
First De-Trend then Attend: Rethinking Attention for Time-Series Forecasting

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Transformer-based models have gained large popularity and demonstrated promis-
2 ing results in long-term time-series forecasting in recent years. In addition to
3 learning attention in time domain, recent works also explore learning attention in
4 frequency domains (e.g., Fourier domain, wavelet domain), given that seasonal pat-
5 terns can be better captured in these domains. In this work, we seek to understand
6 the relationships between attention models in different time and frequency domains.
7 Theoretically, we show that attention models in different domains are equivalent
8 under linear conditions (i.e., linear kernel to attention scores). Empirically, we ana-
9 lyze how attention models of different domains show different behaviors through
10 various synthetic experiments with seasonality, trend and noises, with emphasis
11 on the role of softmax operation therein. Both these theoretical and empirical
12 analyses motivate us to propose a new method: TDformer (Trend Decomposition
13 Transformer), that first applies seasonal-trend decomposition, and then additively
14 combines an MLP which predicts the trend component with Fourier attention which
15 predicts the seasonal component to obtain the final prediction. Extensive experi-
16 ments on benchmark time-series forecasting datasets demonstrate that TDformer
17 achieves state-of-the-art performance against existing attention-based models.

18 1 Introduction

19 Transformer [18] recently gains wide popularity in time-series forecasting, inspired by its success
20 in natural language processing and its ability to capture long-range dependencies [19]. Apart from
21 the vanilla Transformer that calculates attention in time domain, recently variants of Transformer
22 which calculate attention in frequency domains (e.g., Fourier domain or wavelet domain) (Figure 2)
23 [22, 21, 24, 20, 14] have also been proposed to better model global characteristics of time series.

24 Despite the progress made by Transformer-based methods for time series forecasting, there lacks
25 a rule of thumb to select the domain in which attention is best learned. Our work is driven by
26 better understanding the following research question: *Does learning attention in one domain offer
27 better representation ability than the other? If so, how?* We show mathematically that under
28 linear conditions, learning attention in time or frequency domains leads to equivalent representation
29 power. We then show that due to the softmax non-linearity used for normalization, this theoretical
30 linear equivalence does not hold empirically. In particular, attention models in different domains
31 demonstrate different empirical advantages. This finding sheds light on how to best apply attention
32 models under different practical scenarios. We propose TDformer based on these insights and
33 demonstrate that we achieve state-of-the-art performance against current attention-based models.

34 More specifically, we find that (1) for *data with strong seasonality*, frequency-domain attention models
35 are more *sample-efficient* compared with time-domain attention models, as softmax with exponential
36 terms correctly amplify the dominant frequency modes in Fourier space. (2) For *data with trend*,

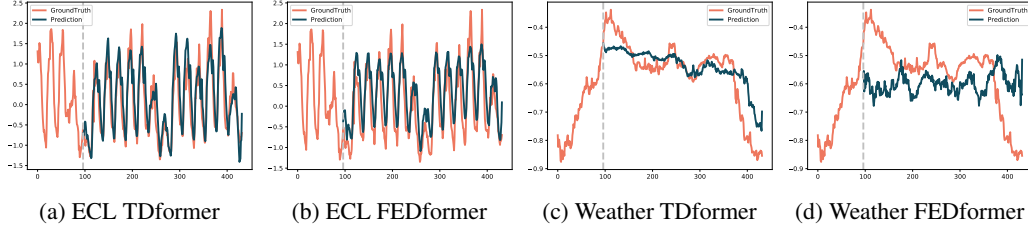


Figure 1: Prediction comparison between TDformer and FEDformer on electricity dataset ((a),(b)) and weather dataset ((c), (d)). We predict the future 336 steps given context 96 steps (the gray dash line). Orange line represents the ground truth, and blue line represents the prediction.

37 attention models generally show inferior *generalizability*, as attention models by nature interpolate
 38 rather than extrapolate the context. This finding of difference in performances of attention models
 39 on various types of time series data emphasizes the importance of seasonal-trend decomposition
 40 module in the attention model framework. (3) For *data with noisy spikes*, frequency-domain attention
 41 models are more *robust* to such spiky data, as large-value spikes in the time domain correspond to
 42 small-amplitude high-frequency modes, whose attention would be filtered out by softmax operations.

43 Due to the different performances of the various attention modules on data with seasonality and
 44 trend, we propose TDformer that first decomposes the context time series into trend and seasonal
 45 components. We use a MLP for predicting the future trend, Fourier attention to predict the future
 46 seasonal part, and add these two components to obtain the final prediction. Extensive experiments
 47 on benchmark forecasting datasets demonstrate the effectiveness of our proposed approach. As a
 48 motivating example, we visualize predictions of TDformer and one of the best performing baselines
 49 FEDformer in Figure 1. On data with strong seasonality (Figure 1a and Figure 1b) TDformer
 50 preserves both the seasonality and trend of the original data, while FEDformer [24] deviates from the
 51 trend of the ground truth. On data with strong trend (Figure 1c and Figure 1d), TDformer generates
 52 predictions that better follow the trend of the original data.

53 In summary, our contributions are:

- 54 • We theoretically show that under linear conditions, attention models in time domain, Fourier
 55 domain and wavelet domain have the same representation power;
- 56 • We empirically analyze attention models in different domains with synthetic data of different
 57 characteristics, given the non-linearity of softmax. We show that frequency-domain attention
 58 performs the best on data with seasonality, and attention models in general have inferior
 59 generalizability on trend data, which motivates the design of a hybrid model based on
 60 seasonal trend decomposition;
- 61 • We propose TDformer that separately models the trend with MLP and seasonality with
 62 Fourier attention, and shows state-of-the-art performance against current attention models
 63 on time-series forecasting benchmarks.

64 2 Related Work

65 **Time-Domain Attention Forecasting Models.** Informer [22] proposes efficient ProbSparse self-
 66 attention mechanism. Autoformer [21] renovates time-series decomposition as a basic inner block and
 67 designs Auto-Correlation mechanism for dependencies discovery. Non-stationary Transformer [14]
 68 proposes Series Stationarization and De-stationary Attention to address over-stationarization.

69 **Frequency-Domain Attention Forecasting Models.** FEDformer [24] proposes Fourier and wavelet
 70 enhanced blocks based on Multiwavelet-based Neural Operator Learning [4] to capture important
 71 structures in time series through frequency domain mapping. ETSformer [20] selects top-K largest
 72 amplitude modes as frequency attention and combines with exponential smoothing attention. Adaptive
 73 Fourier Neural Operator (AFNO) [3] builds upon FNO [13] and proposes an efficient token mixer that
 74 learns to mix in the Fourier domain. FNet [11] replaces the self-attention with Fourier Transform and
 75 promotes efficiency without much loss of accuracy on NLP benchmarks. T-WaveNet [15] constructs
 76 a tree-structured network with each node built with invertible neural network (INN) based wavelet
 77 transform unit for iterative decomposition. Adaptive Wavelet Transformer Network (AWT-Net) [5]

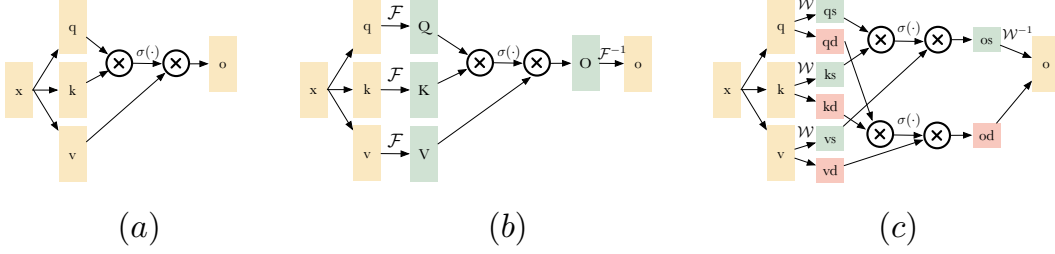


Figure 2: Comparison between (a) time attention, (b) Fourier attention and (c) wavelet attention. For simplicity, we only draw one layer of multiwavelet decomposition/reconstruction, and similar analysis follows for multiple layers. See precise notations in Section 3.

78 generates wavelet coefficients to classify each point into high or low sub-bands components and
 79 exploits Transformer to enhance the original shape features.

80 **Decomposition-Based Forecasting Models** decompose time series into trend and seasonality (with
 81 i.e., STL decomposition [2]). Apart from attention-based Autoformer and FEDformer, N-BEATS [16]
 82 models trend with small-degree polynomials and seasonality with Fourier series. N-HiTS [1] redefines
 83 N-BEATS by enhancing its input decomposition via multi-rate data sampling and its output synthesizer
 84 via multi-scale interpolation. FreDo [17] incorporates frequency-domain features into AverageTile
 85 model that averages history sub-series. FiLM [23] applies Legendre Polynomials projections to
 86 approximate historical information and Fourier projection to remove noise. DeepFS [6] encodes
 87 temporal patterns with self-attention and predicts Fourier series parameters and trend with MLP.

88 Despite the success of attention models in time, Fourier, and wavelet domains, there is still a lack
 89 of notion for understanding their relationships and respective advantages. Decomposition-based
 90 methods also adopt decomposition layers without giving strong reasoning for their necessity. We
 91 propose to fill this gap from both theoretical and empirical perspectives, and based on these analysis
 92 build a new framework that shows better forecasting performance.

93 3 Linear Equivalence of Attention in Various Domains

94 3.1 Formulation of Attention Models

95 We first briefly introduce the canonical Transformers. Denote input queries, keys and values as
 96 $\mathbf{q} \in \mathbb{R}^{L \times D}$, $\mathbf{k} \in \mathbb{R}^{L \times D}$, $\mathbf{v} \in \mathbb{R}^{L \times D}$, which are transformed from input \mathbf{x} through linear embeddings.
 97 Denote output of attention module as $\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) \in \mathbb{R}^{L \times D}$. As shown in Figure 2 (a), The attention
 98 operation in canonical attention is formulated as

$$\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \sigma \left(\frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d_q}} \right) \mathbf{v}, \quad (1)$$

99 where d_q is the dimension for queries that serves as normalization term in attention operation,
 100 and $\sigma(\cdot)$ represents activation function. When $\sigma(\cdot) = \text{softmax}(\cdot)^1$, we have *softmax attention*:
 101 $\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \text{softmax}(\mathbf{q}\mathbf{k}^T / \sqrt{d_q}) \mathbf{v}$. When $\sigma(\cdot) = \text{Id}(\cdot)$ (identity mapping), we have *linear*
 102 *attention*: $\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \mathbf{q}\mathbf{k}^T \mathbf{v}$ (we ignore the normalization term $\sqrt{d_q}$ for simplicity).
 103

104 **Definition 3.1** (Time Attention). Equation 1 refers to time domain attention where $\mathbf{q}, \mathbf{k}, \mathbf{v}$ are all in
 105 original time domain, shown in Figure 2 (a).

106 **Definition 3.2** (Fourier Attention). Fourier attention first converts queries, keys, and values with
 107 Fourier Transform, performs a similar attention mechanism in the frequency domain, and finally
 108 converts the results back to the time domain using inverse Fourier transform, shown in Figure 2 (b).
 109 Let $\mathcal{F}(\cdot)$, $\mathcal{F}^{-1}(\cdot)$ denote Fourier transform and inverse Fourier transform, then Fourier attention is
 110 $\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \mathcal{F}^{-1} \left(\sigma(\mathcal{F}(\mathbf{q})\overline{\mathcal{F}(\mathbf{k})}^T / \sqrt{d_q}) \mathcal{F}(\mathbf{v}) \right)$.

111 **Definition 3.3** (Wavelet Attention). Wavelet transform applies wavelet decomposition and recon-
 112 struction to obtain signals of different scales. Wavelet attention performs attention calculation to

¹ $\text{softmax}(\mathbf{x}) = \frac{e^{x_i}}{\sum_i e^{x_i}}$

113 decomposed queries, keys, and values in each scale, and reconstructs the output from attention results
 114 in each scale, illustrated in Figure 2 (c). Let $\mathcal{W}(\cdot)$, $\mathcal{W}^{-1}(\cdot)$ denote wavelet decomposition and wavelet
 115 reconstruction, then wavelet attention is $\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \mathcal{W}^{-1}\left(\sigma\left(\mathcal{W}(\mathbf{q})\mathcal{W}(\mathbf{k}^T)/\sqrt{d_q}\right)\mathcal{W}(\mathbf{v})\right)$.

116 3.2 Linear Equivalence of Time, Fourier and Wavelet Attention

117 In this section we formally prove that time, Fourier and wavelet attention models are equivalent under
 118 linear attention case.

119 **Lemma 3.1.** When $\sigma(\cdot) = \text{Id}(\cdot)$ (linear attention), time, Fourier and wavelet attention are equivalent.

120 *Proof.* Let $\mathbf{W} = \left(\frac{\omega^{jk}}{\sqrt{L}}\right) \in \mathbb{C}^{L \times L}$, $\omega = e^{-\frac{2\pi j}{L}}$ denote the Fourier matrix, then Fourier transform to
 121 signal $\mathbf{x} \in \mathbb{R}^{L \times D}$ can be expressed as $\mathbf{X} = \mathbf{W}\mathbf{x}$, $\mathbf{X} \in \mathbb{C}^{L \times D}$, and inverse Fourier transform can be
 122 expressed as $\mathbf{x} = \mathbf{W}^H \mathbf{X}$, where \mathbf{W}^H is the Hermitian (conjugate transpose) of \mathbf{W} . Given properties
 123 of Fourier matrix, we could easily show that

$$\mathbf{W}^{-1} = \mathbf{W}^H, \mathbf{W}^T = \mathbf{W}. \quad (2)$$

124 Following this expression, Fourier domain linear attention can be written as

$$\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \mathbf{W}^H[(\mathbf{W}\mathbf{q})(\overline{\mathbf{W}\mathbf{k}})^T(\mathbf{W}\mathbf{v})] = \mathbf{q}\mathbf{k}^T\mathbf{v}. \quad (3)$$

126 Therefore, calculating attention in Fourier domain is equivalent to time-domain attention.

128 For wavelet attention, we take single-scale wavelet decomposition and reconstruction as an example,
 129 and multi-scale wavelet transform follows the same analysis. Using the same notation, let $\mathbf{W} \in$
 130 $\mathbb{R}^{L \times \frac{L}{2}}$, $\mathbf{W}^{-1} \in \mathbb{R}^{\frac{L}{2} \times L}$ denote the wavelet decomposition and reconstruction matrix, then wavelet
 131 decomposition to signal $\mathbf{x} \in \mathbb{R}^{L \times D}$ can be expressed as $\mathbf{X} = \mathbf{W}\mathbf{x}$, $\mathbf{X} \in \mathbb{R}^{\frac{L}{2} \times D}$, and wavelet
 132 reconstruction can be expressed as $\mathbf{x} = \mathbf{W}^{-1}\mathbf{X}$. Since wavelet matrix is orthogonal, we have the
 133 property that $\mathbf{W}^T\mathbf{W} = \mathbf{I}$. Wavelet linear attention is

$$\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \mathbf{W}^{-1}[(\mathbf{W}\mathbf{q})(\mathbf{W}\mathbf{k})^T(\mathbf{W}\mathbf{v})] = \mathbf{q}\mathbf{k}^T\mathbf{v}, \quad (4)$$

134 which is again equivalent to time-domain attention. Therefore, we show that mathematically, time,
 135 Fourier and wavelet attention models are equivalent given linear assumptions. \square

137 4 Investigation on the Role of Softmax

138 Although these attention models are equivalent given linear assumptions, in practice we apply
 139 softmax as normalization, which changes the behavior of different attention models. In this section,
 140 we empirically analyze how softmax causes such performance gaps on datasets with three different
 141 representative properties: seasonality, trend and noises. For all experiments in this section, the task
 142 is to predict the next 96 time steps given history 96 time steps. We implement the wavelet-domain
 143 attention model based on multiwavelet transform model [4].

144 4.1 Data with Seasonality

145 **For data with fixed seasonality, Fourier attention is the most sample-efficient.** We use $\sin(x)$ as
 146 an example of seasonal data (visualized in Figure 3a and Figure 3b). There exist dominant frequency
 147 modes for data with seasonality. We visualize linear attention (Figure 3c) and softmax attention
 148 (Figure 3d) in Fourier space. Attention scores are concentrated on the dominant frequency mode. As
 149 softmax with exponential terms has the ‘‘polarization’’ effect (increasing the gap between large and
 150 small values), softmax attention further concentrates the scores on the dominant frequency, helping
 151 the model to better capture seasonal information. Therefore, we find that frequency-domain attention
 152 models are capable of quickly recognizing the dominant frequency modes (more sample efficient)
 153 compared with time-domain models (Figure 4a).

154 To further illustrate such polarization effect, we also compare softmax attention with polynomial
 155 kernels $\sigma(x) = x_i^d / \sum_i x_i^d$, where d is the degree of polynomials (without loss of generality we
 156 assume $x_i > 0, \forall i$). Polarization effect increases with respect to polynomial degrees. As shown
 157 in Figure 4c, the performance also increases as we increase the polarization effect and approaches

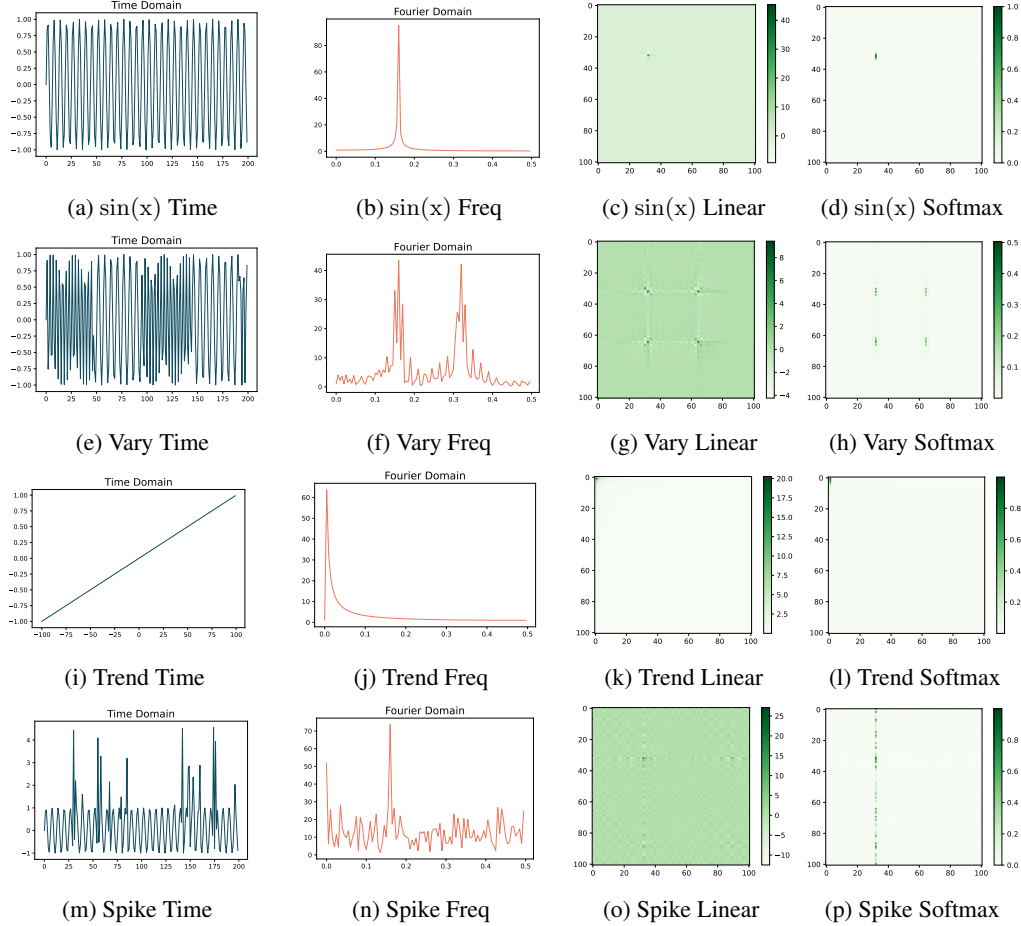


Figure 3: (a)-(d): Data with fixed seasonality: $\sin(x)$. Fourier softmax attention amplifies the correct frequency modes compared with Fourier linear attention. (e)-(h): Data with varying seasonality. Fourier softmax attention amplifies the dominant frequency modes, but also neglects the small-amplitude modes that embed the localized frequency information. (i)-(l): Data with linear trend. Fourier softmax attention incorrectly amplifies the low-frequency modes compared with Fourier linear attention. (m)-(p): Data with spikes as noises. Fourier softmax attention filters out the noisy components and emphasizes the correct frequency modes compared with Fourier linear attention.

Table 1: MSE and MAE of attention models and MLP with linear-trend data.

Metric	Time	Fourier	Wavelet	MLP
MSE	3.157 ± 0.435	8.567 ± 0.487	2.327 ± 0.689	$\mathbf{0 \pm 0}$
MAE	1.741 ± 0.121	2.880 ± 0.073	1.477 ± 0.239	$\mathbf{0.006 \pm 0.003}$

158 the performance of softmax operations. We also notice that apart from the polarization effect from
 159 exponential terms, normalization itself also introduces performance gaps between different attention
 160 models. The possible reason is that it's easier to optimize in the sparse Fourier domain compared
 161 with time domain. We leave this as our future explorations.

162 **For data with varying seasonality, wavelet attention is the most effective.** We use alternating
 163 $\sin(x)$ and $\sin(2x)$ as an example of varying seasonal data (visualized in Figure 3e and Figure 3f).
 164 The Fourier representation has both dominant modes as well as small-amplitude modes, where the
 165 latter embeds the varying-seasonality information. The Fourier softmax attention correctly amplifies
 166 the dominant frequency modes, but at the same time neglects the small-amplitude modes that convey
 167 the information of varying seasonality. By contrast, wavelet attention combines multi-scale time-
 168 frequency representation, and provides better localized frequency information. As shown in Figure 4b,
 169 wavelet attention is the most effective for varying-seasonality data.

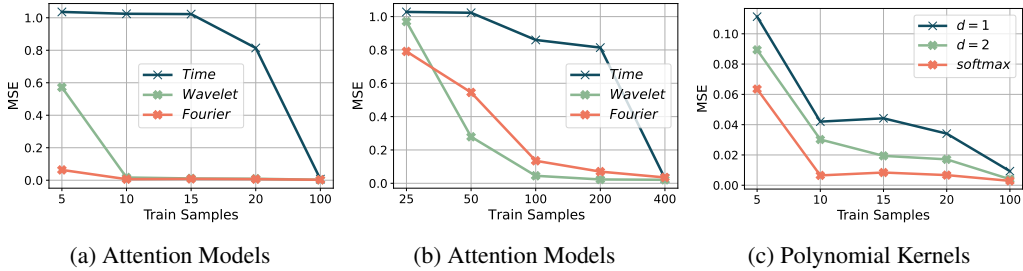


Figure 4: (a): Sample efficiency comparison of time, Fourier, wavelet attention models on data with fixed seasonality ($\sin(x)$). Fourier attention models are more sample-efficient. (b): Sample efficiency comparison on data with varying seasonality (alternating $\sin(x)$ and $\sin(2x)$). Wavelet attention models are more sample-efficient. (c): Sample efficiency comparison between polynomial kernels and softmax. Polarization effect increases with respect to the degree of polynomial kernels and approaches the softmax performance.

Table 2: MSE and MAE of different attention models with spiky data.

Metric	Time	Fourier	Wavelet
MSE	0.303 ± 0.002	0.019 ± 0.003	0.030 ± 0.008
MAE	0.495 ± 0.001	0.111 ± 0.010	0.137 ± 0.021

170 4.2 Data with Trend

171 **For data with trend, all attention models show inferior generalizability, especially Fourier**
 172 **attention.** We take linear trend data as an example (Figure 3i and Figure 3j) and evaluate different
 173 attention models. The first several frequency modes in Fourier space carry large values; the attention
 174 scores hence mostly focus on the first few frequency modes (top-left corner of Figure 3k). With
 175 the polarization effect of softmax, attention scores emphasize even more on these low-frequency
 176 components (Figure 3l) and generate misleading reconstruction results. We evaluate different attention
 177 models in Table 1. Fourier attention, with inappropriate polarization, leads to the largest errors.

178 Moreover, all these attention models fail to extrapolate linear trend well and suffer from large errors,
 179 since attention mechanism by nature works through interpolating the context history. By contrast,
 180 MLP perfectly predicts such trend signals, as shown in Table 1. This motivates us to decompose the
 181 time series into trend and seasonality [21, 24], apply attention mechanism only for seasonality, and
 182 use MLP for modeling trend.

183 4.3 Data with Spikes

184 **For data carrying noises, Fourier attention is the most robust.** We randomly inject large-value
 185 spikes into the training set of $\sin(x)$ as a motivating example (Figure 3m and Figure 3n). Spikes
 186 which have large values in time domain result in small-amplitude frequency components after
 187 Fourier transforms. With the polarization effect of softmax, time-domain softmax attention focuses
 188 incorrectly on large-value spikes while Fourier-domain softmax attention correctly filters out the noisy
 189 components and attends to the dominant frequency modes induced by $\sin(x)$. Comparing Figure 3o
 190 and Figure 3p, linear attention still distributes attention to the noisy frequency modes, while softmax
 191 attention mostly focuses correctly on the dominant frequency modes. Therefore, frequency-domain
 192 attention models are more robust to spikes, as shown in Table 2. All these analysis on datasets with
 193 different characteristics help guide our model design in the next section.

194 5 Our Method: TDformer

195 The performance difference in data with various characteristics motivates our model design. For data
 196 with seasonality, Fourier softmax attention amplifies dominant frequency modes and demonstrates the
 197 best performance. For data with trend, Fourier softmax incorrectly attends to only the low-frequency
 198 modes and produces large errors. Meanwhile, all attention models which work through interpolating
 199 the historical context, do not generalize well on trend data compared with MLP. These analyses

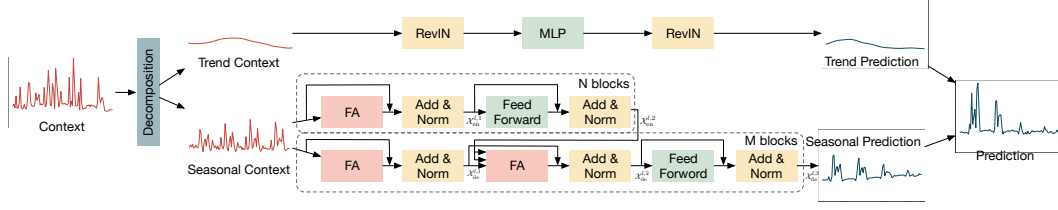


Figure 5: TDformer. We first apply seasonal trend decomposition to decompose the context time series into trend part and seasonal part. We adopt MLP to predict the trend part, and Fourier Attention (FA) model to predict the seasonal part, and add two parts together for final prediction.

200 motivate us to decompose time series into trend and seasonality, use Fourier attention to predict the
 201 seasonal part and MLP to predict the trend part. Figure 5 overviews our proposed model architecture.

202 We first decompose the time series into trend parts and seasonal parts following FEDformer [24].
 203 More specifically, we apply multiple average filters with different sizes to extract different trend
 204 patterns, and apply adaptive weights to combine these patterns into the final trend component. The
 205 seasonal component is acquired by subtracting trend from the original time series:

$$\mathbf{x}_{\text{trend}} = \sigma(w(\mathbf{x})) * f(\mathbf{x}), \mathbf{x}_{\text{seasonal}} = \mathbf{x} - \mathbf{x}_{\text{trend}}, \quad (5)$$

206 where $\sigma, w(x), f(x)$ denote the softmax operation, data-dependent weights and average filters.

207 For the trend component, we use a three-layer MLP to predict the future trend. As reversible
 208 instance normalization (RevIN) proves to be effective to remove and restore the non-stationary
 209 information [7, 23] which mainly resides in trend, we also add RevIN layers before and after MLP:
 210 $\mathcal{X}_{\text{trend}} = \text{RevIN}(\text{MLP}(\text{RevIN}(\mathbf{x}_{\text{trend}})))$. For the seasonal component, we adopt Transformer
 211 architecture but replace time-domain attention with Fourier-domain attention. More specifically, we
 212 first feed the seasonal part to N layers of encoder:

$$\mathcal{X}_{\text{en}}^{l,1} = \text{Norm}(\text{FA}(\mathcal{X}_{\text{en}}^{l-1}) + \mathcal{X}_{\text{en}}^{l-1}), \mathcal{X}_{\text{en}}^{l,2} = \text{Norm}(\text{FF}(\mathcal{X}_{\text{en}}^{l,1}) + \mathcal{X}_{\text{en}}^{l,1}), \mathcal{X}_{\text{en}}^l = \mathcal{X}_{\text{en}}^{l,2}, l = 1, \dots, N, \quad (6)$$

214 where $\mathcal{X}_{\text{en}}^0 = \mathbf{x}_{\text{seasonal}}$, FA and FF are short for Fourier Attention and Feed Forward network.
 215 Fourier Attention computes the attention in Fourier space and converts the output to time domain at
 216 the end (Definition 3.2) with $\sigma(\cdot) = \text{softmax}(\cdot)$:

$$\mathbf{o}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \mathcal{F}^{-1}\{\text{softmax}(\mathcal{F}\{\mathbf{q}\}\overline{\mathcal{F}\{\mathbf{k}\}}^T)\mathcal{F}\{\mathbf{v}\}\}. \quad (7)$$

218 The seasonal part is also zero-padded for the future part and fed into M layers of decoder to obtain
 219 the final seasonal output:

$$\mathcal{X}_{\text{de}}^{l,1} = \text{Norm}(\text{FA}(\mathcal{X}_{\text{de}}^{l-1}) + \mathcal{X}_{\text{de}}^{l-1}), \mathcal{X}_{\text{de}}^{l,2} = \text{Norm}(\text{FA}(\mathcal{X}_{\text{en}}^N, \mathcal{X}_{\text{de}}^{l,1}) + \mathcal{X}_{\text{de}}^{l,1}), \quad (8)$$

$$\mathcal{X}_{\text{de}}^{l,3} = \text{Norm}(\text{FF}(\mathcal{X}_{\text{de}}^{l,2}) + \mathcal{X}_{\text{de}}^{l,2}), \mathcal{X}_{\text{de}}^l = \mathcal{X}_{\text{de}}^{l,3}, l = 1, \dots, M, \quad (9)$$

220 where $\mathcal{X}_{\text{de}}^0 = \text{Padding}(\mathbf{x}_{\text{seasonal}})$. We add the trend prediction from MLP and seasonal prediction
 222 from Transformer to obtain the final output prediction, i.e., $\mathcal{X}_{\text{final}} = \mathcal{X}_{\text{trend}} + \mathcal{X}_{\text{de}}^M$. Optimization is
 223 based on a reconstruction MSE loss between predicted and ground truth future time series.

224 **Remark.** While FEDformer [24] and Autoformer [21] also have seasonal-trend decomposition,
 225 their trend and seasonal components are not disentangled; the trend prediction still comes from
 226 the attention module, which is sub-optimal based on our analysis in Section 4 and our empirical
 227 results in Section 6. By contrast, we apply seasonal-trend decomposition in the beginning, and
 228 apply Fourier attention only on seasonality components. This seemingly simple different way of
 229 decomposition brings significant performance gains to see in the experiment section, with even less
 230 model complexity. Non-stationary Transformer also computes attention for trend data. Moreover,
 231 with RevIN, TDformer has the similar effect of stationarization.

232

Table 3: MSE and MAE of different attention models with real-world seasonal and trend data.

Method	Metric	Traffic				Weather			
		96	192	336	720	96	192	336	720
Time	MSE	0.659	0.671	0.691	0.691	0.332	0.556	0.743	0.888
	MAE	0.358	0.358	0.368	0.363	0.395	0.533	0.622	0.702
Fourier	MSE	0.631	0.629	0.655	0.667	0.774	0.743	0.833	1.106
	MAE	0.338	0.336	0.345	0.350	0.648	0.632	0.659	0.769
Wavelet	MSE	0.622	0.629	0.640	0.655	0.358	0.564	0.815	1.312
	MAE	0.337	0.334	0.338	0.346	0.413	0.535	0.664	0.841

233 6 Experiments

234 6.1 Dataset and Baselines

235 We conduct experiments on benchmark time-series forecasting datasets: ETTm2 [22], electricity²,
 236 exchange [10], traffic³, weather⁴. We quantify the strength of seasonality for each dataset (details
 237 in Appendix). Electricity and traffic are strongly seasonal data, while exchange rate and weather
 238 demonstrate less seasonality and more trend. We compare TDformer with state-of-the-art attention
 239 models: Non-stationary Transformer [14], FEDformer [24], Autoformer [21], Informer [22], Log-
 240 Trans [12], Reformer [9]. As classical models (e.g., ARIMA), RNN-based models and CNN-based
 241 models generate large errors as shown in previous papers [22, 21], here we do not include their
 242 performance in the comparison. We use Adam [8] optimizer with a learning rate of $1e^{-4}$ and batch
 243 size of 32. We split the dataset with 7 : 2 : 1 into training, validation and test set, use validation
 244 set for hyperparameter tuning and report the results on the test set. For all real-world experiments,
 245 we feed the past 96 timesteps as context to predict the next 96, 192, 336, 720 timesteps following
 246 previous works [24, 21]. All experiments are repeated 5 times and we report the mean MSE and
 247 MAE. We implement in Pytorch on NVIDIA V100 16GB GPUs.

248 6.2 Comparing Attention Models on Real-World Datasets

249 As an extension to experiments on synthetic data (Section 4), we also compare attention models
 250 on real-world datasets, and observe consistent results as on synthetic datasets. Note that for a
 251 fair comparison, we directly compare the attention models without additional components like
 252 decomposition blocks or additional learnable transformation kernels [21, 24]. We choose traffic
 253 dataset as data with seasonality and weather dataset as data with trend. As shown in Table 3,
 254 frequency-domain attention models demonstrate better performance with seasonal data, which aligns
 255 with our observations on synthetic datasets. For trend data, Fourier-attention models show larger errors
 256 compared with time and wavelet attention models, which is also consistent with our observations
 257 on synthetic datasets. Compared with the reported performance after seasonal-trend decomposition
 258 as in FEDformer [24] and Autoformer [21], the errors on seasonal data remain similar, while errors
 259 increase significantly on trend data. This emphasizes the importance of seasonal-trend decomposition.

260 6.3 Main Results

261 We compare TDformer with the state-of-the-art baselines and report on MSE and MAE in Table 4.
 262 TDformer consistently demonstrates better performance across different datasets and forecasting
 263 horizons. On average, TDformer reduces the MSE by 9.14% compared with Non-stationary Trans-
 264 former and by 14.69% compared with FEDformer, and we attribute such improvement to our separate
 265 modeling of trend and seasonality with MLP and Fourier attention. As we mention in Remark 5,
 266 trend prediction of FEDformer and Non-stationary Transformer still come from attention modules,
 267 while TDformer decouples the modeling of trend and seasonality, and demonstrates better forecasting
 268 results. See Figure 1 for qualitative comparison.

269 6.4 Ablation Study

270 To separately understand the effect of trend and seasonal modules, we conducted ablation studies.
 271 TDformer-MLP-TA(WA) replaces Fourier attention with time (wavelet) attention for seasonality,

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

³<http://pems.dot.ca.gov>

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

Table 4: MSE and MAE of multivariate time-series forecasting on benchmark datasets with input context length 96 and forecasting horizon $\{96, 192, 336, 720\}$. We **bold** the best performing results.

Methods		TDformer		Non-stat TF		FEDformer		Autoformer		Informer		LogTrans		Reformer	
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Electricity	96	0.160	0.263	0.169	0.273	0.193	0.308	0.201	0.317	0.274	0.368	0.258	0.357	0.312	0.402
	192	0.172	0.275	0.182	0.286	0.201	0.315	0.222	0.334	0.296	0.386	0.266	0.368	0.348	0.433
	336	0.186	0.290	0.200	0.304	0.214	0.329	0.231	0.338	0.300	0.394	0.280	0.380	0.350	0.433
	720	0.215	0.313	0.222	0.32	0.246	0.355	0.254	0.361	0.373	0.439	0.283	0.376	0.340	0.420
Exchange	96	0.089	0.208	0.111	0.237	0.148	0.278	0.197	0.323	0.847	0.752	0.968	0.812	1.065	0.829
	192	0.183	0.305	0.219	0.335	0.271	0.380	0.300	0.369	1.204	0.895	1.040	0.851	1.188	0.906
	336	0.353	0.429	0.421	0.476	0.460	0.500	0.509	0.524	1.672	1.036	1.659	1.081	1.357	0.976
	720	0.932	0.725	1.092	0.769	1.195	0.841	1.447	0.941	2.478	1.310	1.941	1.127	1.510	1.016
Traffic	96	0.545	0.320	0.612	0.338	0.587	0.366	0.613	0.388	0.719	0.391	0.684	0.384	0.732	0.423
	192	0.571	0.329	0.613	0.340	0.604	0.373	0.616	0.382	0.696	0.379	0.685	0.390	0.733	0.420
	336	0.589	0.331	0.618	0.328	0.621	0.383	0.622	0.337	0.777	0.420	0.733	0.408	0.742	0.420
	720	0.606	0.337	0.653	0.355	0.626	0.382	0.660	0.408	0.864	0.472	0.717	0.396	0.755	0.423
Weather	96	0.177	0.215	0.173	0.223	0.217	0.296	0.266	0.336	0.300	0.384	0.458	0.490	0.689	0.596
	192	0.224	0.257	0.245	0.285	0.276	0.336	0.307	0.367	0.598	0.544	0.658	0.589	0.752	0.638
	336	0.278	0.290	0.321	0.338	0.339	0.359	0.380	0.395	0.578	0.523	0.797	0.652	0.639	0.596
	720	0.368	0.351	0.414	0.410	0.403	0.428	0.419	0.428	1.059	0.741	0.869	0.675	1.130	0.792
ETTm2	96	0.174	0.256	0.192	0.274	0.203	0.287	0.255	0.339	0.365	0.453	0.768	0.642	0.658	0.619
	192	0.243	0.302	0.280	0.339	0.269	0.328	0.281	0.340	0.533	0.563	0.989	0.757	1.078	0.827
	336	0.308	0.344	0.334	0.361	0.325	0.366	0.339	0.372	1.363	0.887	1.334	0.872	1.549	0.972
	720	0.400	0.400	0.417	0.413	0.421	0.415	0.422	0.419	3.379	1.338	3.048	1.328	2.631	1.242

Table 5: MSE and MAE of our model ablations. TDformer-MLP-TA replaces Fourier Attention by Time Attention (TA) for seasonality; TDformer-MLP-WA replaces Fourier Attention by Wavelet Attention (WA) for seasonality; TDformer-TA-FA replaces MLP with Time Attention (TA) for trend. TDformer w/o RevIN removes RevIN normalization.

Method	Metric	Traffic				Exchange			
		96	192	336	720	96	192	336	720
TDformer	MSE	0.545	0.571	0.589	0.606	0.089	0.183	0.353	0.932
	MAE	0.320	0.329	0.331	0.337	0.208	0.305	0.429	0.725
TDformer-MLP-TA	MSE	0.573	0.592	0.605	0.630	0.086	0.181	0.340	0.923
	MAE	0.334	0.336	0.340	0.351	0.205	0.303	0.422	0.721
TDformer-MLP-WA	MSE	0.552	0.583	0.599	0.629	0.088	0.185	0.348	0.925
	MAE	0.322	0.330	0.337	0.347	0.208	0.307	0.426	0.721
TDformer-TA-FA	MSE	0.590	0.590	0.617	0.642	0.242	0.349	0.629	0.908
	MAE	0.338	0.336	0.349	0.357	0.327	0.419	0.558	0.720
TDformer w/o RevIN	MSE	0.577	0.595	0.607	0.636	0.093	0.201	0.392	1.042
	MAE	0.320	0.325	0.328	0.339	0.222	0.330	0.474	0.763

272 and shows larger errors especially on seasonal data (traffic), as Fourier attention is more capable of
 273 capturing seasonality. Exchange data is mainly composed of trend, so different attention variants
 274 demonstrate similar performance. We also replace MLP with time attention for trend (TDformer-
 275 TA-FA) and observe large errors, as attention models show inferior generalization ability on trend
 276 data. TDformer w/o RevIN removes RevIN normalization and displays larger errors, which shows
 277 the importance of normalization for non-stationary data.

278 7 Conclusion

279 In this work we are driven by better understanding the relationships and separate benefits of attention
 280 models in time, Fourier and wavelet domains. We show that theoretically these three attention
 281 models are equivalent given linear assumptions. However, empirically due to the role of softmax,
 282 these models have respective benefits when applied to datasets with specific properties. Moreover,
 283 all attention models show inferior generalizability on data with trend. Based on these analyses
 284 of performance differences, we propose TDformer which separately models trend and seasonality
 285 with MLP and Fourier attention models after seasonal trend decomposition. TDformer achieves
 286 state-of-the-art performance against current attention models on time-series forecasting benchmarks.
 287 In the future, we plan to explore more complicated models to predict trend (e.g., autoregressive
 288 models) and explore other seasonal-trend decomposition methods.

289 **References**

- 290 [1] Cristian Challu, Kin G Olivares, Boris N Oreshkin, Federico Garza, Max Mergenthaler, and
291 Artur Dubrawski. N-hits: Neural hierarchical interpolation for time series forecasting. *arXiv*
292 *preprint arXiv:2201.12886*, 2022.
- 293 [2] Robert B Cleveland, William S Cleveland, Jean E McRae, and Irma Terpenning. Stl: A
294 seasonal-trend decomposition. *J. Off. Stat*, 6(1):3–73, 1990.
- 295 [3] John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan
296 Catanzaro. Adaptive fourier neural operators: Efficient token mixers for transformers. *arXiv*
297 *preprint arXiv:2111.13587*, 2021.
- 298 [4] Gaurav Gupta, Xiongye Xiao, and Paul Bogdan. Multiwavelet-based operator learning for
299 differential equations. *Advances in Neural Information Processing Systems*, 34:24048–24062,
300 2021.
- 301 [5] Hao Huang and Yi Fang. Adaptive wavelet transformer network for 3d shape representation
302 learning. In *International Conference on Learning Representations*, 2021.
- 303 [6] Song Jiang, Tahin Syed, Xuan Zhu, Joshua Levy, Boris Aronchik, and Yizhou Sun. Bridging
304 self-attention and time series decomposition for periodic forecasting. 2022.
- 305 [7] Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo.
306 Reversible instance normalization for accurate time-series forecasting against distribution shift.
307 In *International Conference on Learning Representations*, 2021.
- 308 [8] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
309 *arXiv:1412.6980*, 2014.
- 310 [9] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer.
311 *arXiv preprint arXiv:2001.04451*, 2020.
- 312 [10] Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term
313 temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference*
314 *on research & development in information retrieval*, pages 95–104, 2018.
- 315 [11] James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens
316 with fourier transforms. *arXiv preprint arXiv:2105.03824*, 2021.
- 317 [12] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyu Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng
318 Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series
319 forecasting. *Advances in neural information processing systems*, 32, 2019.
- 320 [13] Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya,
321 Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential
322 equations. *arXiv preprint arXiv:2010.08895*, 2020.
- 323 [14] Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers:
324 Rethinking the stationarity in time series forecasting. *arXiv preprint arXiv:2205.14415*, 2022.
- 325 [15] LIU Minhao, Ailing Zeng, LAI Qiuxia, Ruiyuan Gao, Min Li, Jing Qin, and Qiang Xu.
326 T-wavenet: A tree-structured wavelet neural network for time series signal analysis. In *International*
327 *Conference on Learning Representations*, 2021.
- 328 [16] Boris N Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis
329 expansion analysis for interpretable time series forecasting. *arXiv preprint arXiv:1905.10437*,
330 2019.
- 331 [17] Fan-Keng Sun and Duane S Boning. Fredo: Frequency domain-based long-term time series
332 forecasting. *arXiv preprint arXiv:2205.12301*, 2022.
- 333 [18] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
334 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
335 *processing systems*, 30, 2017.

- 336 [19] Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun.
337 Transformers in time series: A survey. *arXiv preprint arXiv:2202.07125*, 2022.
- 338 [20] Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. Etsformer: Exponential smoothing transformers for time-series forecasting. *arXiv preprint arXiv:2202.01381*,
339 2022.
340
- 341 [21] 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.
342
343
- 344 [22] 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.
345
346
347
- 348 [23] Tian Zhou, Ziqing Ma, Qingsong Wen, Liang Sun, Tao Yao, Rong Jin, et al. Film: Frequency improved legendre memory model for long-term time series forecasting. *arXiv preprint arXiv:2205.08897*, 2022.
349
350
- 351 [24] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. *arXiv preprint arXiv:2201.12740*, 2022.
352
353

354 Checklist

- 355 1. For all authors...
- 356 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
357 contributions and scope? [\[Yes\]](#)
- 358 (b) Did you describe the limitations of your work? [\[Yes\]](#)
- 359 (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#) We focus
360 on promoting understanding of existing models, and propose method that improves on
361 forecasting benchmarks.
- 362 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
363 them? [\[Yes\]](#)
- 364 2. If you are including theoretical results...
- 365 (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Section 3
- 366 (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See Section 3
- 367 3. If you ran experiments...
- 368 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
369 mental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See Section 6
- 370 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
371 were chosen)? [\[Yes\]](#) See Section 6
- 372 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
373 ments multiple times)? [\[Yes\]](#)
- 374 (d) Did you include the total amount of compute and the type of resources used (e.g., type
375 of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Section 6
- 376 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 377 (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) See Section 6
- 378 (b) Did you mention the license of the assets? [\[Yes\]](#) See Section 6
- 379 (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
380 We didn’t introduce new datasets.
- 381 (d) Did you discuss whether and how consent was obtained from people whose data you’re
382 using/curating? [\[N/A\]](#) We use public benchmark forecasting datasets.
- 383 (e) Did you discuss whether the data you are using/curating contains personally identifiable
384 information or offensive content? [\[N/A\]](#) We use public benchmark forecasting datasets.

- 385 5. If you used crowdsourcing or conducted research with human subjects...
- 386 (a) Did you include the full text of instructions given to participants and screenshots, if
 387 applicable? [N/A] We didn't use crowdsourcing or involve human subjects in this
 388 paper.
- 389 (b) Did you describe any potential participant risks, with links to Institutional Review
 390 Board (IRB) approvals, if applicable? [N/A] We didn't use crowdsourcing or involve
 391 human subjects in this paper.
- 392 (c) Did you include the estimated hourly wage paid to participants and the total amount
 393 spent on participant compensation? [N/A] We didn't use crowdsourcing or involve
 394 human subjects in this paper.

395 A Appendix

396 We first apply STL decomposition [2] for each dataset

$$X_t = T_t + S_t + R_t, \quad (10)$$

397 where T_t, S_t, R_t respectively represent the trend, seasonal and remainder component. For data with
 398 strong seasonality, the seasonal component would have much larger variation than the remainder
 399 component; while for data with little seasonality, the two variances should be similar. Therefore, we
 400 can quantify the strength of seasonality as

$$S = \max(0, 1 - \frac{\text{Var}(R_t)}{\text{Var}(S_t) + \text{Var}(R_t)}). \quad (11)$$

401 Following this equation, we summarize the seasonality strength of each dataset in Table 6. Electricity
 402 and traffic are strongly seasonal data, while exchange rate and weather demonstrate less seasonality
 403 and more trend.

Table 6: Seasonality strength of benchmark datasets.

Dataset	Electricity	Exchange	Traffic	Weather	ETTM2
Seasonality Strength	0.982	0.284	0.967	0.317	0.711