-\gamma)}{1 + \exp(\alpha\beta)} \in [0, 1].$$
 (9)

where the regulated parameter  $\alpha > 0$  control the gradient of attenuation and  $\beta > 0$  is the distance penalty threshold. This function has following properties:

(1)  $S(\gamma; \alpha, \beta)$  is the smooth approximation of  $\mathcal{H}(\beta - \gamma_{ij})$  for arbitrary  $\gamma \geq 0$  satisfies,

$$\mathcal{S}(0; \alpha, \beta) = 1, \quad \mathcal{S}(+\infty; \alpha, \beta) = 0, \quad \forall \alpha, \beta > 0.$$
 (10)

(2) The cumulative error upper bound of this smooth approximation satisfies,

$$\int_0^{+\infty} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma < \frac{2\log 2}{\alpha} + \frac{1}{1 + \exp \alpha} \to 0^+ \quad (\alpha \to +\infty). \quad (11)$$

where the proof of Theorem 1 is provided in Appendix A.1. Regulated relaxation function takes relative periodic distance matrix  $\Gamma = \{\gamma_{ij}\}_{i,j=1}^P \in \mathbb{R}^{P \times P}$  to generate the attention modulation term. Finally, the formula of our attention with **R**egulated **R**elaxation **F**unction (RRF) is as follows,

RRF 
$$(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \operatorname{Softmax}(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_h}} + \log \mathcal{S}(\mathbf{\Gamma}; \alpha, \beta))\mathbf{V} \in \mathbb{R}^{P \times d_h},$$
 (12)

where 
$$\mathbf{Q} = \mathbf{Z}\mathbf{W}_{Q}^{j} \in \mathbb{R}^{P \times d_{h}}, \mathbf{K} = \mathbf{Z}\mathbf{W}_{K}^{j} \in \mathbb{R}^{P \times d_{h}}, \mathbf{V} = \mathbf{Z}\mathbf{W}_{V}^{j} \in \mathbb{R}^{P \times d_{h}},$$
 (13)

where  $\alpha$  and  $\beta$  are learnable parameters, and  $\mathbf{Q}$ ,  $\mathbf{K}$ , and  $\mathbf{V}$  are the linear projections of  $\mathbf{Z}$ .

Advantages for Periodic Modeling.  $\bullet$  The key distance penalty threshold  $\beta$  is defined as a learnable parameter, which can be adaptively learned from the data. This enables more accurate and dynamic modeling of periodic trend.  $\bullet$  The modulation term varies smoothly with the periodic relative distance, implicitly encoding discriminative positional information of time steps. This effectively addresses the permutation-invariant limitation inherent in the standard attention mechanism.  $\bullet$  The function exhibits a smoother trend and can be flexibly controlled through the learnable parameter  $\alpha$ . This adaptability allows our Transformer architecture to customize personalized modulation behaviors for time series with diverse periodic characteristics, thereby refining the attention process to better capture periodic dependencies.

# 3.3 Forecasting and Optimization

The Final time series prediction values is obtained by summing the flattened and linearly transformed outputs of the Transformer backbone  $\vec{\mathbf{Z}} \in \mathbb{R}^{P \times d}$  as follows:

$$\hat{\mathbf{Y}} = \text{Flatten}(\mathbf{Z})\mathbf{W}_{out} + \mathbf{b}_{out} \in \mathbb{R}^L,$$
 (14)

where  $\mathbf{W}_{out} \in \mathbb{R}^{(P*d) \times L}$  and  $\mathbf{b}_{out} \in \mathbb{R}^L$  are learnable parameters. When the time series involves multiple variables (i.e., C > 1), we adopt channel-independent learning and compute the relative loss weights based on the maximum loss across channels, dynamically adjusting the loss weights during training to promote equal learning with stable convergence. This process is formulated as follows:

$$\mathcal{L}^* = \omega \sum_{c=1}^{C} \frac{\mathcal{L}(\mathbf{Y}_c, \hat{\mathbf{Y}}_c)}{\|\mathcal{L}(\mathbf{Y}_c, \hat{\mathbf{Y}}_c)\|}, \quad \omega = \max \left\{ \|\mathcal{L}(\mathbf{Y}_c, \hat{\mathbf{Y}}_c)\|; c \in \{1, 2, \dots, C\} \right\},$$
(15)

where  $\mathbf{Y}_c$ ,  $\hat{\mathbf{Y}}_c$  are the ground-truth values and prediction values of the channel c.

**Complexity Analysis.** MoFo's computational complexity has a quadratic dependence on period length but is independent of the input sequence length, making its efficiency highly favorable.

For example, in daily-period datasets with hourly sampling (e.g., 24 time steps per period), the computational cost is minimal relative to the total sequence length in long-term forecasting tasks, such as 720 time steps.

# 4 Experiment

### 4.1 Experimental Setup

**Datasets.** We conduct our experiments on widely used real-world time series datasets with periodic pattern from 4 different domains, including ETTh1, ETTh2, ETTm1, ETTm2, Weather, Solar Energy, Electricity, and Traffic. A summary of all datasets is provided in Table 1.

Tabl	e 1:	Statistics	of	used	datasets.
------	------	------------	----	------	-----------

Dataset	ETTh1	ETTh2	ETTm1	ETTm2	Weather	Solar Energy	Electricity	Traffic
# Channels	7	7	7	7	21	137	321	862
# Samples	14,400	14,400	57,600	57,600	52,696	52,560	26,304	17,544
Frequency	1 hour	1 hour	15 mins	15 mins	10 mins	10 mins	1 hour	1 hour
Split ratio	6:2:2	6:2:2	6:2:2	6:2:2	7:1:2	6:2:2	7:1:2	7:1:2

**Settings.** Our experiments are conducted on an NVIDIA A100 GPU with 40GB memory, using PyTorch under Python 3.11.5. We implement our method within the TFB platform [34] to ensure a fair comparison. Following the evaluation protocol in TFB [34], we report the best performance achieved over look-back window lengths  $T \in \{96, 336, 512\}$  and forecasting horizons  $L \in \{96, 192, 336, 720\}$ . Model performance is evaluated using two standard metrics: mean squared error (MSE) and mean absolute error (MAE). To maintain fairness in evaluation, we disable the "Drop Last" batch-sampling trick[18]. We use the Adam optimizer [14] with the  $L_1$  loss function from the FreDF strategy [53].

Table 2: Performance comparisons for LTSF. The **best** and <u>second best</u> are marked in corresponding colors. All experimental results are selected from the best performance under the look-back window length  $T \in \{96, 336, 512\}$ .

_	ethod	Me	oFo urs)		ET (25)	Pl	DF (24)		former (24)		ormer (24)	Cycl (20			Mixer (24)		hTST (23)		former (23)		near (23)
M	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1		$\frac{0.397}{0.407}$	0.413 0.424	0.398 <u>0.414</u>	<b>0.409</b> <u>0.426</u>	0.392 0.418	$0.414 \\ 0.435$	0.424 0.449	0.405 0.440 0.460 0.487	$0.408 \\ 0.438$	$0.415 \\ 0.434$	0.404 0.430	$0.417 \\ 0.429$	0.413 0.438	$0.430 \\ 0.450$	0.409 0.431	$0.425 \\ 0.444$	0.409 0.433	0.438 0.457	0.408 0.440	$0.419 \\ 0.440$
ETTh2		<b>0.327</b> 0.361	<b>0.373</b> <u>0.405</u>	0.332 0.353	0.374 0.397	0.339 0.374	$0.382 \\ 0.406$	0.372 0.388	0.348 0.403 0.417 0.444	0.345 0.378	0.380 0.408	0.341 0.370	$0.385 \\ 0.411$	0.349 0.366	$0.387 \\ 0.413$	0.348 0.377	$0.384 \\ 0.416$	0.723 0.740	0.607 0.628	0.387 0.490	0.423 0.487
ETTm1		0.320 0.347	0.363 0.382	<b>0.320</b> <u>0.348</u>	0.358 0.377	0.321 0.354	$0.364 \\ 0.383$	0.341 0.374	0.353 0.380 0.396 0.430	$0.337 \\ 0.374$	0.363 0.384	$0.332 \\ 0.366$	$0.365 \\ 0.386$	0.335 0.365	$0.372 \\ 0.386$	0.329 0.362	$0.368 \\ 0.390$	0.374 0.413	$0.410 \\ 0.432$	0.336 0.367	$0.366 \\ 0.386$
ETTTm2		0.211 0.258	0.283 0.314	$\frac{0.214}{0.267}$	$\begin{array}{c} \underline{0.287} \\ 0.320 \end{array}$	0.219 0.269	$0.290 \\ 0.330$	0.242 0.282	0.266 0.312 0.337 0.394	0.219 <u>0.267</u>	0.288 0.319	0.214 0.269	0.286 0.322	0.225 0.277	$0.298 \\ 0.332$	0.221 0.276	$0.293 \\ 0.327$	0.369 0.588	$0.416 \\ 0.600$	0.224 0.277	$0.304 \\ 0.337$
Weather		0.186 0.233	<b>0.230</b> <u>0.272</u>	$\frac{0.188}{0.234}$	0.231 0.268	0.193 0.245	$0.240 \\ 0.280$	0.200 0.252	0.207 0.248 0.287 0.336	0.191 0.243	$0.235 \\ 0.274$	0.213 0.262	$0.259 \\ 0.291$	0.192 0.247	$0.243 \\ 0.284$	0.191 0.242	0.239 0.279	0.195 0.254	0.261 0.319	0.216 0.258	0.273 0.307
Solar	336	0.177 0.186	$0.231 \\ 0.238$	0.187 0.199	0.207 0.213	0.199 0.208	$0.257 \\ 0.269$	0.193 0.203	0.244 0.257 0.266 0.281	0.196 0.195	$\frac{0.220}{0.220}$	0.221 0.233	$0.261 \\ 0.269$	0.201 0.190	$0.259 \\ 0.256$	0.204 0.212	$0.302 \\ 0.293$	0.208 0.212	0.226 0.239	0.220 0.234	$0.282 \\ 0.295$
Electricity	192 336	0.140 0.157	$0.234 \\ 0.252$	0.145 0.163	0.235 0.255	0.147 0.165	$0.242 \\ 0.260$	0.154 0.169	0.230 0.250 0.265 0.288	$0.157 \\ 0.170$	$0.253 \\ 0.267$	0.144 0.161	0.239 $0.253$	0.168 0.189	$0.269 \\ 0.291$	0.158 0.168	$0.260 \\ 0.267$	0.146 0.165	$0.243 \\ 0.264$	0.154 0.169	$0.251 \\ 0.268$
Traffic	192 336	0.379 0.390	0.254 0.258	0.383 0.395	<b>0.249</b> <u>0.259</u>	$\frac{0.382}{0.393}$	$0.261 \\ 0.268$	0.384 0.396	0.265 0.273 0.277 0.308	$0.405 \\ 0.424$	$0.257 \\ 0.265$	0.406 0.425	$0.280 \\ 0.291$	0.400 0.407	$0.272 \\ 0.272$	0.386 0.396	$0.269 \\ 0.275$	0.503 0.505	0.263 0.276	0.407 0.417	$0.280 \\ 0.286$

**Baselines.** We compare MoFo with 17 advanced baselines in long-term time series forecasting comprising DUET [37], PDF [7], iTransformer [24], Pathformer [5], CycleNet [20], TimeMixer [59], PatchTST [32], Crossformer [75], DLinear [73], NLinear [73], FITS [68], FiLM [78], MICN [57], FEDformer [79], Triformer [6], Non-stationary Transformer [23], and Informer [77].

# 4.2 Forecasting Performance of MoFo

Experimental results are summarized in Table 2. Due to space constraints, we compare against a larger set of baselines in Appendix D.3. CycleNet explicitly models periodic patterns, yet achieves lower forecasting performance compared to DLinear, which merely employs simple fully connected layers. TimeMixer, on the other hand, utilizes multi-scale modeling techniques to comprehensively capture complex temporal dynamics. Compared to methods such as PatchTST, Pathformer, and other approaches that rely on continuous patching strategies, our model achieves superior forecasting performance. This improvement is attributed to the proposed Periodicity-based Discrete Patching, which enables the model to better capture temporal dependencies across time steps. PDF introduces a multi-scale decomposition framework that models temporal dependencies from both long-term and short-term perspectives. DUET employs bidirectional clustering over both temporal and channel dimensions to adaptively capture spatio-temporal dependencies, achieving the best performance among existing baselines. However, it does not explicitly model periodic patterns and suffers from high computational complexity. Benefiting from the Period-Aware Modulator, our model focuses explicitly on periodic pattern learning, leading to the overall best forecasting accuracy. These results demonstrate the effectiveness of our design choices in capturing long-range periodic dependencies.

### 4.3 Efficiency Analysis of MoFo

We compare the computational efficiency with Transformer-based baselines on the Traffic dataset. As shown in Table 3 and Table 11, MoFo achieves the best forecasting accuracy among all Transformer-based models, while demonstrating the lowest computational complexity and highest efficiency. Similarly, PatchTST, which also adopts a patching strategy, employs an independent channel learning approach that significantly increases model complexity. Its parameter count grows by more than  $10\times$ , FLOPs increase by  $25\times$ , and training speed slows down by  $17\times$ . Pathformer enhances prediction accuracy through dynamic path adaptation and a patching strategy, albeit at the cost of increased learning overhead. Compared to DUET, one of the top-performing baseline models, MoFo reduces the number of parameters by more than  $3\times$ , the computational cost by more than  $10\times$ , and significantly lowers both memory consumption and training time. This is because our model's complexity grows quadratically with the period length, which is significantly shorter than the input sequence length. Moreover, we show that only a single Transformer layer is sufficient for effective learning.

Table 3: Efficiency comparison of MoFo and SOTA baselines with L=720 in Traffic dataset. All results of each model are under the optimal hyperparameters for fair comparison. Parameters: All learnable parameters requiring gradient descent. MACs: multiply–accumulate operations. FLOPs: floating point operations. M: Million  $(10^6)$ . B: Billion  $(10^9)$ . T: Trillion  $(10^{12})$ . MB: Megabyte. s: Second.  $\uparrow$  indicates the relative percentage increasing regarding MoFo.

	Models	MSE	# Parameters	# MACs	# FLOPs	Memory Usage	<b>Epoch Time</b>
	Crossformer	0.552 <sub>↑30.18%</sub>	3.23 M <sub>↑34.58%</sub>	85.03 B <sub>↑13.43%</sub>	92.41 B <sub>↑20.94%</sub>	13,556 MB <sub>↑171.77%</sub>	368 s <sub>↑682.97%</sub>
720]	PatchTST	$0.435_{\uparrow 2.59\%}$	27.8 M <sub>↑1058.75%</sub>	2.02 T <sub>↑2594.77%</sub>	2.08 T <sub>↑2622.16%</sub>	44,782 MB <sub>↑797.79%</sub>	876 s <sub>1763.83</sub> %
ÎII	Pathformer	$0.452_{\uparrow 6.60\%}$	9.61 M <sub>↑300.42%</sub>	110.46 B <sub>↑47.36%</sub>	117.76 B <sub>↑54.12%</sub>	36,602 MB <sub>↑633.80%</sub>	1,081 s <sub>†2200.0%</sub>
T]	iTransformer	$0.445_{\uparrow 4.95\%}$	5.37 M <sub>123.75</sub> %	297.96 B <sub>297.49</sub> %	446.30 B <sub>↑484.09%</sub>	19,608 MB <sub>293.10</sub> %	52 s <sub>10.64</sub> %
ЩС	PDF	0.438	2.45 M <sub>↑2.08%</sub>	637.05 B <sub>↑749.85%</sub>	662.47 B <sub>↑766.99%</sub>	38,014 MB <sub>1662.11%</sub>	76 s <sub>1.70</sub> %
Traffic	DUET	0.435 <sub>12.59</sub> %	11.2 M <sub>↑367.08%</sub>	137.33 B <sub>↑83.20%</sub>	975.96 B <sub>↑1177.27%</sub>	75,616 MB <sub>↑1415.96%</sub>	516 s <sub>↑997.87%</sub>
•	MoFo	0.424	2.40 M	74.96 B	<b>76.74</b> B	<b>4,988</b> MB	<b>47</b> s

### 4.4 Hyperparameters Sensitivity Experiments

We investigate the sensitivity of MoFo to its two core hyperparameters—the number of Transformer layers and model dimensionality—on the ETTh2 and Electricity datasets. For each forecasting horizon  $L \in \{96, 192, 336, 720\}$ , we use the best-performing hyperparameter configuration and vary only the target hyperparameter, as shown in Figure 4.4. We report both MSE and MAE for evaluation. The number of Transformer layers is varied from 1 to 6, and the model dimensionality is tested in

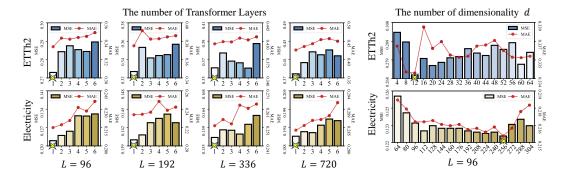


Figure 4: Ablation study on all datasets with forecasting horizons L = 712.

the range of  $4\sim64$  (ETTh2) and  $64\sim304$  (Electricity), with the number of attention heads fixed at H=4. Results show that MoFo achieves strong performance even with a single Transformer layer, while additional layers yield marginal gains. This suggests that the model efficiently captures interpatch dependencies without requiring deep architectures. Furthermore, increasing the dimensionality generally improves performance; however, the gains plateau beyond a certain threshold, indicating diminishing returns from further capacity expansion.

## 4.5 Ablation Study

We design ablation experiments to validate the soundness of the each component of MoFo: '+ Cpatch' which used continuous patch technology; '+Sinuous Pos' and '+ Learnable Pos' which use Sinuous and Learnable position instead of our periodic relative position, respectively; '- Modulator' which removes the period-aware modulator; '+ Mean Loss' which only use L1 loss function. As shown in Figure 3, the ablation study reveals that '-Modulator' variant achieves the worst performance. This is because the Period-Aware Modulator plays a crucial role in guiding the model to focus on extracting periodic patterns. '+cpatch' variant also suffers from higher prediction errors, which can be attributed to the fact that our proposed discrete patching strategy enables direct modeling of dependencies among periodically aligned time steps. In summary, all variants perform worse than MoFo model to varying degrees, demonstrating the effectiveness and necessity of each component in capturing long-range periodic dependencies.

### 4.6 Scalability for Look-back Window of MoFo

We evaluate the scalability of MoFo under varying look-back window lengths. Specifically, we construct an ultra-long look-back setting on the ETTm2 dataset, where the input sequence length reaches up to 10K, one hundred times of the forecasting horizon: L=96, T=96\*(5k),  $k=\{1,2,\ldots,20\}$ . We report the maximum memory consumption during training, per-epoch training time, and FLOPs. As shown in Figure 5, as the look-back window increases, both DUET and PatchTST exhibit a sharp rise in memory usage and training time, indicating that their computational complexity scales closely with the input length. In contrast, even when the input sequence length grows by  $100\times$ , MoFo sees less than a 10% increase in peak memory usage and less than a 2% increase in training time. Its computational cost depends primarily on the period length rather than

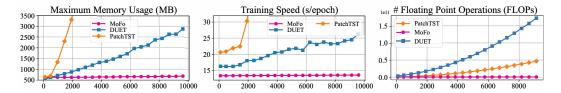


Figure 5: Scalability for look-back window of models on ETTm2 dataset.

the full sequence length. These results highlight the strong scalability of MoFo and demonstrate its great potential for modeling ultra-long sequences with minimal overhead.

### 4.7 Attention Visualization

We extract and visualize the attention coefficient matrices from models trained on four datasets with periodic lengths (P) of 24, 96, 144, and 24, respectively, all spanning a time window of one day. As shown in Figure 6, the coefficients in each attention head exhibit distinct periodic patterns. Notably, when the cyclical distance between two time steps is the largest (i.e., equal to P/2), the attention scores reach their minimum values. These attention coefficients that align with the underlying periodic trends significantly enhance the model's performance in long-term time series forecasting.

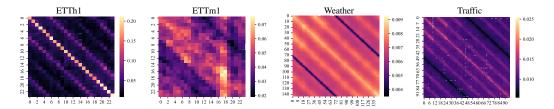


Figure 6: Visualization of attention scores.

## 5 Conclusion

In this paper, we propose MoFo, a novel framework for time series forecasting that leverages the inherent periodic structure of temporal data to explicitly model both periodic correlations and temporal trends. Through our proposed Period-structured patches, the model is able to directly capture correlations among time steps that share the same phase across periods. We further introduce a period-aware modulator, which enhances the attention mechanism with an adaptive inductive bias guided by underlying periodic trends. Experimental results demonstrate that MoFo achieves competitive forecasting performance compared to state-of-the-art methods, while significantly improving computational efficiency and scalability, especially for long input sequences.

# Acknowledgment

This paper is partially supported by the National Natural Science Foundation of China (No.12227901). The AI-driven experiments, simulations and model training were performed on the robotic AI-Scientist platform of Chinese Academy of Sciences.

# References

- [1] David R Anderson, Kenneth P Burnham, and William L Thompson. Null hypothesis testing: problems, prevalence, and an alternative. *The journal of wildlife management*, pages 912–923, 2000.
- [2] Hossein Arsham and Miodrag Lovric. Bartlett's test. *International encyclopedia of statistical science*, 2(2):20–23, 2011.

- [3] George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time series analysis: forecasting and control.* John Wiley & Sons, 2015.
- [4] David Chandler. Introduction to modern statistical. *Mechanics. Oxford University Press, Oxford, UK*, 5(449):11, 1987.
- [5] Peng Chen, Yingying Zhang, Yunyao Cheng, Yang Shu, Yihang Wang, Qingsong Wen, Bin Yang, and Chenjuan Guo. Pathformer: Multi-scale transformers with adaptive pathways for time series forecasting. *arXiv preprint arXiv:2402.05956*, 2024.
- [6] Razvan-Gabriel Cirstea, Chenjuan Guo, Bin Yang, Tung Kieu, Xuanyi Dong, and Shirui Pan. Triformer: Triangular, variable-specific attentions for long sequence multivariate time series forecasting–full version. *arXiv preprint arXiv:2204.13767*, 2022.
- [7] Tao Dai, Beiliang Wu, Peiyuan Liu, Naiqi Li, Jigang Bao, Yong Jiang, and Shu-Tao Xia. Periodicity decoupling framework for long-term series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024.
- [8] Rahul Dey and Fathi M Salem. Gate-variants of gated recurrent unit (gru) neural networks. In 2017 IEEE 60th international midwest symposium on circuits and systems (MWSCAS), pages 1597–1600. IEEE, 2017.
- [9] Petr Hájek and Michal Johanis. Smooth approximations. *Journal of Functional Analysis*, 259 (3):561–582, 2010.
- [10] Qihe Huang, Lei Shen, Ruixin Zhang, Shouhong Ding, Binwu Wang, Zhengyang Zhou, and Yang Wang. Crossgnn: Confronting noisy multivariate time series via cross interaction refinement. Advances in Neural Information Processing Systems, 36:46885–46902, 2023.
- [11] Qihe Huang, Lei Shen, Ruixin Zhang, Jiahuan Cheng, Shouhong Ding, Zhengyang Zhou, and Yang Wang. Hdmixer: Hierarchical dependency with extendable patch for multivariate time series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 12608–12616, 2024.
- [12] Qihe Huang, Zhengyang Zhou, , Yangze Li, Kuo Yang, Binwu Wang, and Yang Wang. Many minds, one goal: Time series forecasting via sub-task specialization and inter-agent cooperation. In Advances in Neural Information Processing Systems, 2025.
- [13] Qihe Huang, Zhengyang Zhou, Kuo Yang, Zhongchao Yi, Xu Wang, and Yang Wang. Timebase: The power of minimalism in efficient long-term time series forecasting. In *Forty-second International Conference on Machine Learning*, 2025.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- [15] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*, pages 3744–3753. PMLR, 2019.
- [16] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyou Zhou, Wenhu Chen, Yu-Xiang Wang, and Xifeng Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *Advances in neural information processing systems*, 32, 2019.
- [17] Xinyu Li, Yuchen Luo, Hao Wang, Haoxuan Li, Liuhua Peng, Feng Liu, Yandong Guo, Kun Zhang, and Mingming Gong. Towards accurate time series forecasting via implicit decoding. *Advances in Neural Information Processing Systems*, 2025.
- [18] Zhe Li, Xiangfei Qiu, Peng Chen, Yihang Wang, Hanyin Cheng, Yang Shu, Jilin Hu, Chenjuan Guo, Aoying Zhou, Qingsong Wen, et al. Foundts: Comprehensive and unified benchmarking of foundation models for time series forecasting. *arXiv* preprint arXiv:2410.11802, 2024.
- [19] Shengsheng Lin, Weiwei Lin, Wentai Wu, Feiyu Zhao, Ruichao Mo, and Haotong Zhang. Segrnn: Segment recurrent neural network for long-term time series forecasting. arXiv preprint arXiv:2308.11200, 2023.

- [20] Shengsheng Lin, Weiwei Lin, Xinyi Hu, Wentai Wu, Ruichao Mo, and Haocheng Zhong. Cyclenet: enhancing time series forecasting through modeling periodic patterns. *Advances in Neural Information Processing Systems*, 37:106315–106345, 2024.
- [21] Chenxi Liu, Hao Miao, Cheng Long, Yan Zhao, Ziyue Li, and Panos Kalnis. Llms meet cross-modal time series analytics: Overview and directions. In *Proceedings of the 19th International Symposium on Spatial and Temporal Data*, pages 101–106, 2025.
- [22] Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In # PLACEHOLDER\_PARENT\_METADATA\_VALUE#, 2022.
- [23] Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. *Advances in neural information processing systems*, 35:9881–9893, 2022.
- [24] Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. arXiv preprint arXiv:2310.06625, 2023.
- [25] Zhipeng Liu, Peibo Duan, Binwu Wang, Xuan Tang, Qi Chu, Changsheng Zhang, Yongsheng Huang, and Bin Zhang. Disms-ts: Eliminating redundant multi-scale features for time series classification. *arXiv* preprint arXiv:2507.04600, 2025.
- [26] Zhipeng Liu, Peibo Duan, Xiaosha Xue, Changsheng Zhang, Wenwei Yue, and Bin Zhang. A dynamic hypergraph attention network: Capturing market-wide spatiotemporal dependencies for stock selection. *Applied Soft Computing*, 169:112524, 2025.
- [27] Donghao Luo and Xue Wang. Moderntcn: A modern pure convolution structure for general time series analysis. In *The twelfth international conference on learning representations*, pages 1–43, 2024.
- [28] Jiaming Ma, Zhiqing Cui, Binwu Wang, Pengkun Wang, Zhengyang Zhou, Zhe Zhao, and Yang Wang. Causal learning meet covariates: Empowering lightweight and effective nationwide air quality forecasting. *International Joint Conference on Artificial Intelligence*, 2025.
- [29] Jiaming Ma, Binwu Wang, Pengkun Wang, Zhengyang Zhou, Xu Wang, and Yang Wang. Bist: A lightweight and efficient bi-directional model for spatiotemporal prediction. *Proceedings of the VLDB Endowment*, 18(6):1663–1676, 2025.
- [30] Jiaming Ma, Binwu Wang, Pengkun Wang, Zhengyang Zhou, Xu Wang, and Yang Wang. Robust spatio-temporal centralized interaction for ood learning. In *Forty-second International Conference on Machine Learning*, 2025.
- [31] Jiaming Ma, Binwu Wang, Pengkun Wang, Zhengyang Zhou, Yudong Zhang, Xu Wang, and Yang Wang. Mobimixer: A multi-scale spatiotemporal mixing model for mobile traffic prediction. *IEEE Transactions on Mobile Computing*, 2025.
- [32] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. arXiv preprint arXiv:2211.14730, 2022.
- [33] Zhen Qin, Xiaodong Han, Weixuan Sun, Dongxu Li, Lingpeng Kong, Nick Barnes, and Yiran Zhong. The devil in linear transformer. *arXiv preprint arXiv:2210.10340*, 2022.
- [34] Xiangfei Qiu, Jilin Hu, Lekui Zhou, Xingjian Wu, Junyang Du, Buang Zhang, Chenjuan Guo, Aoying Zhou, Christian S Jensen, Zhenli Sheng, et al. Tfb: Towards comprehensive and fair benchmarking of time series forecasting methods. *arXiv preprint arXiv:2403.20150*, 2024.
- [35] Xiangfei Qiu, Zhe Li, Wanghui Qiu, Shiyan Hu, Lekui Zhou, Xingjian Wu, Zhengyu Li, Chenjuan Guo, Aoying Zhou, Zhenli Sheng, Jilin Hu, Christian S. Jensen, and Bin Yang. Tab: Unified benchmarking of time series anomaly detection methods. In *Proc. VLDB Endow.*, pages 2775–2789, 2025.

- [36] Xiangfei Qiu, Xingjian Wu, Hanyin Cheng, Xvyuan Liu, Chenjuan Guo, Jilin Hu, and Bin Yang. DBLoss: Decomposition-based loss function for time series forecasting. In *NeurIPS*, 2025.
- [37] Xiangfei Qiu, Xingjian Wu, Yan Lin, Chenjuan Guo, Jilin Hu, and Bin Yang. DUET: Dual clustering enhanced multivariate time series forecasting. In *SIGKDD*, pages 1185–1196, 2025.
- [38] Adrian E Raftery, WR Gilks, S Richardson, and D Spiegelhalter. Hypothesis testing and model. *Markov chain Monte Carlo in practice*, 1:165–87, 1995.
- [39] Fred L Ramsey. Characterization of the partial autocorrelation function. *The Annals of Statistics*, pages 1296–1301, 1974.
- [40] Noam Shazeer. Glu variants improve transformer. arXiv preprint arXiv:2002.05202, 2020.
- [41] Michael Smithson. Confidence intervals. Number 140. Sage, 2003.
- [42] Peiwang Tang and Weitai Zhang. Unlocking the power of patch: Patch-based mlp for long-term time series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 12640–12648, 2025.
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [44] Binwu Wang, Yudong Zhang, Jiahao Shi, Pengkun Wang, Xu Wang, Lei Bai, and Yang Wang. Knowledge expansion and consolidation for continual traffic prediction with expanding graphs. *IEEE Transactions on Intelligent Transportation Systems*, 2023.
- [45] Binwu Wang, Yudong Zhang, Xu Wang, Pengkun Wang, Zhengyang Zhou, Lei Bai, and Yang Wang. Pattern expansion and consolidation on evolving graphs for continual traffic prediction. In *Proc. of KDD*, 2023.
- [46] Binwu Wang, Jiaming Ma, Pengkun Wang, Xu Wang, Yudong Zhang, Zhengyang Zhou, and Yang Wang. Stone: A spatio-temporal ood learning framework kills both spatial and temporal shifts. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 2948–2959, 2024.
- [47] Binwu Wang, Pengkun Wang, Yudong Zhang, Xu Wang, Zhengyang Zhou, Lei Bai, and Yang Wang. Towards dynamic spatial-temporal graph learning: A decoupled perspective. In *Proc. of AAAI*, 2024.
- [48] Chengsen Wang, Qi Qi, Jingyu Wang, Haifeng Sun, Zirui Zhuang, and Jianxin Liao. Unlocking the potential of linear networks for irregular multivariate time series forecasting. *arXiv* preprint *arXiv*:2505.00590, 2025.
- [49] Chengsen Wang, Qi Qi, Jingyu Wang, Haifeng Sun, Zirui Zhuang, Jinming Wu, Lei Zhang, and Jianxin Liao. Chattime: A unified multimodal time series foundation model bridging numerical and textual data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 12694–12702, 2025.
- [50] Guoxin Wang, Yijuan Lu, Lei Cui, Tengchao Lv, Dinei Florencio, and Cha Zhang. A simple yet effective learnable positional encoding method for improving document transformer model. In *Findings of the Association for Computational Linguistics: AACL-IJCNLP 2022*, pages 453–463, 2022.
- [51] Hao Wang, Zhiyu Wang, Yunlong Niu, Zhaoran Liu, Haozhe Li, Yilin Liao, Yuxin Huang, and Xinggao Liu. An accurate and interpretable framework for trustworthy process monitoring. *IEEE Transactions on Artificial Intelligence*, 5(5):2241–2252, 2023.
- [52] Hao Wang, Haoxuan Li, Xu Chen, Mingming Gong, Zhichao Chen, et al. Optimal transport for time series imputation. In *The Thirteenth International Conference on Learning Representations*, 2025.

- [53] Hao Wang, Lichen Pan, Yuan Shen, Zhichao Chen, Degui Yang, Yifei Yang, Sen Zhang, Xinggao Liu, Haoxuan Li, and Dacheng Tao. Fredf: Learning to forecast in the frequency domain. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [54] Hao Wang, Licheng Pan, Zhichao Chen, Xu Chen, Qingyang Dai, Lei Wang, Haoxuan Li, and Zhouchen Lin. Time-o1: Time-series forecasting needs transformed label alignment. Advances in Neural Information Processing Systems, 2025.
- [55] Haotian Wang, Haoxuan Li, Hao Zou, Haoang Chi, Long Lan, Wanrong Huang, and Wenjing Yang. Effective and efficient time-varying counterfactual prediction with state-space models. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [56] Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei. Deepnet: Scaling transformers to 1,000 layers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [57] Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. Micn: Multi-scale local and global context modeling for long-term series forecasting. In *The eleventh international conference on learning representations*, 2023.
- [58] Lei Wang, Shanshan Huang, Chunyuan Zheng, Jun Liao, Xiaofei Zhu, Haoxuan Li, and Li Liu. Mitigating data imbalance in time series classification based on counterfactual minority samples augmentation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery* and Data Mining V. 2, pages 2962–2973, 2025.
- [59] Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y Zhang, and Jun Zhou. Timemixer: Decomposable multiscale mixing for time series forecasting. arXiv preprint arXiv:2405.14616, 2024.
- [60] Yihe Wang, Nan Huang, Taida Li, Yujun Yan, and Xiang Zhang. Medformer: A multi-granularity patching transformer for medical time-series classification. arXiv preprint arXiv:2405.19363, 2024.
- [61] Eric W Weisstein. Normal distribution. https://mathworld. wolfram. com/, 2002.
- [62] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in neural information processing systems*, 34:22419–22430, 2021.
- [63] Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. arXiv preprint arXiv:2210.02186, 2022.
- [64] Xingjian Wu, Xiangfei Qiu, Hanyin Cheng, Zhengyu Li, Jilin Hu, Chenjuan Guo, and Bin Yang. Enhancing time series forecasting through selective representation spaces: A patch perspective. In *NeurIPS*, 2025.
- [65] Xingjian Wu, Xiangfei Qiu, Hongfan Gao, Jilin Hu, Bin Yang, and Chenjuan Guo. K<sup>2</sup>VAE: A koopman-kalman enhanced variational autoencoder for probabilistic time series forecasting. In ICML, 2025.
- [66] Xingjian Wu, Xiangfei Qiu, Zhengyu Li, Yihang Wang, Jilin Hu, Chenjuan Guo, Hui Xiong, and Bin Yang. CATCH: Channel-aware multivariate time series anomaly detection via frequency patching. In *ICLR*, 2025.
- [67] Mingyuan Xia, Chunxu Zhang, Zijian Zhang, Hao Miao, Qidong Liu, Yuanshao Zhu, and Bo Yang. Timeemb: A lightweight static-dynamic disentanglement framework for time series forecasting. arXiv preprint arXiv:2510.00461, 2025.
- [68] Zhijian Xu, Ailing Zeng, and Qiang Xu. Fits: Modeling time series with 10k parameters. arXiv preprint arXiv:2307.03756, 2023.
- [69] Kangjia Yan, Chenxi Liu, Hao Miao, Xinle Wu, Yan Zhao, Chenjuan Guo, and Bin Yang. Deciphering invariant feature decoupling in source-free time series forecasting with proxy denoising. *arXiv* preprint arXiv:2510.05589, 2025.

- [70] Wenzhen Yue, Xianghua Ying, Ruohao Guo, DongDong Chen, Ji Shi, Bowei Xing, Yuqing Zhu, and Taiyan Chen. Sub-adjacent transformer: Improving time series anomaly detection with reconstruction error from sub-adjacent neighborhoods. *arXiv* preprint arXiv:2404.18948, 2024.
- [71] Wenzhen Yue, Yong Liu, Hao Wang, Haoxuan Li, Xianghua Ying, Ruohao Guo, Bowei Xing, and Ji Shi. Olinear: A linear model for time series forecasting in orthogonally transformed domain. Advances in Neural Information Processing Systems, 2025.
- [72] Wenzhen Yue, Yong Liu, Xianghua Ying, Bowei Xing, Ruohao Guo, and Ji Shi. Freeformer: Frequency enhanced transformer for multivariate time series forecasting. *arXiv preprint arXiv:2501.13989*, 2025.
- [73] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 11121–11128, 2023.
- [74] Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [75] Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *The eleventh international conference on learning representations*, 2023.
- [76] Zheng Zhao, Weihai Chen, Xingming Wu, Peter CY Chen, and Jingmeng Liu. Lstm network: a deep learning approach for short-term traffic forecast. *IET Intelligent Transport Systems*, 2017.
- [77] 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 Proceedings of the AAAI conference on artificial intelligence, volume 35, pages 11106–11115, 2021.
- [78] Tian Zhou, Ziqing Ma, Qingsong Wen, Liang Sun, Tao Yao, Wotao Yin, Rong Jin, et al. Film: Frequency improved legendre memory model for long-term time series forecasting. *Advances in neural information processing systems*, 35:12677–12690, 2022.
- [79] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International conference on machine learning*, pages 27268–27286. PMLR, 2022.

# **A** Mathematics Justification

### A.1 The Proof of Theorem 1

**Proof of (1):** For arbitrary  $\alpha, \beta > 0$ , there are  $\mathcal{S}(0; \alpha, \beta) = 1, \mathcal{S}(+\infty; \alpha, \beta) = 0$ . In fact

$$S(0; \alpha, \beta) = \frac{1}{1 + \exp(\alpha(0 - \beta))} + \frac{\exp(-0)}{1 + \exp(\alpha\beta)} = \frac{1}{1 + \exp(-\alpha\beta)} + \frac{1}{1 + \exp(\alpha\beta)}$$

$$= \frac{\exp(\alpha\beta)}{1 + \exp(\alpha\beta)} + \frac{1}{1 + \exp(\alpha\beta)} = 1.$$
(16)

and

$$S(+\infty; \alpha, \beta) = \frac{1}{1 + \exp(\alpha(+\infty - \beta))} + \frac{\exp(-\infty)}{1 + \exp(\alpha\beta)} = \frac{1}{1 + \infty} + \frac{0}{1 + \exp(\alpha\beta)} = 0. \quad (17)$$

**Proof of (2):** The upper error bound of the smooth approximation satisfies,

$$\int_{0}^{+\infty} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma \le \frac{\log 2}{\alpha} - \frac{1}{1 + \exp(\alpha)} \to 0 \quad (\alpha \to +\infty).$$
 (18)

In fact, the Heaviside step function  $\mathcal{H}(\beta - \gamma)$  can be viewed as the differentiation of the maximum value function  $\max\{0, \beta - \gamma\}$  as follows,

$$\mathcal{H}(\beta - \gamma) = \frac{\mathrm{d}}{\mathrm{d}(\beta - \gamma)} \max\{0, \beta - \gamma\} = -\frac{\mathrm{d}}{\mathrm{d}\gamma} \max\{0, \beta - \gamma\}. \tag{19}$$

Our sigmoidal attenuation function can be seen as follows,

$$S(\gamma; \alpha, \beta) = -\frac{\mathrm{d}}{\mathrm{d}\gamma} \frac{1}{\alpha} \log \left(1 + \exp\left(\alpha(\beta - \gamma)\right)\right) - \frac{\mathrm{d}}{\mathrm{d}\gamma} \frac{\exp\left(-\gamma\right)}{1 + \exp\left(\alpha\beta\right)}.$$
 (20)

Hence, the cumulative error  $\int_0^{+\infty} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| d\gamma$  satisfies,

$$\int_{0}^{+\infty} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma$$

$$= \int_{0}^{\beta} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma + \int_{\beta}^{+\infty} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma$$

$$= \int_{0}^{\beta} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma + \int_{\beta}^{+\infty} |\mathcal{H}(\beta - \gamma) + \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma$$

$$= \int_{0}^{\beta} -\frac{d}{d\gamma} \max\{0, \beta - \gamma\} + \frac{d}{d\gamma} \frac{1}{\alpha} \log(1 + \exp(\alpha(\beta - \gamma))) + \frac{d}{d\gamma} \frac{\exp(-\gamma)}{1 + \exp(\alpha\beta)} \, d\gamma$$

$$+ \int_{\beta}^{+\infty} \frac{d}{d\gamma} \max\{0, \beta - \gamma\} - \frac{d}{d\gamma} \frac{1}{\alpha} \log(1 + \exp(\alpha(\beta - \gamma))) - \frac{d}{d\gamma} \frac{\exp(-\gamma)}{1 + \exp(\alpha\beta)} \, d\gamma$$

$$= \left( -\max\{0, \beta - \gamma\} + \frac{1}{\alpha} \log(1 + \exp(\alpha(\beta - \gamma))) + \frac{\exp(-\gamma)}{1 + \exp(\alpha\beta)} \right) \Big|_{\beta}^{\beta}$$

$$= \left( -0 + \frac{\log 2}{\alpha} + \frac{\exp(-\beta)}{1 + \exp(\alpha\beta)} + \beta - \frac{1}{\alpha} \log(1 + \exp(\alpha\beta)) - \frac{1}{1 + \exp(\alpha\beta)} \right)$$

$$+ \left( 0 - 0 - 0 - 0 + \frac{\log 2}{\alpha} + \frac{\exp(-\beta)}{1 + \exp(\alpha\beta)} \right)$$

$$= \frac{2\log 2}{\alpha} + \frac{2\exp(-\beta) - 1}{1 + \exp(\alpha\beta)} + \beta - \frac{1}{\alpha} \log(1 + \exp(\alpha\beta)).$$

Since,

$$\beta - \frac{1}{\alpha} \log \left( 1 + \exp \left( \alpha \beta \right) \right) < \beta - \frac{1}{\alpha} \log \left( \exp \left( \alpha \beta \right) \right) = \beta - \frac{\alpha \beta}{\alpha} = \beta - \beta = 0, \tag{22}$$

and

$$\frac{2\log 2}{\alpha} + \frac{2\exp\left(-\beta\right) - 1}{1 + \exp\left(\alpha\beta\right)} \le \frac{2\log 2}{\alpha} + \frac{1}{1 + \exp\left(\alpha\beta\right)} \le \frac{2\log 2}{\alpha} + \frac{1}{1 + \exp\left(\alpha\right)},\tag{23}$$

then the the cumulative error has upper bound satisfies

$$\int_{0}^{+\infty} |\mathcal{H}(\beta - \gamma) - \mathcal{S}(\gamma; \alpha, \beta)| \, d\gamma < \frac{2\log 2}{\alpha} + \frac{1}{1 + \exp(\alpha)} \to 0^{+} \quad (\alpha \to +\infty).$$
 (24)

# A.2 Autocorrelation Function and Bartlett's Test with Null Hypothesis

Autocorrelation Function (ACF). The Autocorrelation Function [3] measures the linear correlation between a time series and a lagged version of itself. For the complete observable time series values  $\mathbf{X} \in \mathbb{R}^{N_t}$  with total observable time steps  $N_t$ , the autocorrelation function  $\mathcal{A}_k$  of  $\mathbf{X}$  at lag  $k \geq 0$  is the ratio of the estimator of the covariance between the time series and the series lagged by k to the estimator of the variance of the time series as follows:

$$\mathcal{A}_{k} = \frac{\sum_{t=1}^{N_{t}-k} (\mathbf{X}_{t} - \mu)(\mathbf{X}_{t+k} - \mu)}{\sum_{t=1}^{N_{t}} (\mathbf{X}_{t} - \mu)^{2}} \in [-1, 1],$$
(25)

where  $\mu = \frac{1}{N_t} \sum_{t=1}^{N_t} \mathbf{X}_t \in \mathbb{R}$  is the expected value of  $\mathbf{X}$ .

Bartlett's Test with Null Hypothesis. The autocorrelation coefficients  $\mathcal{A}_k$  can be viewed as random variables. However, even in the time series consisting of pure random noise, there may exist non-zero autocorrelation coefficients at some lags [39]. Hence, we require a method to ascertain whether an observed autocorrelation  $\mathcal{A}_k$  represents a truly non-zero population or is merely due to this inherent randomness. This is typically achieved by performing confidence intervals derived from Bartlett's test [2] on a null hypothesis [38]. The null hypothesis denoted as  $\mathcal{H}_0$  is a fundamental concept in statistical inference [1]. Its primary purpose is to serve as a base assumption for hypothesis testing. In the context of ACF, the relevant null hypothesis is  $\mathcal{H}_0: \mathcal{A}_k = 0$  that the ground-true autocorrelation coefficient  $\mathcal{A}_k$  at a specific lag k is zero. The Bartlett's test constructs confidence intervals to test this null hypothesis for individual lags. These intervals are based on an estimate of the standard deviation of  $\mathcal{A}_k$  under the assumption that  $\mathcal{H}_0$  is true. Specifically, the  $1-\alpha$  confidence interval<sup>2</sup> for the autocorrelation  $\mathcal{A}_k$  under the null hypothesis  $\mathcal{H}_0: \mathcal{A}_k = 0$  is centered at 0, with boundaries given by  $\pm Z_{\alpha/2} \cdot \mathrm{SE}(\mathcal{A}_k)$ . Using Bartlett's formula for the variance of  $\mathcal{A}_k$ , the standard error  $\mathrm{SE}(\mathcal{A}_k)$  is approximated by:

$$SE(\mathcal{A}_k) \approx \sqrt{\frac{1}{N_t} (1 + 2\sum_{j=1}^{k-1} \mathcal{A}_j^2)},$$
(26)

where  $\mathcal{A}_j$  are the autocorrelations for lags  $j=1,2,\ldots,k-1$  and  $Z_{\alpha/2}$  is the  $1-\alpha/2$  quantile of the standard normal distribution (e.g.,  $Z_{0.025}\approx 1.96$  for a 95% confidence interval, corresponding to  $\alpha=0.05$ ) [61]. Thus, the approximate  $1-\alpha$  confidence interval boundaries for testing  $H_0:\mathcal{A}_k=0$  are:

$$\left[ -Z_{\alpha/2} \sqrt{\frac{1}{N_t} (1 + 2\sum_{j=1}^{k-1} \mathcal{A}_j^2), +Z_{\alpha/2} \sqrt{\frac{1}{N_t} (1 + 2\sum_{j=1}^{k-1} \mathcal{A}_j^2)} \right].$$
 (27)

Any autocorrelation value  $A_k$  that falls within this confidence interval is considered consistent with the null hypothesis ( $\mathcal{H}_0: A_k = 0$ ). In such cases, we do not have sufficient statistical evidence to

<sup>&</sup>lt;sup>2</sup>It is important to note that the alpha  $\alpha$  here is different from the learnable parameter alpha mentioned in the RRF in the MoFo. The alpha  $\alpha$  here is a specific parameter notation used in statistics [41]. We maintain consistency here to reduce potential confusion with specialized terminology.

conclude that the true autocorrelation at lag k is different from zero; the observed  $A_k$  is likely due to random variation inherent in a noise process. Conversely, if  $A_k$  exceeds these boundaries, we reject the null hypothesis and conclude that the autocorrelation at lag k is statistically significant.

We present in the Fig 7 the average ACF values and from some selective channel for all datasets used in this study, calculated across each channel with a maximum lag of 512, and report the null hypothesis region with 95% confidence interval. We observe that at the same positions within each period of the time series, the ACF values exhibit significant peaks. Therefore, it is highly rational to apply patching based on periodic positions for the time steps. The orange region represents the 95% confidence interval for the null hypothesis from Bartlett's Test [2]. Correlation coefficients that fall within this interval are not statistically significant and cannot be rejected as noise [4]

### **B** Related Work

Time series modeling serves as a core task in numerous domains, such as transportation and atmospheric science [21, 29, 44, 45, 46, 47, 48, 49, 69]. Early approaches primarily relied on recurrent neural networks (RNNs) and temporal convolutional networks (TCNs). In recent years, models based on multilayer perceptrons (MLPs) have garnered significant attention due to their lightweight and efficient performance.

**RNN-based Models.** RNNs, among the earliest deep learning architectures for sequential data, have been widely adopted for long-term time series forecasting, with notable variants such as LSTM [76] and GRU [8]. To mitigate the problem of too many recurrent steps, SegRNN [19] introduces a segmented recurrence mechanism combined with parallel multi-step prediction, substantially cutting down the number of iterations.

TCN-Based Model. TCNs employ convolutional operations to effectively model local contextual patterns in time series, offering a good trade-off between computational efficiency and forecasting accuracy. Recent advances have extended TCNs to better capture long-range temporal dependencies. For instance, ModernTCN [27] adopts large convolutional kernels to greatly expand the receptive field, allowing the model to capture broader temporal structures. Likewise, Pyraformer [22] integrates TCN layers with a Transformer framework; it uses stacked TCN layers for downsampling to obtain coarse-grained time series representations, which are then processed by the Transformer to enhance both scalability and performance.

MLP-Based Model. MLP-based models, when thoughtfully designed, have shown strong performance in time series forecasting. DLinear [73] illustrates this by using a moving average kernel to decompose the input series into trend and seasonal parts, each modeled separately by dedicated linear layers. PatchMLP [42] adopts a patching strategy that incorporates channel mixing to improve cross-variable information exchange. Extending this idea, HDMixer [11] employs adaptive patch lengths to capture both intra-patch short-term dynamics and inter-patch long-term dependencies while modeling intricate variable interactions. Meanwhile, FITS [68] operates MLPs in the frequency domain, leveraging spectral analysis to emphasize dominant signal components and better capture global temporal relationships.

### C The Transformer Layer in MoFo

We use the pre-norm Transformer Layer [43, 56] of multi-head attention with Regulated Relaxation Function (Eq. 12).

Vanilla Transformer Layer. Transformer [43] consists of the self-attention function  $\operatorname{MultiHead}(\cdot)$  with feedforward networks  $\operatorname{FFN}(\cdot)$ , and two distinct normalization layers  $\operatorname{Norm}_i(\cdot)$ . Assuming the input hidden representation is  $\mathbf{Z} \in \mathbb{R}^{P \times d}$  with period length P and model dimensionality d, the output hidden representation  $\mathbf{Z} \in \mathbb{R}^{n \times d}$  of one Transformer layer is as follows,

$$\vec{\mathbf{Z}} = \text{Multi-Head}(\text{Norm}_2(\bar{\mathbf{Z}})) + \bar{\mathbf{Z}},$$
  
 $\bar{\mathbf{Z}} = \text{FFN}(\text{Norm}_1(\mathbf{Z})) + \mathbf{Z}.$  (28)

Here we depict the pre-norm structure [56]. The multi-head attention mechanism is used in Transformer to improve the representation performance. Let  $\tilde{\mathbf{Z}} = \operatorname{Norm}(\bar{\mathbf{Z}})$ , the multi-head attention

function is a weighted combination of outputs from different head as follows,

$$Multi-Head(\tilde{\mathbf{Z}}) = Concat(head_1, head_2, \dots, head_H) \mathbf{W}_O \in \mathbb{R}^{P \times d},$$

$$head_j = Attention(\tilde{\mathbf{Z}}\mathbf{W}_Q^j, \tilde{\mathbf{Z}}\mathbf{W}_K^j, \tilde{\mathbf{Z}}\mathbf{W}_V^j) \in \mathbb{R}^{P \times d_h},$$
(29)

where H is the number of heads.  $\mathbf{W}_Q^j, \mathbf{W}_K^j, \mathbf{W}_V^j \in \mathbb{R}^{d \times d_h}$  and  $\mathbf{W}_O \in \mathbb{R}^{d \times d}$  are learnable projections parameters with head dimensionality  $d_h = d/H$ . And self-attention function in vanilla Transformer is defined as follows.

Attention 
$$(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \operatorname{Softmax}(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_h}})\mathbf{V} \in \mathbb{R}^{P \times P},$$
 (30)

where Softmax is an exponential activation with  $l_1$  normalization [9] in the last dimensionality. And the attention scores between the i-th token and all other tokens after softmax operation are as follows,

$$\operatorname{Softmax}\left(\frac{\mathbf{Q}_{i}\mathbf{K}^{\top}}{\sqrt{d_{h}}}\right) = \frac{\exp\left(\frac{\mathbf{Q}_{i}\mathbf{K}^{\top}}{\sqrt{d_{h}}}\right)}{\sum_{j=1}^{n} \exp\left(\frac{\mathbf{Q}_{i}\mathbf{K}_{j}^{\top}}{\sqrt{d_{h}}}\right)} \in \mathbb{R}^{P}.$$
(31)

**Modification in MoFo.** However, we leverage Regulated Relaxation Function to instead the attention function in MoFo introduced in Section 3.2.2 as follows,

Attention 
$$(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{RRF}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_h}} + \log \mathcal{S}(\mathbf{\Gamma}; \alpha, \beta))\mathbf{V} \in \mathbb{R}^{P \times d_h},$$
  
where  $\mathbf{Q} = \mathbf{Z}\mathbf{W}_O^j \in \mathbb{R}^{P \times d_h}, \mathbf{K} = \mathbf{Z}\mathbf{W}_K^j \in \mathbb{R}^{P \times d_h}, \mathbf{V} = \mathbf{Z}\mathbf{W}_V^j \in \mathbb{R}^{P \times d_h},$  (32)

The feedforward networks in MoFo are gated linear units [40] with Swish activation.Let  $\dot{\mathbf{Z}} = \operatorname{Norm}_1(\mathbf{Z})$ , the  $\operatorname{FFN}(\cdot)$  is defined as follows,

$$FFN(\dot{\mathbf{Z}}) = \left(SwiGLU(\dot{\mathbf{Z}}\mathbf{W}_1 + \mathbf{b}_1) \odot (\dot{\mathbf{Z}}\mathbf{W}_2 + \mathbf{b}_2)\right)\mathbf{W}_3 + \mathbf{b}_3 \in \mathbb{R}^{P \times d},\tag{33}$$

where weight matrices  $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times 4d}, \mathbf{W}_3 \in \mathbb{R}^{4d \times d}$  and bias parameters  $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{4d}, \mathbf{W}_3 \in \mathbb{R}^d$  are learnable. The normalization layer is root mean square normalization [74] as follows,

$$\operatorname{Norm}_{1}(\mathbf{Z}) = \frac{\mathbf{Z} \odot \mathbf{g}_{1}}{\sqrt{\frac{1}{d} \sum_{k=1}^{d} \mathbf{Z}_{:,k}}} \in \mathbb{R}^{P \times d}, \quad \operatorname{Norm}_{2}(\bar{\mathbf{Z}}) = \frac{\bar{\mathbf{Z}} \odot \mathbf{g}_{2}}{\sqrt{\frac{1}{d} \sum_{k=1}^{d} \bar{\mathbf{Z}}_{:,k}}} \in \mathbb{R}^{P \times d}, \quad (34)$$

where  $\mathbf{g}_1, \mathbf{g}_2 \in \mathbb{R}^d$  are learnable scale parameters of normalization

# **D** Experiments

# D.1 Dataset Analysis

We further visualize the temporal correlations across multiple datasets. As shown in Figure 7, most datasets exhibit clear periodic patterns, with relatively high correlations between time steps separated by fixed intervals. This suggests that commonly adopted sequential input strategies—such as those used in patch-based methods—tend to group temporally adjacent but semantically unrelated time steps, thereby hindering the model's ability to capture intrinsic periodic structures.

### **D.2** Settings

Our experiment is based on the TFB platform [34] for a fair comparison. Following the settings in TFB [34], we report the best performance within the optional historical sequence length  $T \in \{96, 336, 512\}$  of the multiple forecasting length  $L \in \{96, 192, 336, 720\}$  with two common generic metrics including the mean square error (MSE) and mean absolute error (MAE) to judge the performance of our model. Our experiments are executed on an NVIDIA A100 with 40GB memory. Our code environment is based on the PyTorch framework using Python 3.11.5. The 'Drop Last' trick is closed to ensure a fair comparison [18]. We adopt Adam [14] optimizer. The training process is guided by the  $L_1$  loss function of the FreDF strategy [53]. The penalization distance parameter  $\beta > 0$  is restricted in (0, P). We utilize only one layer of Transformer with attention head H = 4 for each setting in all datasets. LTSF datasets often have multiple channels (or variates), and we adopt the channel independence approach [32] to simultaneous independent learning of all channels.

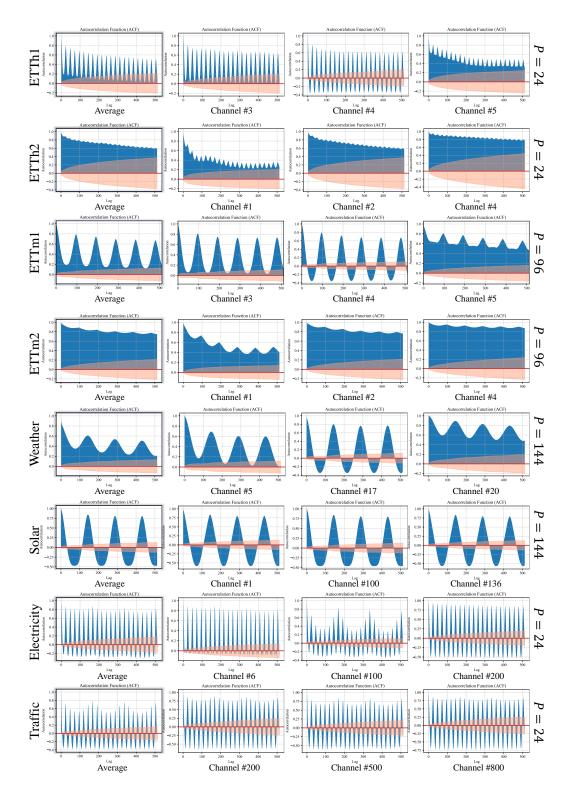


Figure 7: The autocorrelation visualization of all datasets.

Table 4: Performance comparisons for LTSF. The **best** results are marked in corresponding colors. All experimental results are selected from the best performance under the historical sequence length  $T \in \{96, 336, 512\}$ .

М	ethod		oFo urs)	FI (20	TS (23)		near (23)		esNet 23)	FEDfe (20	ormer (22)	Trifo (20		MI (20		FiI (20		Statio (20	onary 22)	Info	rmer 21)
M	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1			0.413 0.424	0.400 0.419	$0.418 \\ 0.435$	0.422 0.431		0.440 0.523	$0.443 \\ 0.487$	0.420 0.458	$0.444 \\ 0.466$	0.444 0.492	$0.449 \\ 0.479$	$0.400 \\ 0.428$	$0.430 \\ 0.447$	$0.405 \\ 0.434$	$0.416 \\ 0.435$	$0.615 \\ 0.632$	$0.540 \\ 0.551$	0.574 0.588	0.444 0.492
ETTh2	192 336	0.273 0.327 0.361 0.379	<b>0.373</b> 0.405	0.331 <b>0.350</b>	0.379 <b>0.396</b>	0.345 0.368	0.408	0.404 0.389	$0.413 \\ 0.435$	0.415 0.389	$0.428 \\ 0.457$	1.290 1.325	$0.768 \\ 0.781$	$0.419 \\ 0.474$	$0.439 \\ 0.475$	$0.358 \\ 0.372$	$0.401 \\ 0.425$	$0.379 \\ 0.358$	$0.418 \\ 0.413$	0.448 0.464	1.290 1.325
ETTm1	192 336	0.286 0.320 0.347 0.388	0.363 0.382	0.337 0.368	$0.365 \\ 0.384$	0.355 0.372	$0.379 \\ 0.385$	$0.392 \\ 0.423$	$0.404 \\ 0.426$	0.575 0.618	$0.516 \\ 0.544$	0.387 0.426	$0.410 \\ 0.446$	$0.336 \\ 0.370$	$0.369 \\ 0.391$	$0.339 \\ 0.374$	$0.365 \\ 0.385$	0.494 0.577	$0.451 \\ 0.490$	0.480 0.531	0.387 0.426
ETTm2	336	0.155 0.211 0.258 0.342	0.283 0.314	0.219 0.272	0.291 0.326	0.218 0.273	0.326	0.254 0.313	$0.310 \\ 0.345$	0.297 0.366	$0.360 \\ 0.400$	0.473 0.692	$0.453 \\ 0.549$	$0.232 \\ 0.303$	$0.313 \\ 0.367$	$0.220 \\ 0.277$	0.291 0.329	$0.338 \\ 0.432$	$0.373 \\ 0.416$	0.365 0.414	0.473 0.692
Weather	192 336		0.230 0.272	0.215 0.261	$0.261 \\ 0.295$	0.218 0.266		0.219 0.278	$0.262 \\ 0.302$	0.265 0.330	$0.334 \\ 0.372$	0.216 0.272	$0.277 \\ 0.324$	$0.214 \\ 0.259$	$0.271 \\ 0.309$	$0.218 \\ 0.266$	$0.263 \\ 0.295$	$0.240 \\ 0.322$	$0.290 \\ 0.328$	0.300 0.332	0.216 0.272
Solar	336		0.231 0.238	0.229 0.241	$0.267 \\ 0.273$	0.223 0.238		$0.206 \\ 0.208$	$0.276 \\ 0.284$	0.415 1.008	$0.477 \\ 0.839$	0.250 0.261	$0.295 \\ 0.297$	$0.226 \\ 0.259$	$0.284 \\ 0.308$	$0.226 \\ 0.241$	$0.257 \\ 0.265$	$0.400 \\ 0.414$	$0.386 \\ 0.394$	0.388 0.420	0.250 0.261
Electricity	336		0.234 0.252	0.154 0.170	$0.250 \\ 0.268$	0.155 0.171		$0.180 \\ 0.204$	$0.280 \\ 0.304$	0.203 0.221	$0.316 \\ 0.333$	0.209 0.225	$0.307 \\ 0.323$	$0.175 \\ 0.184$	$0.287 \\ 0.296$	$0.168 \\ 0.189$	$0.261 \\ 0.284$	$0.180 \\ 0.204$	$0.283 \\ 0.305$	0.362 0.416	0.209 0.225
Traffic	192 336	0.362 0.379 0.390 0.424	0.254 0.258	0.418 0.433	$0.294 \\ 0.308$	0.407 0.417	$0.277 \\ 0.282$	$0.613 \\ 0.626$	$0.322 \\ 0.332$	0.614 0.627	$0.381 \\ 0.389$	0.597 0.617	$0.325 \\ 0.332$	$0.526 \\ 0.545$	$0.302 \\ 0.307$	$0.415 \\ 0.430$	$0.285 \\ 0.299$	$0.611 \\ 0.628$	$0.338 \\ 0.342$	1.280 0.477	0.597 0.617

### **D.3** Performance Comparison with More Baselines

Considering readability, we only compared our approach with some representative SOTA baselines in Section 4.2. Here, we include more additional LTSF baselines to provide a more comprehensive evaluation of the performance of MoFo. Specifically, we add the following baselines: MLP-based models including FITS [68] NLinear [73] and FiLM [78]; TCN-based Models including TimesNet [63] and MICN [57]; Transformer-based models including FEDformer [79], Triformer [6], Non-stationary Transformer (Stationary) [23] and Informer [77]. As shown in Table 4, NLinear remains a powerful baseline as a linear model. FITS and FiLM extract time series representation in the frequency domain and in combining Legendre memory models, respectively. MICN utilizes convolutional networks to capture local and global contexts, while FEDformer enhances the Transformer in the frequency domain. Triformer introduces efficient triangular attention with convolutional down-sampling on coarse-time series, and Non-stationary Transformer (Stationary) focuses on addressing the nonstationarity of time series. Finally, Informer effectively tackles the computational challenges of long sequence prediction through its ProbSparse attention and distillation techniques. However, MoFo demonstrates superiority across almost all metrics comparing to all the baselines by reasonable Period-based Discrete Patching strategy. The Modulator in MoFo not only dynamically models the periodicity of time series data but also addresses permutation invariance, enhancing the representation capabilities of the Transformer in MoFo.

### D.4 Performance on None Periodicity Datasets

To further evaluate the generalization ability of MoFo beyond strictly periodic signals, we investigate its performance on datasets that lack explicit periodicity of a publicly available Influenza-Like Illness (ILI) dataset released by the U.S. Centers for Disease Control and Prevention (CDC), which contains weekly reports of the proportion of ILI-related visits from 2002 to 2021. The ILI dataset exhibits weak or no clear seasonality, making it an appropriate dataset for testing models under

non-periodic temporal conditions. We forecast future prediction using four prediction horizons:  $L \in \{24, 36, 48, 60\}$  with optional look-back window  $T \in \{36, 104\}$ .

To emphasize the modeling of non-periodic time series, we follow the preprocessing strategy of TimesNet [63] by identifying the main pseudo-period as the reciprocal of the dominant frequency derived via Fast Fourier Transform (FFT) for each input sequence. As shown in Table 5, although MoFo was originally designed to leverage explicit periodic structures, its superior results on ILI demonstrate strong adaptability and robustness in capturing complex temporal dependencies even in the absence of clear periodic patterns. This highlights that MoFo not only excels in periodic forecasting but also generalizes effectively to irregular, non-stationary time series domains.

Table 5: Performance comparison on non-periodicity dataset ILI	Table 5: Performance	comparison	on non-periodicity	dataset ILI
--	----------------------	------------	--------------------	-------------

Methods	Mo	Fo	FI	TS	TimesNet		
Metrics	MSE	MAE	MSE	MAE	MSE	MAE	
24	2.113	0.927	2.176	0.928	2.255	0.936	
36	1.952	0.924	2.166	0.993	2.132	0.940	
48	1.714	0.824	2.011	0.928	2.182	0.944	
60	1.800	0.906	2.010	0.967	2.169	0.940	

### **D.5** Ablation Study on Padding Strategy

Our padding strategy, detailed in Section 3.1 formalized in Eq. 1, is designed to preserve temporal continuity when the sequence length T is not an integer multiple of the detected period length P. Specifically, we start from the current time step and move backward to delineate periods of length P. Any prefix that does not form a complete period is left-padded with the leftmost elements of the nearest full period on its right. This scheme ensures that all time steps participate in subsequent computations without discarding boundary information. To evaluate the effect of this design, we introduce a variant termed '+ **Zero Padding**', where incomplete periods are instead filled with zeros. We conduct a comparative analysis on the ETTm1 and ETTm2 datasets under identical settings. The results are summarized in Table 6, where the best-performing metrics are highlighted in bold.

Table 6: Ablation Study on Padding Strategy

Me	ethods	Mo	oFo .	+ zeros	padding
M	Ietric	MSE	MAE	MSE	MAE
	96	0.286	0.335	0.292	0.345
Œ	192	0.320	0.363	0.328	0.372
ETTm1	336	0.347	0.382	0.348	0.386
Ш,	720	0.388	0.411	0.401	0.435
	96	0.155	0.240	0.156	0.247
E	192	0.211	0.283	0.215	0.293
ETTm2	336	0.258	0.314	0.261	0.318
Щ	720	0.342	0.368	0.345	0.372

Empirically, our proposed period-aware padding consistently outperforms zero padding across both benchmarks. We attribute this improvement to its ability to re-use the immediately preceding historical patterns, thereby maintaining local temporal coherence and facilitating smoother periodic transitions. In contrast, zero padding introduces abrupt discontinuities that disrupt the temporal rhythm, leading to inferior generalization. These findings confirm that our padding mechanism not only preserves data integrity but also serves as an implicit temporal regularizer that enhances long-term forecasting stability.

# **D.6** Look-back Window Sensitivity Experiments

In time series forecasting, the look-back window length—i.e., the number of historical steps fed into the model—is a critical hyperparameter that directly affects performance. Different architectures exhibit varying sensitivities to the length of historical context: models emphasizing long-term dependencies may benefit from longer input sequences, while those designed for short-term dynamics

might suffer from redundant or noisy inputs when the window is excessively long. Therefore, treating the look-back length as a tunable hyperparameter rather than a fixed setting is essential for a fair and comprehensive evaluation.

We perform experiments under multiple look-back window configurations to systematically assess this sensitivity. In the original setup, look-back lengths are predefined for each dataset, and models are trained and tested under all candidate configurations. The best-performing results are then reported, reflecting each model's optimal temporal receptive field. This protocol ensures a fair comparison among models with heterogeneous design principles and differing dependency ranges.

In line with this methodology, we examine the performance of our model and baselines under varying look-back window lengths. For the forecasting horizon of 720, we adopt  $\{96, 336, 512\}$  as candidate look-back windows, corresponding to commonly used temporal spans in long-term forecasting. The experimental results, summarized in Table 7, reveal that our model maintains consistently strong performance across different window lengths, demonstrating both its robustness and its capacity to adaptively leverage available historical information. These findings justify our choice of treating the look-back window length as a tunable hyperparameter in the main experiments and highlight the stability of MoFo under diverse temporal contexts.

Table 7: Look-back window sensitivity experiments of all optional look-back window length  $T \in \{96, 336, 512\}$  on the forecasting length setting L = 720.

Method	Me	oFo	FI	TS	DLi	near	TimesNet		
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
$ \begin{array}{c c} \hline E & T = 96 \\ T = 336 \\ T = 512 \end{array} $	0.447	<b>0.454</b>	0.547	0.518	0.515	0.511	0.521	0.495	
	0.459	<b>0.469</b>	0.475	0.487	0.471	0.493	0.542	0.519	
	0.443	0.463	<b>0.435</b>	<b>0.458</b>	0.464	0.488	0.560	0.531	
T = 96 $T = 336$ $T = 512$	0.416	0.442	0.439	0.452	0.650	0.571	0.434	0.448	
	0.393	0.428	0.397	0.431	0.704	0.597	0.472	0.480	
	0.379	0.425	0.382	0.425	0.786	0.623	0.480	0.468	

## **D.7** Period Sensitivity Experiments

The accurate identification of periodic structures plays a crucial role in MoFo's ability to capture long-term dependencies and recurrent temporal dynamics. However, in real-world time series, period selection is often uncertain due to noise, seasonal drift, or data heterogeneity. To comprehensively examine MoFo's robustness to such variations, we conduct three complementary sensitivity analyses:

• adjusting the given period length to test the impact of under- and over-estimation of period length,

• introducing multiple coexisting periodicities to evaluate the model's adaptability to mixed-period

- introducing multiple coexisting periodicities to evaluate the model's adaptability to mixed-period settings, and ② comparing robustness under complex dynamic period settings. The other experimental settings are kept consistent with the main experiments in Section 4.
- **O** Sensitivity to period length P. To assess the importance of accurate period calibration in MoFo, we conduct a series of sensitivity experiments on the given period length P. We first design two variants—"+ Half" and "+ Double"—in which the detected period is halved or doubled, respectively. As summarized in Table 8, results on three representative datasets with distinct periodicities demonstrate that a well-calibrated period is crucial for forecasting accuracy. When P is substantially under- or over-estimated, the model's ability to capture intrinsic temporal regularities degrades significantly, while the original configuration preserves the correct periodic structure and yields optimal performance. To further verify the robustness of this observation, we introduce four additional fine-grained perturbations: "+5%" and "+15%", where P is increased by 5% and 15%, and "-5%" and "-15%", where P is decreased by 5% and 15%. As shown in Table 8, even minor deviations from the detected period lead to measurable performance drops across datasets with different dominant periods. These results collectively highlight that the accuracy of period estimation plays a pivotal role in MoFo's ability to effectively model periodic dependencies in time series data.
- **Q** Sensitivity to multiple coexisting periods. In addition to single-period sensitivity, we further investigate MoFo's behavior under multiple periodic structures. We adopt the Traffic dataset, which exhibits possible two major periods—daily (P=24) and weekly (P=168) patterns. In practical multi-period scenarios, we typically use the shorter period as the baseline configuration. Our

Table 8: Sensitivity experiments of period length P.

Me	thod	Me	oFo	+5	5%	+1	5%	-5	%	-1:	5%	+ F	Ialf	+ Do	ouble
M	etric	MSE	MAE												
ETTm1	96	0.286	0.335	0.292	0.335	0.297	0.338	0.291	0.337	0.297	0.341	0.302	0.346	0.295	0.359
	192	0.320	0.363	0.329	0.368	0.336	0.370	0.333	0.369	0.344	0.386	0.335	0.379	0.339	0.388
	336	0.347	0.382	0.357	0.388	0.372	0.391	0.354	0.385	0.376	0.402	0.374	0.404	0.372	0.408
	720	0.388	0.411	0.406	0.420	0.403	0.416	0.402	0.415	0.411	0.420	0.414	0.425	0.419	0.429
Weather	96	0.141	0.186	0.143	0.192	0.155	0.206	0.157	0.221	0.146	0.188	0.152	0.199	0.157	0.206
	192	0.186	0.230	0.196	0.236	0.202	0.252	0.199	0.240	0.201	0.242	0.193	0.239	0.194	0.234
	336	0.233	0.272	0.238	0.277	0.245	0.286	0.235	0.280	0.243	0.286	0.241	0.279	0.244	0.284
	720	0.312	0.331	0.335	0.358	0.323	0.358	0.332	0.352	0.335	0.355	0.348	0.363	0.342	0.355
Traffic	96	0.362	0.247	0.388	0.266	0.395	0.265	0.383	0.255	0.382	0.265	0.376	0.258	0.370	0.252
	192	0.379	0.254	0.397	0.275	0.408	0.287	0.398	0.284	0.409	0.288	0.380	0.262	0.388	0.275
	336	0.390	0.258	0.395	0.260	0.405	0.275	0.395	0.266	0.417	0.294	0.406	0.286	0.395	0.283
	720	0.424	0.281	0.424	0.281	0.443	0.298	0.437	0.283	0.443	0.303	0.434	0.297	0.438	0.302

experiments show that this choice enables MoFo to capture fine-grained temporal variations while retaining stable performance across longer periods, as shown in Table 9. To explore whether integrating multiple MoFo models could further improve prediction, we introduce three variants: (1) MoFo-1, trained with a 24-hour (daily) period (default setting of MoFo); (2) MoFo-7, trained with a 7-day (weekly) period; and (3) Mix-MoFo, which ensembles both models and combines their outputs through a learnable fusion layer. The results indicate that using the shorter period configuration provides the most precise periodic alignment, while the ensemble model further stabilizes performance under complex multi-period signals. These findings confirm MoFo's robustness and adaptability when modeling time series with heterogeneous or nested periodic patterns.

Table 9: Sensitivity experiments of multiple coexisting periods.

Me	ethod	Mo	Fo-1	Mix-	MoFo	Mol	Fo-7
M	etric	MSE	MAE	MSE	MAE	MSE	MAE
Traffic	96 192 336 720	0.362 0.379 0.390 0.424	0.247 0.254 0.258 0.281		0.254 0.255 0.263 0.285		0.257 0.262 0.272 0.289

**Sensitivity to dynamic periods.** To further evaluate robustness of MoFo under complex temporal settings, we conduct experiments on datasets with mixed and time-varying periodicities. To systematically investigate this issue, we construct a synthetic dataset named Mixed-ETT by combining two standard datasets, ETTh1 (period = 24) and ETTm1 (period = 96). Specifically, each dataset is evenly divided into four temporal segments, and these segments are alternately concatenated along the time axis to form the new sequence as Mixed-ETT =  $\{\text{ETTh1}[: \frac{1}{4}], \text{ETTm1}[\frac{1}{4}: \frac{1}{2}], \text{ETTh1}[\frac{1}{2}: \frac{1}{2}: \frac{1}{2}], \text{ETTh1}[\frac{1}{2}: \frac{1}{2}: \frac{1}{2}], \text{ETTh1}[\frac{1}{2}: \frac{1}{2}: \frac{1}{2}], \text{ETTh1}[\frac{1}{2}: \frac{1}{2}: \frac{1}{2}: \frac{1}{2}: \frac{1}{2}], \text{ETTh1}[\frac{1}{2}: \frac{1}{2}: \frac{1$  $\frac{3}{4}$ , ETTm1 $\left[\frac{3}{4}:1\right]$ . This design produces a dataset with alternating periodic structures of 24, 96, 24, 96, mimicking the temporal heterogeneity commonly observed in real applications. The other experimental settings are kept consistent with the main experiments in Section 4. We compare MoFo against two representative baselines, DUET and DLinear. As shown in Table 10, experimental results demonstrate that our straightforward MoFo implementation achieved strong performance. When dealing with multi-period time series, we set the smallest period length as our baseline configuration. The "Mix-MoFo" variant added assumptions regarding additional periods, which increased the risk of overfitting and consequently led to a decline in performance. These results confirm that MoFo can effectively handle dynamic and mixed periodic behaviors without explicit retraining or manual period adjustment.

### D.8 Efficiency Analysis

We compare the computational efficiency of Transformer-based baseline models on the Solar dataset with L=96. As shown in Table 11, MoFo achieves the least MSE among all Transformer-based models, while demonstrating the least parameters number with fastest training speed. All the benefits of MoFo arise from its effective Period-based Discrete Patching strategy, which reduces the complexity of the Transformer to quadratic in relation to the period length, while utilizing only

Table 10: Sensitivity experiments of period-varying synthetic dataset.

Method	Mo	ъFо	DU	ET	DLinear			
Metric	MSE	MAE	MSE	MAE	MSE	MAE		
H 96 H 192 336 720	0.178 0.185 0.180 0.191	0.223 0.232 0.237 0.253	0.186 0.184 0.187 0.218	0.229 0.241 0.247 0.278	0.194 0.190 0.196 0.219	0.232 0.247 0.256 0.282		

a single layer of the Transformer, which is sufficient. Similarity, the competitive baseline DUET with dual clusting strategy on both temporal dimension and channel dimension exhibits over  $14\times$  increase in parameters number,  $4\times$  higher FLOPs, and slower training speed. Compared to PatchTST, a classic Transformer-based model that employs a successive patching strategy on the time series, MoFo reduces the number of MACs by more than  $80\times$ , accelerates the training speed by over  $10\times$ , and significantly decreases both the parameter requirements and memory usage. This is primarily because PatchTST stacks multiple layers of the Transformer, which is necessary for its architecture, yet lacks a reasonable positional encoding for time series data. Although iTransformer, FEDformer, and Informer have fewer computations than MoFo, they require longer training times as well as larger parameters, and their performance lags behind MoFo by up to  $1.8\times$  since their complexity architectures.

Table 11: Efficiency comparison of MoFo and SOTA baselines with L=720 in Traffic dataset. All results of each model are under the optimal hyperparameters for fair comparison. Parameters: All learnable parameters requiring gradient descent. MACs: multiply–accumulate operations. FLOPs: floating point operations. M: Million  $(10^6)$ . B: Billion  $(10^9)$ . T: Trillion  $(10^{12})$ . MB: Megabyte. s: Second.  $\uparrow$  indicates the relative percentage increasing regarding MoFo and  $\downarrow$  indicates the relative percentage decreasing.

	Models	MSE	# Parameters	# MACs	# FLOPs	Memory Usage	Epoch Time
$\mathrm{Solar}\left[L=96\right]$	Informer	0.368 117.75%	2.26 M <sub>↑527.78%</sub>	7.13 B <sub>↓2.72%</sub>	7.18 B <sub>\$\psi_38.58\%</sub>	852 MB <sub>↓75.36%</sub>	58 s <sub>↑222.22%</sub>
	Stationary	$0.365_{\uparrow 115.97\%}$	11.2 M <sub>↑3011.11%</sub>	39.56 B <sub>↑439.70%</sub>	41.32 B <sub>↑253.46%</sub>	1,710 MB <sub>\$\psi_50.54\%</sub>	24 s <sub>†33.33%</sub>
	Triformer	$0.225_{\uparrow 33.14\%}$	1.62 M <sub>↑350.00%</sub>	33.56 B <sub>↑357.84%</sub>	38.45 B <sub>↑228.91%</sub>	11,714 MB <sub>238.75</sub> %	147 s <sub>116.67</sub> %
	FEDformer	$0.485_{\substack{\uparrow186.98\%}}$	3.63 M <sub>↑908.33%</sub>	1.86 B <sub>↓74.62%</sub>	1.52 B <sub>\$\to\$86.99\%</sub>	858 MB <sub>\$\psi75.19\%</sub>	204 s <sub>1033.33</sub> %
	Crossformer	0.183	3.82 M <sub>↑961.11%</sub>	159.3 B <sub>2073.26</sub> %	166.1 B <sub>↑1320.87%</sub>	10,698 MB <sub>209.37</sub> %	101 s <sub>↑461.11%</sub>
	PatchTST	$0.170_{\uparrow 0.59\%}$	2.62 M <sub>↑627.78%</sub>	607.4 B <sub>↑8186.49%</sub>	644.9 B <sub>↑5416.68%</sub>	29,726 MB <sub>↑759.63%</sub>	215 s <sub>1094.44</sub> %
	Pathformer	0.218 +28.99%	5.72 M <sub>1488.89</sub> %	24.84 B <sub>238.88</sub> %	27.99 B <sub>↑139.44%</sub>	12,754 MB <sub>268.83</sub> %	614 s <sub>†3311.11%</sub>
	iTransformer	$0.190_{\uparrow 12.42\%}$	0.51 M <sub>↑41.67%</sub>	2.32 B <sub>\$\delta 68.34\%</sub>	$2.66~B_{\downarrow 77.24\%}$	788 MB <sub>↓77.21%</sub>	18 s <sub>↑0.00%</sub>
	PDF	$0.181_{\uparrow 7.10\%}$	5.82 M <sub>↑1516.67%</sub>	204.1 B <sub>2684.45</sub> %	208.0 B <sub>↑16790.30%</sub>	5,616 MB <sub>162.41%</sub>	25 s <sub>†38.89%</sub>
	DUET	0.169 <sub>\tau0.00\%</sub>	5.64 M <sub>1466.67%</sub>	19.23 B <sub>↑162.35</sub> %	60.88 B <sub>↑420.79%</sub>	4,422 MB <sub>27.88</sub> %	35 s <sub>†94.44%</sub>
	MoFo	0.169	0.36 M	<b>7.33</b> B	11.69 B	3,458 MB	18 s

# **E** Discussion and Future Work

MoFo performs discrete patching based on periodic position, grouping the most correlevant time steps within the same patch for prioritized interaction. However, our periodic positional embedding provided by the Regulated Relaxation Function is currently only applicable to vanilla self-attention mechanisms, which exhibit quadratic complexity (although our specific implementation's complexity is quadratic with respect to the period length P rather than the sequence length T with  $P \ll T$ ). Exploring how to adapt this method to attention mechanisms with linear complexity is a direction worthy of future investigation. Simultaneously, applying the core ideas of MoFo within time series LLM foundational models to empower their development in time series learning is another promising avenue for subsequent research.

# **NeurIPS Paper Checklist**

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

# IMPORTANT, please:

- Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",
- · Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.

#### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We propose a novel transformer-based model in long-term time series forecasting.

# Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

# 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We conducted a discussion about the limitation as future works at the end of the Appendix.

### Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We theoretical prove the the effectiveness of the regulater relaxation function in our model. The proofs are fully demonstrated in the first section in Appendix.

#### Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We have open-sourced our code via an anonymous link for reproducibility, and provide detailed experimental settings in the corresponding section.

### Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We have open-sourced our code via an anonymous link for reproducibility in the corresponding section.

### Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
  possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
  including code, unless this is central to the contribution (e.g., for a new open-source
  benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We report the detailed settings and dataset processing details in the experiments section.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
  material.

### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We only report the mean results from multiple experiments for all experiments for ensuring readability.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We implement our proposed model on an 40GB NVIDIA A100 GPU with Pytorch.

### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <a href="https://neurips.cc/public/EthicsGuidelines">https://neurips.cc/public/EthicsGuidelines</a>?

Answer: [Yes]

Justification: The datasets involved in the paper are all open source and widely used datasets.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The proposed model significantly enhances performance in long-term time seires forecasting scenarios, offering positive implications for a wide range of downstream applications. No notable negative side effects are observed.

#### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

# 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper does not refer to high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets).

### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
  necessary safeguards to allow for controlled use of the model, for example by requiring
  that users adhere to usage guidelines or restrictions to access the model or implementing
  safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

# 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All datasets and code used in this study are publicly available.

#### Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Guidelines:

Justification: Yes, our assets (code and data) are accessible via an anonymous link during the review process. Upon acceptance, they will be made publicly available for open access.

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.

 At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

# 14. Crowdsourcing and research with human subjects

Ouestion: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No crowdsourcing nor research with human subjects are involved.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Ouestion: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLM is only used for language polishing of papers to improve readability.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.