
FReDF: LEARNING TO FORECAST IN FREQUENCY DOMAIN

A PREPRINT

Hao Wang¹, Licheng Pan¹, Zhichao Chen¹, Degui Yang², Sen Zhang³, Yifei Yang⁴, Xinggao Liu¹, Haoxuan Li^{5†},
Dacheng Tao^{3†}

¹Zhejiang University

²Central South University

³Nanyang Technological University

⁴Shanghai Jiaotong University

⁵Peking University

†Corresponding author: lihx@stu.pku.edu.cn, dacheng.tao@gmail.com

February 6, 2024

ABSTRACT

Time series modeling is uniquely challenged by the presence of autocorrelation in both historical and label sequences. Current research predominantly focuses on handling autocorrelation within the historical sequence but often neglects its presence in the label sequence. Specifically, emerging forecast models mainly conform to the direct forecast (DF) paradigm, generating multi-step forecasts under the assumption of conditional independence within the label sequence. This assumption disregards the inherent autocorrelation in the label sequence, thereby limiting the performance of DF-based models. In response to this gap, we introduce the Frequency-enhanced Direct Forecast (FreDF), which bypasses the complexity of label autocorrelation by learning to forecast in the frequency domain. Our experiments demonstrate that FreDF substantially outperforms existing state-of-the-art methods and is compatible with a variety of forecast models. The scripts and log files will be released at <https://github.com/Master-PLC/FreDF>¹.

1 Introduction

Time series modeling aims to encode historical sequence to predict future data, which is crucial in diverse applications: long-term forecast in weather prediction [3, 40], short-term prediction in industrial maintenance [24, 7, 35], and missing data imputation in healthcare [30]. A key challenge in time series modeling, distinguishing it from canonical regression tasks, is the presence of autocorrelation. It refers to the dependence between time steps, which exists in *both* the input and label sequences.

To accommodate autocorrelation in input sequences, diverse forecast models have been developed [28, 5, 8], exemplified by recurrent [29], convolution [37] and graph neural networks [25, 4, 11]. Recently, Transformer-based models, utilizing self-attention mechanisms to dynamically assess autocorrelation, have gained prominence in this line of work [20, 26, 13, 38]. Concurrently, there is a growing trend of incorporating frequency analysis into forecast models [41, 21]. By representing input sequence in the frequency domain, the complexity of portraying autocorrelation is bypassed, which improves forecast performance of Transformers [47, 23], GNNs [42] and MLPs [43]. These pioneering works highlight the importance of autocorrelation and frequency analysis in advanced time series modeling.

Another critical aspect is the autocorrelation in the label sequence, where each step in the label sequence is autoregressively generated. This phenomenon, known as *label autocorrelation*, is often neglected in current forecast techniques. Specifically, modern methods predominantly adopt the direct forecast (DF) paradigm [20, 27], which generates multi-step forecasts simultaneously using a multi-output head [18, 16] and seeks to minimize the forecast errors across all

¹If you have any queries, feel free to contact us through haohaow@zju.edu.cn and 22132045@zju.edu.cn.

steps concurrently. This approach implicitly assumes the *step-wise independence* within the label sequence given input sequence, thereby ignoring the existence of label autocorrelation. Such discrepancy between model assumption and data characteristic results in suboptimal forecast quality, underscoring a huge limitation in existing direct forecast paradigm.

To handle this limitation, we propose Frequency-enhanced Direct Forecast (FreDF), a simple yet effective approach that refines the DF paradigm by aligning the forecasts and label sequence in the frequency domain. By transforming to the frequency domain where bases are orthogonal and independent, the impact of autocorrelation is found to be effectively diminished. Therefore, FreDF bypasses the discrepancy between the assumption of DF and the existence of label autocorrelation, while retaining DF’s benefits such as sample efficiency and implementation simplicity. The main contributions in this work can be summarized below.

- We recognize and formulate the impact of label autocorrelation ignored by current DF paradigm in forecasting.
- We propose FreDF for time series forecasting. As an embarrassingly simple update to DF, it handles label autocorrelation by learning to forecast in the frequency domain. To our knowledge, it is the first attempt to employ frequency analysis for enhancing forecast paradigms.
- We verify the efficacy of FreDF through extensive experiments, where it outperforms state-of-the-art methods substantially and supports various forecast models.

2 Preliminaries and Related Work

2.1 Problem Definition

In this study, uppercase letters (e.g., Y) denote random matrix, with subscripts (e.g., $Y_{i,j}$) indicating matrix entries. An uppercase letter followed by parentheses (e.g., $Y(n)$) represents a sampled observation of the random matrix.

A multi-variate time series can be represented as a sequence $[X(1), X(2), \dots, X(N)]$, where $X(n) \in \mathbb{R}^{1 \times D}$ is the sample at the n -th timestamp with D covariates [20, 37]. Define an input sequence $L \in \mathbb{R}^{L \times D}$ and a label sequence $Y \in \mathbb{R}^{T \times D}$ where L and T are sequence lengths. At the n -th step, these sequences are observed as $L(n) = [X(n-L+1), \dots, X(n)]$ and $Y(n) = [X(n+1), \dots, X(n+T)]$. The goal of time series forecast is identifying a model $g: \mathbb{R}^{L \times D} \rightarrow \mathbb{R}^{T \times D}$ within a model family \mathcal{G} (e.g., decision trees, neural networks) that generates forecasts $\hat{Y} = g(L)$ approximating the label sequence Y .

There are two critical aspects to accommodate autocorrelation in this task: (1) selecting a model family \mathcal{G} that encodes autocorrelation in input sequences, which underscores the design of model architectures; (2) generating forecasts that respect label autocorrelation, which highlights the efficacy of forecast paradigms. Our survey concentrates on examining both aspects for accommodating autocorrelation.

2.2 Model Architecture

To exploit autocorrelation in the input sequences, diverse architectures have been developed. Initial statistical methods include VAR [36] and ARIMA [1]. Subsequently, neural networks became increasingly prominent for their ability to automate feature interaction and capture nonlinear correlations. Exemplars include RNNs (e.g., DeepAR [29], S4 [10]), CNNs (e.g., TimesNet [37], SCINet [17]), and GNNs (e.g., MTGNN [25]), each designed to effectively encode autocorrelation. Current progress reaches a debate between Transformer-based and MLP-based architectures, each with its advantages and limitations [19, 22, 6]. Transformers (e.g., PatchTST [27], iTransformer [20], CrossFormer [45]) excel in encoding autocorrelation but come with high computational costs, while MLPs (e.g., DLinear [44], TSMixer [9, 6]) are more efficient but less adept at autocorrelation encoding.

An emerging approach is representing sequence in the frequency domain. This method, in comparison to modeling autocorrelation in the temporal domain, manages autocorrelation effectively with limited cost. A prominent example is FedFormer [47], which computes attention scores in the frequency domain, leading to improved efficiency, efficacy, and noise reduction capabilities. The success of this technique extends to various architectures like Transformers [47, 39], MLPs [43] and GNNs [42, 4], which makes it a versatile plugin in the design of neural networks for time series forecast.

2.3 Forecast Paradigm

There are two paradigms to generate multi-step forecast: iterative forecast (IF) and direct forecast (DF) [18]. The IF paradigm forecasts one step at a time, using previous predictions as input for subsequent forecasts. This recursive approach respects label autocorrelation in forecast generation, widely used by early-stage methods [14, 29]. However, IF is prone to high variance due to error propagation, which significantly impairs performance in long-term forecasts [31].

Handling the error propagation problem of IF, DF generates multi-step forecasts simultaneously using a multi-output head, featured by fast inference, implementation ease and superior accuracy. As a result, starting with Informer [16], DF has been dominant for multi-step forecast, continuing to be employed in recent works such as TimesNet [37], PatchTST [27] and iTransformer [20].

Significance of this work. Our work augments the DF paradigm by enabling forecast in the frequency domain. Unlike recent advancements [43, 47, 42] that incorporate frequency analysis to refine model architectures for managing input autocorrelation, our work focuses on improving the *forecast paradigm for managing label autocorrelation*, which is an unexplored and innovative aspect of time series modeling.

3 Direct Forecast meets Autocorrelation

3.1 Oversight of Label Autocorrelation

In this section, we delineate the DF paradigm and its neglect of label autocorrelation. DF employs a multi-output model $g_\theta : \mathbb{R}^{L \times D} \rightarrow \mathbb{R}^{T \times D}$ for generating T-step forecasts $\hat{Y} = g(L)$. Let $Y_t \in \mathbb{R}^{1 \times D}$ be the t -th step of Y , $Y_t(n)$ be the n -th sampled observation of it; the model parameters θ are optimized by minimizing the mean squared error (MSE):

$$\mathcal{L}^{(\text{tmp})} := \sum_{n=1}^N \|Y(n) - g_\theta(L(n))\|_2^2, = \sum_{n,t=1}^{N,T} \|Y_t(n) - \hat{Y}_t(n)\|_2^2. \quad (1)$$

The DF paradigm computes the forecast error at each step independently, treating them as separate tasks. This method, while practical, overlooks the autocorrelation present within Y . We provide a theoretical rationale for this limitation in Theorem 3.1, framed from a maximum likelihood estimation perspective. The theorem posits that the likelihood of θ can be accurately depicted by (1) only if the assumption of conditional independence holds: different steps in Y are mutually independent given L . This assumption contradicts the presence of label autocorrelation, which results in a biased likelihood and a deviation from the maximum likelihood principle during model training.

Theorem 3.1. *Given input sequence L and label sequence Y , the negative log-likelihood of the model g parameterized by θ can be expressed as (1) only if Y_t is independent on $Y_{t'}$ for any $1 \leq t, t' \leq T$ given L , expressed as $Y_t \perp\!\!\!\perp Y_{t'} \mid L$. The proof is provided in Appendix B.*

This issue can be visualized in Figure 1 (a). Specifically, label sequence is autoregressively generated, with the value of Y_{t+1} being highly dependent on Y_t as indicated by the blue arrows. In contrast, the learning objective (1) presumes conditional independence as indicated by the black arrows, neglecting the label autocorrelation as indicated by the blue arrows. Such discrepancy between model assumption and data characteristic limits forecast performance.

3.2 Empirical Evidence

In this section, we empirically verify the presence of label autocorrelation in time series data. As depicted in Figure 1 (a), verifying label autocorrelation entails quantifying the causal relationship $Y_t \rightarrow Y_{t+1}$. However, this quantification is complex due to the confounding effect of L , which induces pseudo-correlation. This pseudo-correlation obscures the actual causal relationship, rendering correlation measures like Pearson correlation ineffective².

To handle the confounding effect, we utilize tools from causal inference. Specifically, we treat L as the confounder to adjust, Y_t as the treatment, and Y_{t+1} as the outcome. Double machine learning (DML), a reliable method in causal inference, is employed for accurately quantifying the causation $Y_t \rightarrow Y_{t+1}$ while eliminating the confounding effect. For verification efficiency and without loss of generality, we focus on the last feature in Y_t , as the analysis does not

²The fork structure $Y_t \leftarrow L \rightarrow Y_{t+1}$ produces psudeo-correlations between Y_t and Y_{t+1} which confounds the true autocorrelation $Y_t \rightarrow Y_{t+1}$. This is a known issue namely confounding effect that has been intensively investigated in causal inference [34, 33, 15].

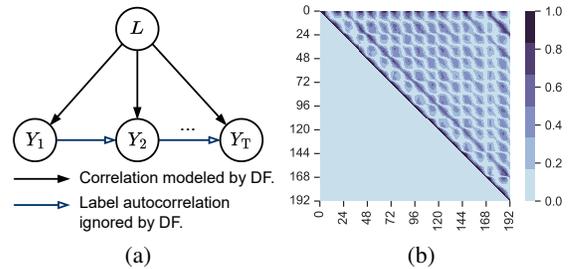


Figure 1: Visualization of autocorrelation in the forecast window Y . (a) Data generation process of time-series with dependencies depicted as arrows. (b) The autocorrelation identified with DML, where each element at the i -th row and j -th column indicates the absolute causation strength of $Y_i \rightarrow Y_j$. The lower part is masked to preserve causation.

depend on feature dimension. Experiments are conducted on the Weather dataset with $T = 192$. The results are presented in Figure 1 (b) with key observations below.

- Diagonal elements consistently show values of 1, which is expected as the treatment and outcome are identical. The outcome’s change mirrors the treatment’s change.
- Non-diagonal elements demonstrate significant values, with nearly 37.5% exceeding 0.3. It means that different steps in Y are interdependent given L , affirming the presence of label autocorrelation. Additionally, the autocorrelation strength displays a regular variation pattern, evidenced by alternating light and dark areas in the figure, which likely suggests a periodic nature in the series³.

In summary, we have verified the existence of autocorrelation in the label sequence, which contradicts the independence assumption of the DF paradigm in Theorem 3.1.

4 Proposed Method

4.1 Bypass Autocorrelation with Domain Transform

As established in Section 3, the canonical DF paradigm suffers from suboptimal performance due to its oversight of label autocorrelation. A promising strategy to overcome this limitation involves representing the label sequence in a transformed domain formed with orthogonal bases, denoted as $F = \mathcal{F}(Y)$. Specifically, it can be effectively implemented using the Fourier transform in Definition 2, which projects the sequence onto orthogonal bases associated with different frequencies. By transforming the label sequence into this orthogonal frequency domain, the dependence from label autocorrelation could be effectively mitigated.

Definition 4.1. *The Fourier transform of a sequence $\mathbf{x} = [x_0, \dots, x_{T-1}]$ is defined as its projection onto a set of orthogonal Fourier bases with different frequencies. The projection associated with frequency k is computed as*

$$x_k^{(F)} = \sum_{t=0}^{T-1} x_t \exp(-j(2\pi/T)kt), \quad 0 \leq k \leq T-1, \quad (2)$$

where i is the imaginary unit which is defined as the square root of -1 , $\exp(\cdot)$ is the Fourier basis associated with the frequency k which is orthogonal for different k values. Fourier transform refers to the projections associated with frequencies $0 \leq k \leq T-1$, denoted as $\mathbf{x}^{(F)} = \mathcal{F}(\mathbf{x})$, which can be computed via the fast Fourier transform (FFT) algorithm with complexity $\mathcal{O}(L \log L)$.

To substantiate this claim, we employ DML to assess the dependencies between frequency components in the transformed representation F . In this context, L is treated as the confounder, F_k (the component of F corresponding to frequency k) as the treatment, and $F_{k'}$ as the outcome.

According to Figure 2, most non-diagonal elements show negligible values, with merely about 3.6% exceeding 0.1, suggesting that frequency components of F are almost independent given L ³. Such independence implies that representing label sequence in the frequency domain bypasses the dependency raised by autocorrelation in the time domain, aligning with DF’s independence assumption in Theorem 1. This observed alignment underscores the potential of learning to forecast in the frequency domain below.

4.2 Model Implementation

In this section, we construct FreDF, a simple yet effective update to the current DF training paradigm. The core technique contribution is aligning the forecasts generated by DF and the label sequences in the frequency domain.

The workflow is depicted in Figure 3. At a given time-stamp n , the historical sequence $L(n)$ is input into the model to generate T -step forecasts, denoted as $\hat{Y}(n) = g(L(n))$. The forecast error in the time domain $\mathcal{L}^{(\text{tmp})}$ is calculated

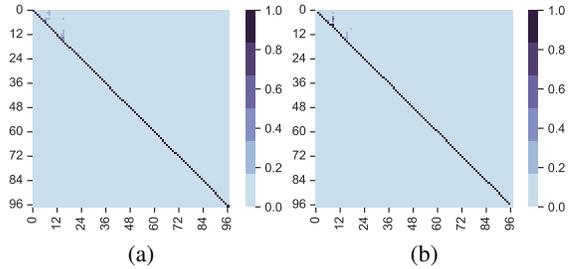


Figure 2: Dependencies between frequency components, measured by the real part (a) and imaginary part (b). Due to the symmetry of Fourier transform, only the half of F is visualized for clarity. The element at the i -th row and j -th column indicates the absolute causation strength $F_i \rightarrow F_j$.

³More implementation details, empirical evidence and formal analysis are provided in Appendix A.

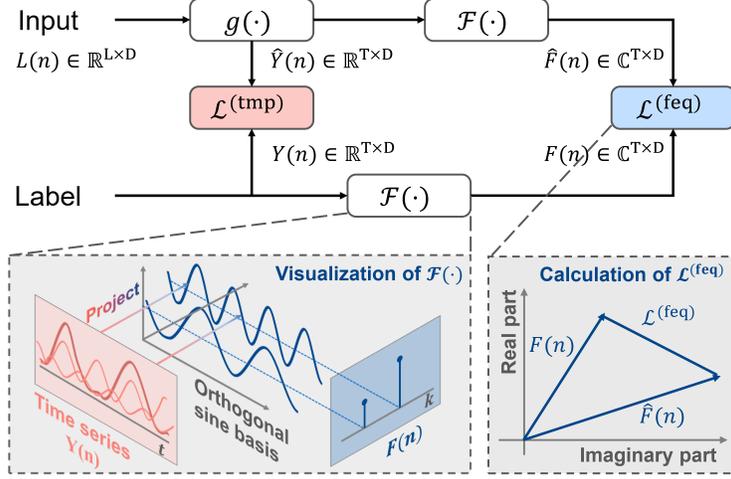


Figure 3: The main components of the Fourier-enhanced DF approach. Key operations in the time and frequency domains are highlighted in red and blue, respectively. The output of model g is firstly used to compute the temporal loss $\mathcal{L}^{(\text{tmp})}$, and subsequently conducted FFT to calculate the frequency loss $\mathcal{L}^{(\text{feq})}$.

Table 1: Long-term forecasting performance averaged over forecast lengths. Full results are present in Table 7.

Models	FreDF (Ours)		iTransformer (2024)		FreTS (2023)		TimesNet (2023)		Crossformer (2023)		TiDE (2023)		DLinear (2023)		FEDformer (2022)		Autoformer (2021)		Transformer (2017)		TCN (2017)		LSTM (1998)	
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	0.392	0.399	0.415	0.416	0.407	0.415	0.413	0.418	0.558	0.532	0.419	0.419	0.404	0.407	0.440	0.451	0.596	0.517	0.943	0.733	0.891	0.632	0.656	0.567
ETTm2	0.278	0.319	0.294	0.335	0.335	0.379	0.297	0.332	1.633	0.782	0.358	0.404	0.344	0.396	0.302	0.348	0.326	0.366	1.322	0.814	3.411	1.432	1.757	0.979
ETTh1	0.437	0.435	0.449	0.447	0.488	0.474	0.478	0.466	0.628	0.574	0.541	0.507	0.462	0.458	0.441	0.457	0.476	0.477	0.993	0.788	0.763	0.636	0.763	0.636
ETTh2	0.371	0.396	0.390	0.410	0.550	0.515	0.413	0.426	2.136	1.130	0.611	0.550	0.558	0.516	0.430	0.447	0.478	0.483	3.296	1.419	3.325	1.445	1.817	1.029
ECL	0.170	0.259	0.176	0.267	0.209	0.297	0.214	0.307	0.182	0.279	0.251	0.344	0.225	0.319	0.229	0.339	0.228	0.339	0.274	0.367	0.617	0.598	0.329	0.406
Traffic	0.421	0.279	0.428	0.286	0.552	0.348	0.636	0.335	0.553	0.292	0.760	0.473	0.673	0.419	0.611	0.379	0.637	0.399	0.680	0.376	1.001	0.652	0.890	0.487
Weather	0.254	0.274	0.281	0.302	0.255	0.299	0.262	0.288	0.262	0.324	0.271	0.320	0.265	0.317	0.311	0.361	0.349	0.391	0.632	0.552	0.584	0.572	0.306	0.357

according to (1). Subsequently, FreDF transform both the forecasts and the label sequences into the frequency domain, denoted as $\hat{F} = \mathcal{F}(\hat{Y})$ and $F = \mathcal{F}(Y)$. The forecast error in the frequency domain is computed below

$$\mathcal{L}^{(\text{feq})} := \sum_{n=1}^N \left| F(n) - \hat{F}(n) \right|, \quad (3)$$

where each term in the summation is a matrix of complex numbers; for a matrix $A \in \mathbb{C}^{N \times N}$, $|A|$ denotes the operation of computing and summing the modulus of each element in the matrix, with the modulus of a complex number $a = a_r + ia_i$ calculated as $\sqrt{a_r^2 + a_i^2}$. Notably, we do not use the squared loss form, as is typical in (1), due to the distinct numerical characteristics of the label sequence in the frequency domain. Specifically, different frequency components often exhibit vastly varying magnitudes, with lower frequencies having higher volumes by several orders of magnitude compared to higher frequencies, which renders squared loss methods unstable.

Finally, the forecast error in the time and frequency domains are fused as follow, where $0 \leq \alpha \leq 1$ controls the relatively strength of frequency-domain alignment:

$$\mathcal{L}^\alpha := \alpha \cdot \mathcal{L}^{(\text{feq})} + (1 - \alpha) \cdot \mathcal{L}^{(\text{tmp})}. \quad (4)$$

By aligning generated forecasts and label sequence in the frequency domain, FreDF bypasses the autocorrelation effect while preserving the benefits of DF such as efficient inference and multi-task capabilities. A notable property of FreDF

Table 2: Short-term forecast performance averaged over forecast lengths. The full results are present in Table 8.

Models	FreDF (Ours)	iTransformer (2024)	FreTS (2023)	Crossformer (2023)	DLinear (2023)	Fedformer (2022)
SMAPE	12.112	12.298	<u>12.169</u>	71.332	12.480	12.734
MASE	1.648	1.680	<u>1.660</u>	16.626	1.674	1.702
OWA	0.877	0.893	<u>0.883</u>	6.977	0.898	0.914

is its model and transformation agnosticism. It is compatible with various forecast models g (e.g., Transformers and MLPs) and transformations \mathcal{F} (e.g., Chebyshev and Legendre transforms). This flexibility significantly broadens the potential application scope of FreDF.

5 Experiments

5.1 Setup

Datasets. The datasets for long-term forecast and imputation include ETT (4 subsets), ECL, Traffic and Weather following [39]. The dataset for short-term forecast is M4 following [37]. Each dataset is divided chronologically for training, validation and test. Detailed dataset descriptions are provided in Appendix D.1.

Baselines. Our baselines include various established models in the time series field, which can be grouped into three categories: (1) Transformer-based methods: Transformer [32], Autoformer [39], FEDformer [47], Crossformer [45], iTransformer [20]; (2) MLP-based methods: DLinear [44], TiDE [8], FreTS [43]; (3) other notable models: LSTM [46], TimesNet [37], TCN [2]. Notably, iTransformer [20] is the state-of-the-art baseline released in ICLR-24.

Implementation. The baseline models are reproduced using the scripts sourced from TimesNet [37]. They are trained with Adam [12] optimizer to minimize the MSE loss. When integrating FreDF to enhance an established model, we respect the original hyperparameter settings, merely tuning α and learning rate. Roughly adjusting the learning rate is essential since the magnitude of MSE loss diverges by multiple orders between the time and frequency domains. More implementation details are provided in Appendix D.

5.2 Overall Performance

5.2.1 Long-term Forecast

The performance on the long-term forecast task is present in Table 1, where we select iTransformer as the forecast model g and enhance it with FreDF paradigm. The results are averaged over forecasting lengths $T=96, 192, 336$ and 720 , with the best results **bolded** and the second best results underlined. The main observations are summarized below.

- FreDF improves the performance of iTransformer substantially. For instance, on the ETTm1 dataset, FreDF decreases the MSE of iTransformer by 0.019. This improvement is comparable to the advancement observed in the dataset over 1.5 years, from Fedformer in 2022 to TimesNet in 2023, with a MSE reduction of 0.017. Similar gains are evident in other datasets, which can be attributed to reconciliation of label autocorrelation with the DF paradigm, validating efficacy of FreDF.
- FreDF achieves leading performance compared to a range of competitive baseline models. Notably, FreDF enhances the performance of iTransformer to surpass even those models that originally outperformed iTransformer on some datasets. It indicates that the improvements by FreDF exceed those achievable through dedicated architectural design alone, emphasizing the importance of label autocorrelation management and FreDF.

5.2.2 Short-term Forecast

In this section, we extend our scope to short-term forecast task. The results are summarized in Table 2, where we employ the FreDF paradigm to enhance FreTS which is identified as the best baseline in this task. Empirically, FreDF retains efficacious in this task, improving FreTS across three key metrics. The volume of improvement is commendable since the performance difference among competing models is typically slight in this task.

Table 3: Imputation performance averaged over missing ratios. The full results are present in Table 9.

Models	FreDF (Ours)		iTransformer (2024)		TiDE (2023)		Crossformer (2023)		Fedformer (2023)	
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	0.002	0.032	<u>0.003</u>	<u>0.041</u>	0.014	0.083	0.011	0.078	0.069	0.179
ETTm2	0.003	0.035	<u>0.003</u>	<u>0.036</u>	0.024	0.103	0.033	0.133	0.187	0.301
ETTh1	0.002	0.029	<u>0.003</u>	<u>0.037</u>	0.003	0.038	0.013	0.080	0.139	0.276
ETTh2	<u>0.003</u>	<u>0.033</u>	0.004	0.043	0.005	0.047	0.002	0.026	0.356	0.417
ECL	0.001	0.019	<u>0.002</u>	<u>0.033</u>	0.003	0.037	0.010	0.082	0.169	0.298
Weather	0.001	0.013	<u>0.001</u>	0.015	0.002	<u>0.013</u>	0.007	0.061	0.057	0.153

Table 4: Results of system-level ablation study averaged over forecast lengths. Full results are present in Table 10.

$\mathcal{L}^{(tmp)}$	$\mathcal{L}^{(freq)}$	Weather		ETTm1		ETTh1	
		MSE	MAE	MSE	MAE	MSE	MAE
✓	✗	0.2810	0.3021	0.4146	0.4156	0.4491	0.4467
✗	✓	<u>0.2573</u>	<u>0.2766</u>	<u>0.3929</u>	<u>0.3996</u>	<u>0.4379</u>	<u>0.4353</u>
✓	✓	0.2538	0.2739	0.3920	0.3989	0.4374	0.4351

5.2.3 Missing Data Imputation

In this section, we investigate missing data imputation task. iTransformer, identified as the best baseline for imputation tasks, is selected as the testbed for FreDF. All models are trained in an autoencoding manner: given input sequences with missing entries, the models are tasked with reconstructing the non-missing entries in the training phase, and employed to impute the missing entries in the inference phase.

The results in Table 3 demonstrate the efficacy of FreDF in this task: it improves the performance of iTransformer significantly, outperforming other competitive methods. A unique aspect of this task is that the label sequences are irregularly sampled due to missing entries, which disrupts the physical semantics associated with the Fourier transform. This implies that the principal strength of FreDF lies beyond the semantics of Fourier transform. Instead, its efficacy is rooted in its capability to align the data property and the model assumption underlying DF paradigm.

5.2.4 Showcases

In this section, we visualize the forecast sequences to highlight the improvements of FreDF in forecast quality. A ETTm2 snapshot with T=336 is depicted in Figure 4. While the model without FreDF can follow the general trends of the label sequence, it struggles to capture the sequence’s high-frequency components, resulting in a forecast with a visibly lower frequency. Additionally, the forecast sequence exhibits numerous burrs. These issues reflect the limitations of forecasting in the time domain, namely the difficulty in capturing high-frequency components and the neglect of autocorrelation between steps.

FreDF addresses these limitations effectively. The forecasts generated under FreDF not only keep pace with the label sequence, accurately capturing high-frequency components, but also exhibit a smoother appearance with fewer irregularities, due to its awareness of autocorrelation.

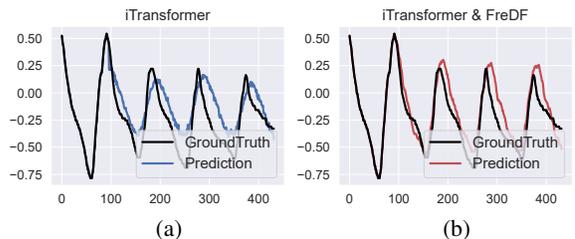
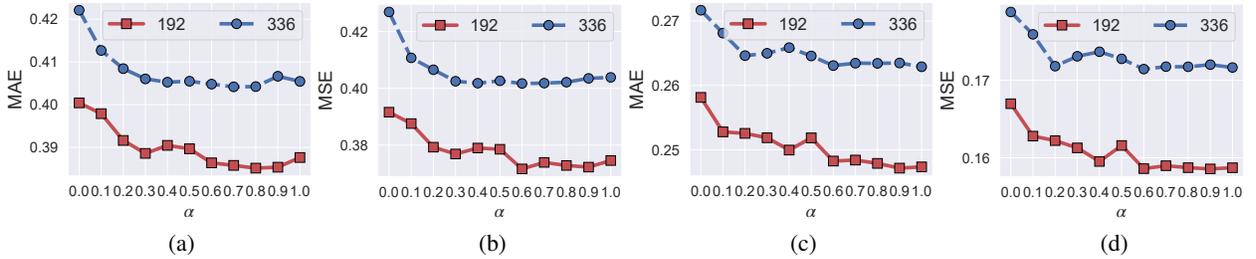


Figure 4: Forecast sequences with and without FreDF.

Table 5: Results of component-level ablation study averaged over forecast lengths.

Amp.	Pha.	ECL		ETTm1		ETTh1	
		MSE	MAE	MSE	MAE	MSE	MAE
✓	✗	0.3356	0.4060	0.5936	0.5169	0.7303	0.5968
✗	✓	0.1836	0.2752	0.4204	0.4173	0.4751	0.4487
✓	✓	0.1698	0.2594	0.3920	0.3989	0.4374	0.4351

Figure 5: Performance given varying strength of frequency loss α . Experiments are conducted on the ETTm1 (a-b) and ECL (c-d) datasets with forecast lengths being 192 and 336.

5.3 Ablation Studies

Temporal loss v.s. frequency loss. In this section, we focus on dissecting the contributions of the temporal ($\mathcal{L}^{(tmp)}$) and frequency domain ($\mathcal{L}^{(freq)}$) loss components. Utilizing the iTransformer as the forecast model, the ablation study’s results are detailed in Table 4. The main observations are summarized below.

- Transforming forecasts to the frequency domain yields consistent improvements. It is evidenced by the huge performance gains across all datasets when replacing $\mathcal{L}^{(tmp)}$ with $\mathcal{L}^{(freq)}$. The underlying rationale is that label autocorrelation can be effectively managed in the frequency domain, aligning better with the conditional independence assumption inherent in DF.
- Learning to forecast in both domains generally showcase improvement compared to relying solely on one domain. However, the improvement over $\mathcal{L}^{(freq)}$ is marginal. Hence, exclusively focusing on frequency domain forecasting emerges as a viable strategy in most cases, offering promising performance without the complexity of balancing learning objectives.

Amplitude v.s. Phase Characteristics. In this analysis, we explore the impact of amplitude and phase alignment on FreDF. Minimizing the frequency loss (3) ensures alignment of both amplitude and phase characteristics between the forecast and actual label sequences in the frequency domain. In the field of signal processing, both amplitude and phases are foundational for accurately representing the dynamics of signals.

Our findings from Table 5 indicate that both amplitude and phase characteristics are essential for FreDF. Notably, phase alignment emerges as particularly crucial. Solely aligning amplitude characteristics without phase alignment results in poor performance. This outcome is reasonable, since minor deviations in phase characteristics could correspond significant discrepancies in the time domain.

5.4 Hyperparameter Sensitivity

In this section, we vary the frequency loss strength α on the efficacy of FreDF. The results are summarized in Figure 5. Overall, increasing α from 0 to 1 results in a reduction of forecast error, albeit with a slight increase towards the end of this range. For instance, on the ECL dataset with $T=192$, both MAE and MSE decrease from approximately 0.258 and 0.167 to 0.247 and 0.158, respectively. Such trend of diminishing error seems consistent across different prediction lengths and datasets, supporting the benefit of learning to forecast in the frequency domain. Interestingly, the optimal reduction in forecast error typically occurs at α values near 1, such as 0.8 for the ETTh1 dataset, rather than at the absolute value of 1. Therefore, unifying supervision signals from both time and frequency domains brings performance improvement. More empirical evidence on other datasets and forecast models are extensively provided in Appendix E.

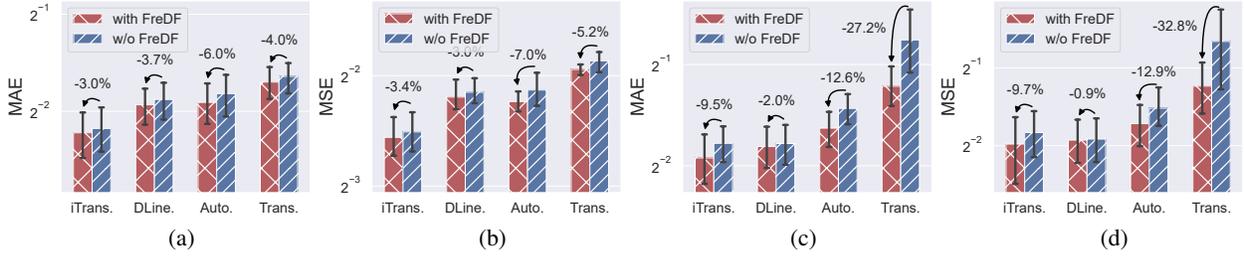


Figure 6: Performance of forecast models with and without FreDF on the ECL (a-b) and Weather (c-d) datasets. Relative reduction of forecast error is reported in percentage⁹.

5.5 Generalization Studies

5.5.1 Generalization to model specifications

In this section, we investigate the generality of FreDF in augmenting diverse neural forecasting models. Specifically, we assess the impact of FreDF on four well-known models: iTransformer, DLinear, Autoformer, and Transformer, with the outcomes illustrated in Figure 6⁴.

Overall, FreDF exhibits remarkable capacity to enhance the performance of these forecasting models. In particular, Transformer-based models, such as the Autoformer and Transformer, benefit substantially from FreDF. For instance, on the ECL dataset, FreDF empowers Autoformer, a model initially introduced in 2021, to outperform DLinear, a cutting-edge model developed in 2023. Such compelling evidence of FreDF’s generality is further detailed in Appendix E. These results underscore the broad applicability of FreDF in enhancing various neural forecast models, establishing its potential as a universally applicable training approach within the realm of time series forecasting.

5.5.2 Generalization to FFT implementations

In this section, we generalize the concept of label autocorrelation: label correlation exists not only between different steps, but also among different variables in multivariate forecasting. Therefore, we implement FFT along the time and variable dimension to handle the corresponding correlations, with the outcomes illustrated in Figure 7.

In general, conducting FFT along the time and variable axis brings similar performance gain, which showcases the existence of correlation between different steps and variables, respectively. In particular, performing FFT on the time axis exhibits slight performance gain, which underscores the relative importance of auto-correlation in the label sequence. Finally, a strategic approach is viewing the multivariate sequence as an image, performing 2-dimensional FFT on both time and variable axis, which further improves the performance of FreDF since it accommodates the correlations between different steps and variables simultaneously.

5.5.3 Generalization to other transforms

In this section, we extend the applicability of FreDF by employing domain transformations based on a diverse array of established polynomial sets. Each polynomial set is adept at capturing specific data patterns, such as trends and periodicity that are challenging to learn in the time domain. The results are summarized in Figure 8.

The performance generally aligns with practical expectations. Notably, projections onto Legendre and Fourier bases demonstrate superior performance. This superiority is attributed to their mutual orthogonality, a feature not guaranteed by other polynomial sets without typical weighting factors, as detailed in Appendix C. It highlights the significance of orthogonality in the selection of basis sets for the FreDF paradigm which is pivotal in effectively managing autocorrelations.

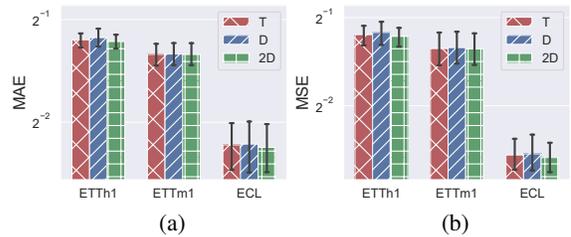


Figure 7: Performance of FreDF with different FFT implementations. *T* and *D* refers to using 1-dimensional FFT along the time and feature dimension, respectively; *2D* refers to using 2-dimensional FFT on both dimensions⁹.

⁴The forecast errors are averaged over prediction lengths with error bars representing 95% confidence intervals.

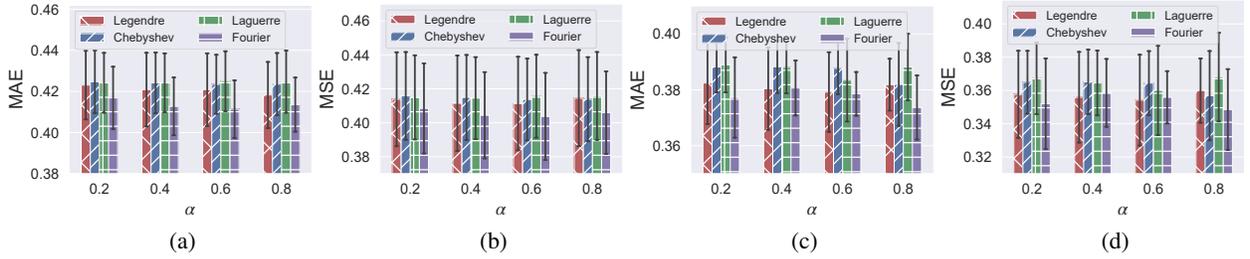


Figure 8: Performance given different implementation of domain transformation. Experiments are conducted on ETTh1 (a-b) and ETTm1 (c-d) with forecast lengths being 192 and 336.

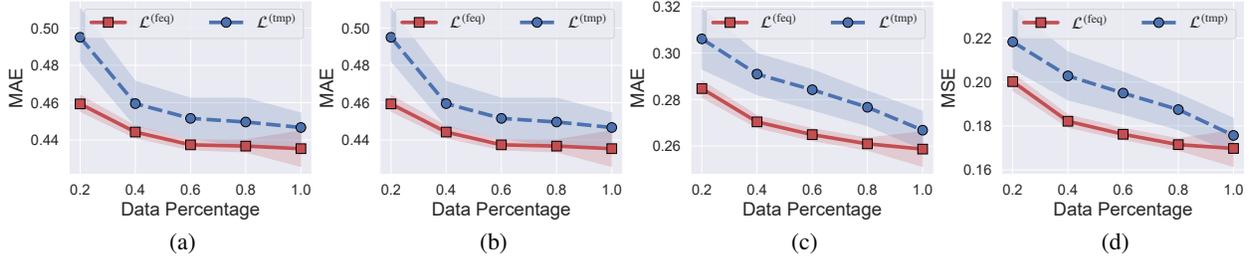


Figure 9: Learning curve on the ETTh1 (a-b) and ECL (c-d) datasets.

5.6 Learning-curve Analysis

In this section, we investigate the sample efficiency of learning in the time versus frequency domains, with the corresponding learning curves showcased in Figure 9. Notably, given limited training data, learning in the frequency domain demonstrates remarkable efficacy. Specifically, with only 30% of the training data, it achieves performance comparable to time domain learning using the full training dataset.

The underlying reason for this enhanced sample efficiency can be attributed to the consistent and more straightforward nature of the data representation. For instance, a sliding window on a sine signal yields a set of distinct sequences in the time domain. However, in the frequency domain, these sequences present a similar pattern: a prominent spike at a specific frequency and negligible values elsewhere. This uniformity simplifies the learning process, as the patterns are more consistent and straightforward to decipher, thereby reducing the reliance on extensive training datasets.

6 Conclusion

This study investigates the label correlation issue in the field of time series modeling, which leads to a discrepancy between the model assumption and data properties. To handle this issue, we develop FreDF, which bypasses label autocorrelation by learning to forecast in the frequency domain. In the frequency domain where bases are orthogonal and independent, a significant reduction of autocorrelation effect is observed. Experiments showcase FreDF’s capability and adaptability across various tasks and forecast models.

Limitation & future works. In this work, we mainly employ the Fourier transform for domain transformation. Despite empirical efficacy, the predefined set of sine bases lacks the ability to adapt to specific data properties. Alternative transforms such as PCA can produce orthogonal bases that better align with data properties, representing a valuable avenue for future research. Additionally, the issue of label autocorrelation extends beyond time series, affecting diverse contexts involving structural labels, such as 3D point clouds, speech, and images. The potential of FreDF to enhance performance in these contexts warrants further exploration.

Broader Impact

Time series modeling is a fundamental field in machine learning, with diverse potential applications in the real world, none of which we feel must be specifically highlighted here. This study contributes to advancing the field by addressing the effects of label correlations, a factor we believe to be pivotal for both the theoretical understanding and practical application for time series modeling. We hold the belief that the issue of label autocorrelation is not confined solely to time series data, pervading various fields where structural labels play a critical role: 3D point clouds, speech, and images. A common oversight in these domains is the treatment of interconnected components—such as pixels in vision tasks—as independent entities *within the learning objective*, which neglects the inherent correlations between these components and therefore limiting the performance. The FreDF paradigm, a significant stride towards mitigating this label autocorrelation issue, has potential to enhance various aspects of machine learning.

References

- [1] Dimitros Asteriou and Stephen G Hall. Arima models and the box–jenkins methodology. *Appl. Econ.*, 2(2):265–286, 2011.
- [2] Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- [3] Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, 619(7970):533–538, 2023.
- [4] Defu Cao, Yujing Wang, Juanyong Duan, Ce Zhang, Xia Zhu, Congrui Huang, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong, et al. Spectral temporal graph neural network for multivariate time-series forecasting. In *NeurIPS*, volume 33, pages 17766–17778, 2020.
- [5] Cristian Challu, Kin G Olivares, Boris N Oreshkin, Federico Garza, Max Mergenthaler, and Artur Dubrawski. N-hits: Neural hierarchical interpolation for time series forecasting. *arXiv preprint arXiv:2201.12886*, 2022.
- [6] Si-An Chen, Chun-Liang Li, Nate Yoder, Sercan Ö. Arik, and Tomas Pfister. Tsmixer: An all-mlp architecture for time series forecasting. *TMLR*, 2023.
- [7] Zhichao Chen, Leilei Ding, Zhixuan Chu, Yucheng Qi, Jianmin Huang, and Hao Wang. Monotonic neural ordinary differential equation: Time-series forecasting for cumulative data. In *CIKM*, pages 4523–4529. ACM, 2023.
- [8] Abhimanyu Das, Weihao Kong, Andrew Leach, Rajat Sen, and Rose Yu. Long-term forecasting with tide: Time-series dense encoder. *arXiv preprint arXiv:2304.08424*, 2023.
- [9] Vijay Ekambaram, Arindam Jati, Nam Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. Tsmixer: Lightweight mlp-mixer model for multivariate time series forecasting. In *SIGKDD*, page 459–469, 2023.
- [10] Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *ICLR*, 2021.
- [11] 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. In *NeurIPS*, 2023.
- [12] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [13] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *ICLR*, 2020.
- [14] Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *SIGIR*, 2018.
- [15] Haoxuan Li, Kunhan Wu, Chunyuan Zheng, Yanghao Xiao, Hao Wang, Zhi Geng, Fuli Feng, Xiangnan He, and Peng Wu. Removing hidden confounding in recommendation: A unified multi-task learning approach. In *NeurIPS*, 2023.
- [16] Jianxin Li, Xiong Hui, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *AAAI*, 2021.
- [17] Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. Scinet: time series modeling and forecasting with sample convolution and interaction. In *NeurIPS*, 2022.
- [18] Shiyu Liu, Rohan Ghosh, and Mehul Motani. Towards better long-range time series forecasting using generative forecasting. *CoRR*, abs/2212.06142, 2022.
- [19] 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. 2021.

- [20] Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In *ICLR*, 2024.
- [21] Yong Liu, Chenyu Li, Jianmin Wang, and Mingsheng Long. Koopa: Learning non-stationary time series dynamics with koopman predictors. In *NeurIPS*, 2023.
- [22] Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Rethinking the stationarity in time series forecasting. In *NeurIPS*, 2022.
- [23] Zhaoran Liu, Yizhi Cao, Hu Xu, Yuxin Huang, Qunshan He, Xinjie Chen, Xiaoyu Tang, and Xinggao Liu. Hidformer: Hierarchical dual-tower transformer using multi-scale mergence for long-term time series forecasting. *Expert Syst Appl.*, 239:122412, 2024.
- [24] Xin Ma, Dehao Wu, Shaoxu Gao, Tongze Hou, and Youqing Wang. Autocorrelation feature analysis for dynamic process monitoring of thermal power plants. *IEEE Trans. Cybern.*, 53(8):5387–5399, 2023.
- [25] Gonzalo Mateos, Santiago Segarra, Antonio G. Marques, and Alejandro Ribeiro. Connecting the dots: Identifying network structure via graph signal processing. *IEEE Signal Process. Mag.*, 36(3):16–43, 2019.
- [26] Zelin Ni, Hang Yu, Shizhan Liu, Jianguo Li, and Weiyao Lin. Basisformer: Attention-based time series forecasting with learnable and interpretable basis. In *NeurIPS*, 2023.
- [27] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *ICLR*, 2023.
- [28] Boris N Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. In *ICLR*, 2019.
- [29] David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *Int. J. Forecast.*, 36(3):1181–1191, 2020.
- [30] Yajuan Si, Mari Palta, and Maureen Smith. Bayesian profiling multiple imputation for missing hemoglobin values in electronic health records. *Ann. Appl. Stat.*, 14(4):1903, 2020.
- [31] Souhaib Ben Taieb and Amir F Atiya. A bias and variance analysis for multistep-ahead time series forecasting. *IEEE Trans. Neural. Netw. Learn. Syst.*, 27(1):62–76, 2015.
- [32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- [33] Hao Wang, Tai-Wei Chang, Tianqiao Liu, Jianmin Huang, Zhichao Chen, Chao Yu, Ruopeng Li, and Wei Chu. ESCM2: entire space counterfactual multi-task model for post-click conversion rate estimation. In *SIGIR*, pages 363–372, 2022.
- [34] Hao Wang, Zhichao Chen, Jiajun Fan, Haoxuan Li, Tianqiao Liu, Weiming Liu, Quanyu Dai, Yichao Wang, Zhenhua Dong, and Ruiming Tang. Optimal transport for treatment effect estimation. In *NeurIPS*, 2023.
- [35] 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 Trans. Artif. Intell.*, 2023.
- [36] Mark W. Watson. Vector autoregressions and cointegration. *Working Paper Series, Macroeconomic Issues*, 4, 1993.
- [37] Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *ICLR*, 2023.
- [38] Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Flowformer: Linearizing transformers with conservation flows. In *ICML*, 2022.
- [39] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with Auto-Correlation for long-term series forecasting. In *NeurIPS*, 2021.
- [40] Haixu Wu, Hang Zhou, Mingsheng Long, and Jianmin Wang. Interpretable weather forecasting for worldwide stations with a unified deep model. *Nat. Mach. Intell.*, pages 1–10, 2023.
- [41] Zhijian Xu, Ailing Zeng, and Qiang Xu. Fits: Modeling time series with $10k$ parameters. In *ICLR*, 2024.
- [42] Kun Yi, Qi Zhang, Wei Fan, Hui He, Liang Hu, Pengyang Wang, Ning An, Longbing Cao, and Zhendong Niu. Fouriergnn: Rethinking multivariate time series forecasting from a pure graph perspective. In *NeurIPS*, 2023.
- [43] Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian, Longbing Cao, and Zhendong Niu. Frequency-domain mlps are more effective learners in time series forecasting. In *NeurIPS*, 2023.
- [44] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *AAAI*, 2023.

-
- [45] Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *ICLR*, 2023.
 - [46] Zheng Zhao, Weihai Chen, Xingming Wu, Peter CY Chen, and Jingmeng Liu. Lstm network: a deep learning approach for short-term traffic forecast. *IET Intell. Transp. Syst.*, 11(2):68–75, 2017.
 - [47] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *ICML*, 2022.

A Overview of DML for Causation Strength Estimation

A.1 Motivation

In this section, we explore the rationale for employing causal inference methods, particularly double machine learning (DML), to quantify the correlations of interest as discussed in Sections 3.2 and 4.1. According to Figure 1, our focus is on the autocorrelation represented by $Y_t \rightarrow Y_{t'}$ where $0 \leq t < t' < T$. However, the exogenous fork structure $Y_t \leftarrow L(n) \rightarrow Y_{t'}$ creates a pseudo correlation between $Y_{t'}$ and Y_t , as established in probabilistic graph analysis. In this case, the autocorrelation $Y_t \rightarrow Y_{t'}$ is confounded by the pseudo correlations from the fork structure, rendering traditional correlation measures, such as Pearson correlation, ineffective for quantifying the autocorrelation $Y_t \rightarrow Y_{t'}$.

To effectively address the confounding effect introduced by the fork structure, it is essential to employ causal inference methods, which excel in quantifying the causation while eliminating the confounding effect. Among various causal inference approaches, DML is chosen for three key reasons: (1) its model-agnostic nature, which does not depend on specific machine learning model specifications; (2) its ability to handle continuous treatments, which is crucial since Y_t is continuous; and (3) its ease of implementation and independence from exhaustive hyperparameter tuning. DML offers a robust and reliable quantification to the autocorrelation that we care about.

A.2 Method

In this section, we detail the implementation of DML, a two-step procedure designed for estimating causal effects. To align with standard causal inference notation, we define $\mathcal{T} \in \mathbb{R}$ as the treatment variable, $\mathcal{Y} \in \mathbb{R}$ as the outcome variable, $\mathcal{X} \in \mathbb{R}^D$ as the confounder variable that needs to be controlled. The implementation of DML is depicted in Figure 10 (b) which consists of two steps below.

- **Orthogonalization.** This step involves orthogonalizing both the outcome (\mathcal{Y}) and the treatment (\mathcal{T}) with respect to the confounders (\mathcal{X}). To this end, we first use two machine learning models, namely ϕ and ψ , to predict the outcome and the treatment based on covariates. These predictions aim to capture the components in \mathcal{Y} and \mathcal{T} that are influenced by the confounder \mathcal{X} . Subsequently, such impact of \mathcal{X} can be eliminated by calculating the residuals:

$$\begin{aligned}\tilde{\mathcal{Y}} &= \mathcal{Y} - \phi(\mathcal{X}), \\ \tilde{\mathcal{T}} &= \mathcal{T} - \psi(\mathcal{X}).\end{aligned}\tag{5}$$

- **Regression.** This step involves regressing the orthogonalized outcome $\tilde{\mathcal{Y}}$ on the orthogonalized treatment $\tilde{\mathcal{T}}$. A linear regression model is utilized for this purpose:

$$\tilde{\mathcal{Y}} = \beta \tilde{\mathcal{T}} + \epsilon,\tag{6}$$

where ϵ is the error term; β is the model coefficient that can be identified via ordinary least squares. The β can be identified in a supervised learning manner, with objective to minimize the MSE of the prediction and real values. The identified β is the estimated causal effect of the treatment on the outcome, which has eliminated the confounding effect.

By regressing the orthogonalized outcome on the orthogonalized treatment, DML captures the direct effect of the treatment on the outcome without the influence of confounding variables, as depicted in Figure 10 (c). That is, DML isolates the desired causal effect $\mathcal{T} \rightarrow \mathcal{Y}$ from the confounding correlation $\mathcal{T} \leftarrow \mathcal{X} \rightarrow \mathcal{Y}$.

A.3 Experimental Settings

In this section, we outline the experimental settings implemented to employ DML for quantifying the correlations of interest.

General settings. For the base learners ϕ and ψ , we opt for a linear regression model optimized using ordinary least squares for its efficiency⁵. Following Appendix A.1, we treat the history sequence L as the confounder to adjust, and simplify the process by considering the last step in L as representative. Moreover, we focus exclusively on the correlations within the last feature of each dataset⁶. This focus makes Y a scalar value within the real number space rather than a D -dimensional vector in this experiment.

⁵The linear regression model, chosen for its computational efficiency, is crucial in managing the experiment’s scale, where the total number of DML estimators can be exceedingly high (e.g., 36,864 for $T=192$). This selection is justified as other more complex models, like random forests, do not significantly alter the results in our experiments.

⁶This focus is aligned with the study’s objective of analyzing autocorrelations instead of inter-feature correlations, which simplifies the interpretation of results.

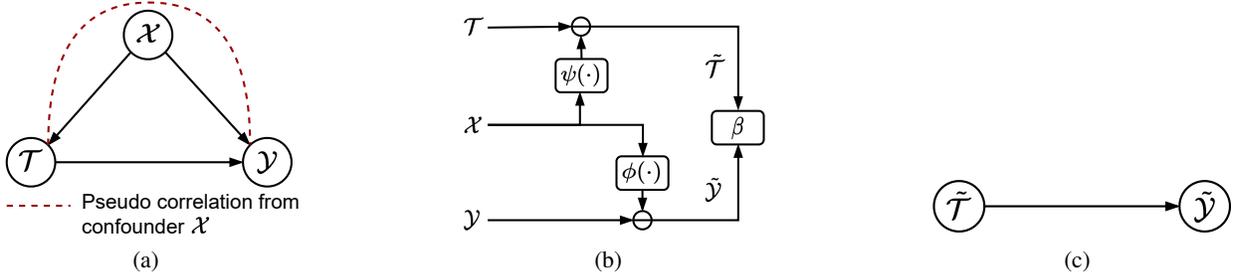


Figure 10: Visualization of confounding effect and DML approach for causation quantification. (a) The causal graph where the pseudo correlation is caused by the fork structure $\mathcal{T} \leftarrow \mathcal{X} \rightarrow \mathcal{Y}$. (b) The implementation of DML, where β is the identified strength of the causation $\mathcal{T} \rightarrow \mathcal{Y}$. (c) The correlations identified by DML, which is indicative to the expected causation $\mathcal{T} \rightarrow \mathcal{Y}$.

Specifications for identifying time-domain causations. To assess the causation $Y_t \rightarrow Y_{t'}$, we treat Y_t as the treatment and $Y_{t'}$ as the outcome. The DML model is trained using a set of N observations: $\{L(n)\}_{n=1:N}$, $\{Y_t(n)\}_{n=1:N}$, and $\{Y_{t'}(n)\}_{n=1:N}$. The coefficient β derived from the DML model is interpreted as the strength of the causation $Y_{n,t} \rightarrow Y_{n,t'}$.

Specifications for identifying frequency-domain causations. To quantify the causation $F_k \rightarrow F_{k'}$, we treat F_k as the treatment and $F_{k'}$ as the outcome. The DML model is trained using a set of N observations: $\{L(n)\}_{n=1:N}$, $\{F_k(n)\}_{n=1:N}$, and $\{F_{k'}(n)\}_{n=1:N}$. The coefficient β derived from the DML model is interpreted as the strength of the causation $F_k \rightarrow F_{k'}$. A notable complexity arises due to F_k being a complex number. Since DML and similar causal inference methods are typically designed for real numbers instead of complex numbers, the identification in this context entails separate consideration of the real and imaginary parts of F_k .

A.4 More Experimental Results

In this section, we provide comprehensive results of the identified causation strengths, which mirrors the autocorrelation effect in the time and frequency domain. We first present the results on three different datasets: Traffic, ETTh1, and ECL in Figure 11, with prediction length set to 192. Subsequently, we present the results given varying prediction lengths: 48, 96, 192, 336 in Figure 12, based on the ECL dataset.

The experimental results show similar patterns with those reported in the main text. Specifically, the non-diagonal elements in Figure 11 (a-c) and Figure 12 (a-d) demonstrate significant values, which affirms the presence of label autocorrelation in the time domain. In contrast, the non-diagonal elements in Figure 11 (d-i) and Figure 12 (e-l) show negligible values, which suggests that frequency components of F are almost independent given L .

In a nutshell, these findings verify the existence of label autocorrelation in the time domain which contradicts the independence assumption of the DF paradigm. By transforming to the frequency domain, the dependency raised by label autocorrelation is largely bypassed, which aligns with DF’s independence assumption as per Theorem 1.

B Theoretical Justification

Theorem B.1. *Given input sequence L and label sequence Y , the negative log-likelihood of the model g parameterized by θ can be expressed as (1) only if Y_t is independent on $Y_{t'}$ for any $1 \leq t, t' \leq T$ given L , expressed as $Y_t \perp\!\!\!\perp Y_{t'} \mid L$.*

Proof. The proof follows establishing the Mean Squared Error (MSE) loss from the maximum likelihood principle, a common approach in statistical analysis. What is new in this proof lies in considering multiple outputs with interdependence, an aspect often neglected in prior research. Since the focus is on autocorrelation between time steps, we consider the case $D = 1$ without loss of generality.

Firstly, define a dataset with N pairs of samples:

$$\mathcal{D} = [(L(1), Y(1)), (L(2), Y(2)), \dots, (L(N), Y(N))], \quad (7)$$

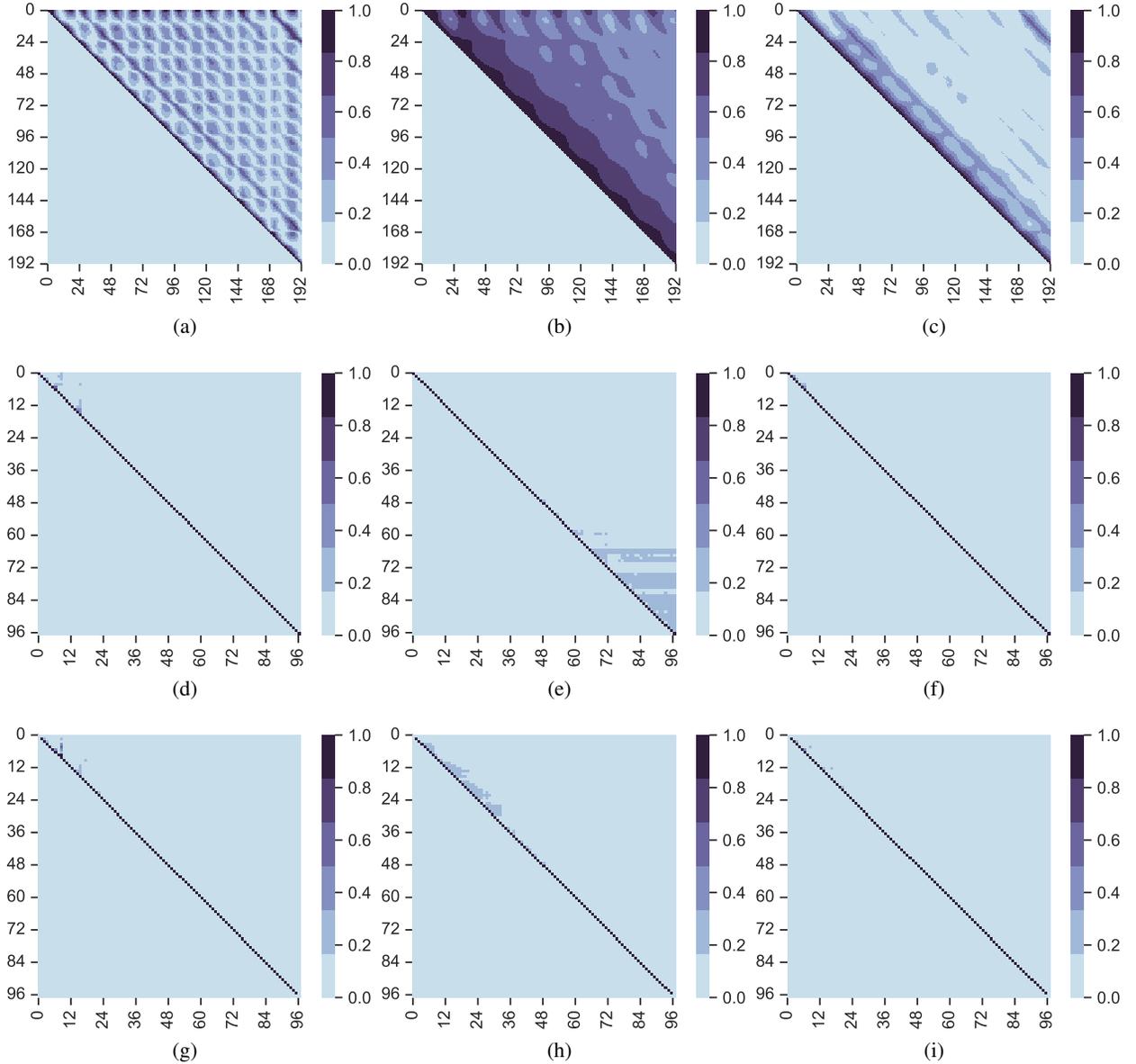


Figure 11: More comprehensive visualizations of label autocorrelation in different domains and datasets, with columns representing different datasets: Traffic, ETTh1, and ECL, from left to right. Panels (a-d) depict the label correlation in the time domain, where each matrix element at the i -th row and j -th column quantifies the absolute causation strength from Y_i to Y_j . Panels (e-h) illustrate the label correlation in the frequency domain, specifically the absolute causation strength from F_i to F_j , calculated using the real parts of the frequency components. Panels (i-l) showcase label correlation in the frequency domain, utilizing the imaginary parts of the frequency components to calculate the absolute causation strength from F_i to F_j . In these experiments, each element is elucidated using DML.

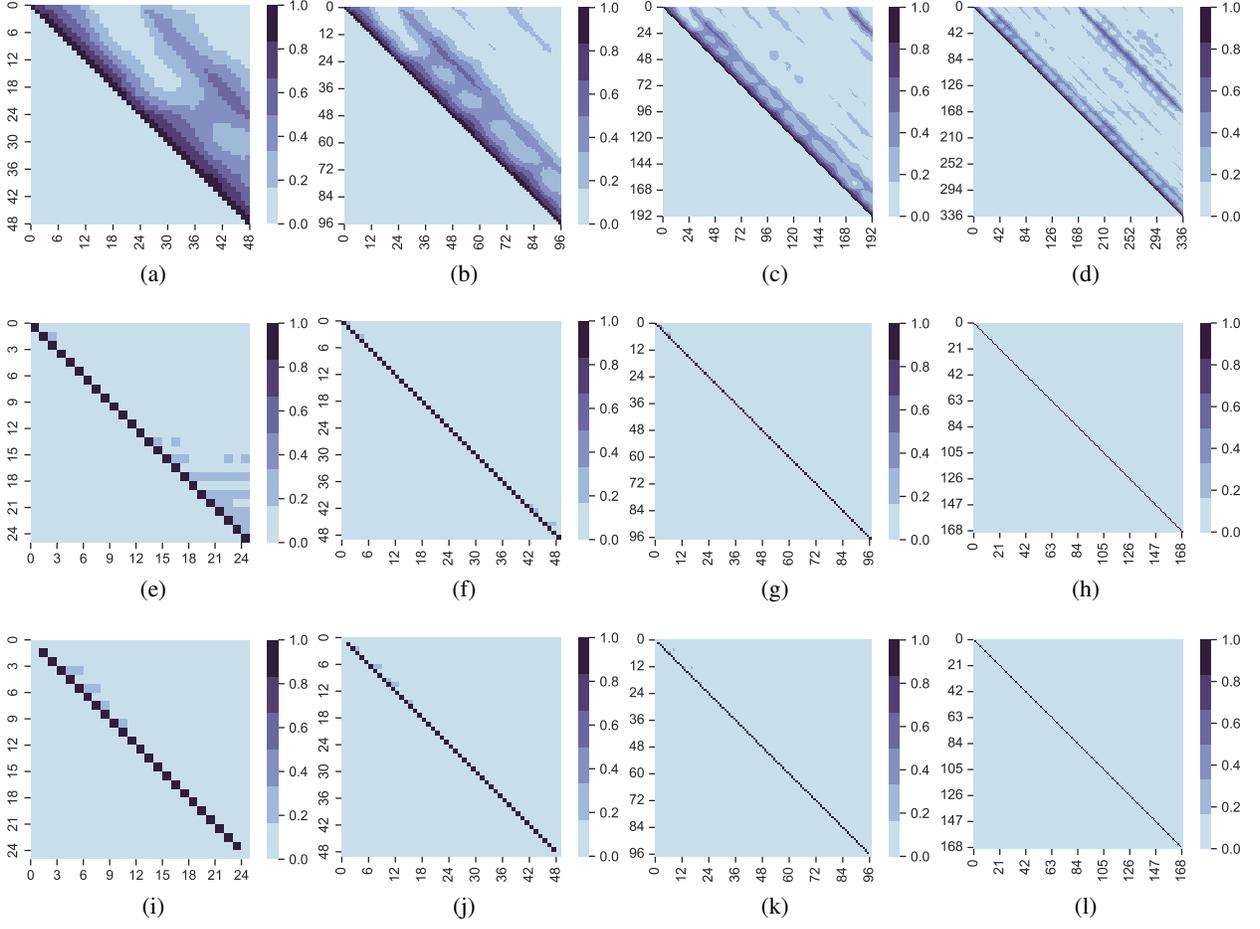


Figure 12: More comprehensive visualizations of label autocorrelation in different domains and label lengths, with columns representing label lengths $H=48, 96, 192, 336$ from left to right. Panels (a-d) depict the label correlation in the time domain, where each matrix element at the i -th row and j -th column quantifies the absolute causation strength from Y_i to Y_j . Panels (e-h) illustrate the label correlation in the frequency domain, specifically the absolute causation strength from F_i to F_j , calculated using the real parts of the frequency components. Panels (i-l) showcase label correlation in the frequency domain, utilizing the imaginary parts of the frequency components to calculate the absolute causation strength from F_i to F_j . In these experiments, each element is elucidated using DML.

where $L(n) \in \mathbb{R}^L$ and $Y(n) \in \mathbb{R}^T$ are the n -th observations of the input and label sequence, respectively. The DF paradigm generates forecasts as $\hat{Y}(n) = g_\theta(L(n))$. The likelihood of θ with respect to data can be expressed as

$$\begin{aligned}
 p(\mathcal{D} | \theta) &= p([(L(1), Y(1)), \dots, (L(N), Y(N))] | \theta) \\
 &= \prod_{n=1}^N p((L(n), Y(n)) | \theta).
 \end{aligned} \tag{8}$$

To identify θ using mean squared loss, assuming that given L , Y obeys a normal distribution $\mathcal{N}(g(L), \zeta)$ with mean $g(L) \in \mathbb{R}^T$ and covariance matrix $\zeta \in \mathbb{R}^{T \times T}$. The negative log-likelihood can be written as:

$$\begin{aligned}
 -\log p(\mathcal{D} | \theta) &= -\log \prod_{n=1}^N p((L(n), Y(n)) | \theta) \\
 &\stackrel{(1)}{=} -\log \prod_{n=1}^N \frac{1}{(2\pi)^{0.5T} |\zeta|^{0.5}} \exp\left(-\frac{1}{2}(Y(n) - g_\theta(L(n)))^\top \zeta^{-1} (Y(n) - g_\theta(L(n)))\right) \\
 &\stackrel{(2)}{=} -\log \prod_{n=1}^N \frac{1}{(2\pi)^{0.5T} \sigma} \exp\left(-\frac{1}{2\sigma^2} \|Y(n) - g_\theta(L(n))\|_2^2\right) \\
 &\stackrel{(3)}{=} c_1 + c_2 * \sum_{n=1}^N \|Y(n) - g_\theta(L(n))\|_2^2,
 \end{aligned} \tag{9}$$

with step-by-step derivations as follows:

- (1) can be derived by specifying p as the probability density function of the normal distribution $\mathcal{N}(g(L(n)), \zeta)$;
- (2) can be derived by assuming $\zeta = \sigma^2 I$ with I being identical matrix, which inherently assumes independence among the elements in Y given L , represented as $Y_t \perp\!\!\!\perp Y_{t'} | L^T$;
- (3) simplifies the equation as a MSE loss with two ignorable constants $c_1 = N \log((2\pi)^{0.5T} \sigma)$ and $c_2 = 1/2\sigma^2$.

Therefore, deriving the equivalence between the MSE loss and negative log-likelihood hinges on the assumption of conditional independence of Y given L (as per the derivation of step (2)); conversely, given conditional independence of Y given L , we can derive the equivalence between the MSE loss and negative log-likelihood. The proof is therefore completed.

The proof can be extended to multivariate time series ($D \neq 1$). Overall, the derivation is identical to (9), with the mean vector being $T \cdot D$ -dimensional vector $g(L(n)) \in \mathbb{R}^{T \cdot D}$ and the covariance matrix being $\zeta \in \mathbb{R}^{T \cdot D \times T \cdot D}$. We assume that different variates in the label sequence are independent⁸, making ζ a block diagonal matrix $\zeta = [\zeta_1, \zeta_2, \dots, \zeta_D]$ where $\zeta_d \in \mathbb{R}^{T \times T}$ denotes the covariance of different steps for the d -th variate. Then, similar to the step (2) in (9), conditional independence assumption is necessary to make $\zeta_d = \sigma^2 I$ for subsequent derivation. \square

C Discussion of Transformation onto Different Bases

Transforming time series data onto predefined spaces is a fundamental aspect of signal processing and data analysis, with various strategies available depending on the choice of bases. The transformation is implemented by projecting the original signal onto a different set of predefined bases, such as the Fourier bases, Legendre bases, and Chebyshev bases. These bases are known for their mutual orthogonality, and the selection of bases depends on the specific characteristics and requirements of the analysis. We provide some formal definition of prevalent transformations below, where we formulate signals as continuous functions for the ease of demonstration.

Fourier transform. It employs sinusoidal functions as bases which prove to be mutually orthogonal. These polynomials are particularly effective for analyzing periodic signals or signals with a strong frequency component. Let k be the frequency, the associated basis function and projection onto it can be formulated as follows:

$$\begin{aligned}
 f_k(t) &= \exp(-j(2\pi/L)kt), \\
 F_k &= \int_{-\infty}^{\infty} x(t) f_k(t) dt
 \end{aligned} \tag{10}$$

Legendre transform. It uses the Legendre polynomials as bases which prove to be mutually orthogonal on the interval $[-1, 1]$. These polynomials are particularly useful for representing functions defined on a finite interval, which

⁷The diagonals can be different values $\sigma_1^2, \sigma_2^2, \dots, \sigma_T^2$. We set them to σ for clarity of derivation.

⁸The assumption also adheres to current DF paradigm which simply adds the forecast error for different variates [16, 18].

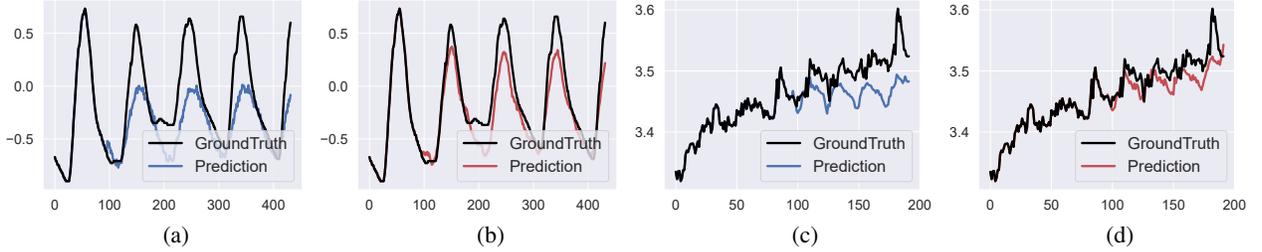


Figure 13: Forecast sequences generated by iTransformer on the snapshots where trend (a-b) and periodicity (c-d) are dominant. (a, c) indicates models trained using canonical DF paradigm; (b, d) indicates models trained with FreDF. The bases are selected as Fourier (b) and Legendre (d) polynomials.

makes them suitable for certain types of data smoothing and approximation tasks. The k -th Legendre polynomial and the associated projection can be formulated as follows:

$$f_k(t) = \frac{1}{2^k k!} \frac{d^k}{dt^k} [(t^2 - 1)^k],$$

$$F_k = \int_{-1}^1 x(t) f_k(t) dt$$
(11)

Chebyshev transform. It uses the Chebyshev polynomials as bases. These bases are NOT originally orthogonal, but can be proved mutually orthogonal on the interval $[-1, 1]$ with respect to the weight $1/\sqrt{1-t^2}$. These polynomials are particularly useful for approximating functions with rapid variations. The k -th Chebyshev polynomial and the associated projection can be formulated as follows, where weighting factor accounts for the varying density of Chebyshev nodes, making this basis well-suited for numerical computations and function approximations.

$$f_k(t) = \cos(k \arccos(t)),$$

$$F_k = \int_{-1}^1 \frac{x(t) f_k(t)}{\sqrt{1-t^2}} dt$$
(12)

Laguerre transform. It uses the Laguerre polynomials as bases. These bases are NOT originally orthogonal, but can be proved mutually orthogonal on the interval $[0, \infty]$ with respect to the exponential weight $\exp(-t)$. These polynomials are particularly useful in quantum mechanics and other fields involving exponential decay. The k -th Laguerre polynomial and the associated projection can be formulated as follows:

$$f_k(t) = \exp(t) \frac{d^k}{dt^k} (\exp(-t) t^k),$$

$$F_k = \int_0^\infty \frac{x(t) f_k(t)}{\exp(t)} dt$$
(13)

These polynomial sets are adept at capturing specific data patterns, such as trends and periodicity that are challenging to learn in the time domain. Their efficacy in FreDF is depicted in Figure 13. Specifically, learning in the time domain fails to capture the increasing trends or follow the high-frequency periods. The involvement of FreDF largely handles the issues and improves the forecast quality.

In summary, the choice of an orthogonal basis for transforming time series data—whether it be Fourier, Legendre, or Chebyshev—depends on the nature of the data and the specific objectives of the analysis. Each basis has unique properties that make it suitable for different types of applications. Understanding these properties is crucial for effectively employing these transformation strategies in time series analysis.

D Reproduction Details

D.1 Dataset Descriptions

The datasets utilized in this study encompass a wide range of time series data, each with its unique characteristics and temporal resolutions:

Table 6: Detailed dataset descriptions. D denotes the number of variates. *Forecast Length* denotes the prediction lengths investigated in this dataset. *Frequency* denotes the sampling interval of time points. *Train, Validation, Test* denotes the number of samples employed in each split. The taxonomy and statistic are built upon recent works [37, 20].

Dataset	D	Forecast Length	Train / validation / test	Frequency	Domain
ETTh1	7	96, 192, 336, 720	8545/2881/2881	Hourly	Health
ETTh2	7	96, 192, 336, 720	8545/2881/2881	Hourly	Health
ETTh1	7	96, 192, 336, 720	34465/11521/11521	15min	Health
ETTh2	7	96, 192, 336, 720	34465/11521/11521	15min	Health
Weather	21	96, 192, 336, 720	36792/5271/10540	10min	Weather
ECL	321	96, 192, 336, 720	18317/2633/5261	Hourly	Electricity
Traffic	862	96, 192, 336, 720	12185/1757/3509	Hourly	Transportation

- ETT [16] comprises data on 7 factors related to electricity transformers, collected from July 2016 to July 2018. This dataset is divided into four subsets: ETTh1 and ETTh2, with hourly recordings, and ETTm1 and ETTm2, documented every 15 minutes.
- Weather [39] includes 21 meteorological variables gathered every 10 minutes throughout 2020 from the Weather Station of the Max Planck Biogeochemistry Institute.
- ECL (Electricity Consumption Load) [39] presents hourly electricity consumption data for 321 clients.
- Traffic [39] features hourly road occupancy rates from 862 sensors in the San Francisco Bay area freeways, spanning from January 2015 to December 2016.

Data processing and the division into training, validation, and testing sets adhere to the protocol established by TimesNet [37]. This approach ensures chronological order division to prevent data leakage. Regarding forecast settings, the length of the lookback series is standardized at 96 across the ETT, Weather, ECL, and Traffic datasets, with varying prediction lengths of 96, 192, 336, and 720. Further dataset specifics are delineated in Table 6.

D.2 Implementation Details

The baseline models for this study were meticulously reproduced using training scripts obtained from the TimesNet Repository [37] after reproducibility verification. Models were trained employing the Adam optimizer [12], with learning rates selected from the set $10^{-3}, 5 \times 10^{-4}, 10^{-4}$ to minimize the MSE loss. A consistent batch size of 32 was employed across all models. The training regime was capped at a maximum of 10 epochs, incorporating an early stopping mechanism that was activated upon a lack of improvement in validation performance over 3 epochs.

For the integration of the FreDF paradigm into existing models, we closely adhered to the original hyperparameter configurations as specified in their respective publications. The only parameters finetuned were the learning rate and the relative strength of frequency-domain alignment in $[0, 1]$. Finetuning the learning rate was essential to accommodate huge disparities in the magnitude of MSE loss observed between the time and frequency domains. Fine-tuning was conducted to minimize the MSE averaged across all prediction lengths on the validation dataset. In fact, more practical approach is finetuning for each prediction length separately, while we omit it since the efficacy of FreDF does not rely on dedicate hyperparameter configurations, and current results suffice to showcase the efficacy of FreDF.

E More Experimental Results

E.1 Overall Performance

Long-term forecast. We provide comprehensive performance comparison on the long-term forecast task in Table 7. The iTransformer model is employed to operationalize the FreDF paradigm. Despite the iTransformer’s existing performance gap compared to other baseline models, the incorporation of FreDF enhances its performance in the majority of cases, securing the lowest MSE in 30 out of 35 cases and MAE in all 35 cases. The consistent improvement across nearly all scenarios underscores FreDF’s robustness. The few instances where FreDF does not achieve the lowest MSE is attributed to the inherent advantages of other models over the iTransformer in specific contexts (for example, FreTS versus iTransformer on the Weather dataset).

Table 7: Full results on the long-term forecasting task with forecast lengths T=96, 192, 336 and 720. The length of history window is set to 96 for all baselines. Avg indicates the results averaged over forecasting lengths.

Models	FreDF (Ours)		iTransformer (2024)		FreTS (2023)		TimesNet (2023)		Crossformer (2023)		TiDE (2023)		DLinear (2023)		FEDformer (2022)		Autoformer (2021)		Transformer (2017)		TCN (2017)		LSTM (1998)		
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
ETTm1	96	0.324	0.362	0.346	0.379	0.339	0.374	0.338	0.379	0.375	0.415	0.364	0.387	0.345	0.372	0.389	0.427	0.468	0.463	0.591	0.549	0.887	0.613	0.639	0.557
	192	0.373	0.385	0.392	0.400	0.382	0.397	0.389	0.400	0.453	0.474	0.398	0.404	0.381	0.390	0.402	0.431	0.573	0.509	0.704	0.629	0.877	0.626	0.624	0.546
	336	0.402	0.404	0.427	0.422	0.421	0.426	0.429	0.428	0.548	0.526	0.428	0.425	0.414	0.414	0.438	0.451	0.596	0.527	1.171	0.861	0.890	0.636	0.655	0.568
	720	0.469	0.444	0.494	0.461	0.485	0.462	0.495	0.464	0.857	0.713	0.487	0.461	0.473	0.451	0.529	0.498	0.749	0.569	1.307	0.893	0.911	0.653	0.706	0.598
	Avg	0.392	0.399	0.415	0.416	0.407	0.415	0.413	0.418	0.558	0.532	0.419	0.419	0.404	0.407	0.440	0.451	0.596	0.517	0.943	0.733	0.891	0.632	0.656	0.567
ETTm2	96	0.173	0.252	0.184	0.266	0.190	0.282	0.185	0.264	0.267	0.349	0.207	0.305	0.195	0.294	0.194	0.284	0.240	0.319	0.317	0.408	3.125	1.345	1.488	0.875
	192	0.241	0.298	0.257	0.315	0.260	0.329	0.254	0.307	0.472	0.479	0.290	0.364	0.283	0.359	0.264	0.324	0.300	0.349	1.069	0.758	3.130	1.350	1.514	0.897
	336	0.298	0.334	0.315	0.351	0.373	0.405	0.314	0.345	0.919	0.702	0.377	0.422	0.384	0.427	0.319	0.359	0.339	0.375	1.325	0.869	3.185	1.375	1.608	0.942
	720	0.398	0.393	0.419	0.409	0.517	0.499	0.434	0.413	4.874	1.601	0.558	0.524	0.516	0.502	0.430	0.424	0.423	0.421	2.576	1.223	4.203	1.658	2.417	1.203
	Avg	0.278	0.319	0.294	0.335	0.335	0.379	0.297	0.332	1.633	0.782	0.358	0.404	0.344	0.396	0.302	0.348	0.326	0.366	1.322	0.814	3.411	1.432	1.757	0.979
ETTth1	96	0.382	0.400	0.390	0.410	0.399	0.412	0.422	0.433	0.441	0.457	0.479	0.464	0.396	0.410	0.377	0.418	0.423	0.441	0.796	0.691	0.767	0.633	0.767	0.633
	192	0.430	0.427	0.443	0.441	0.453	0.443	0.465	0.457	0.521	0.503	0.525	0.492	0.449	0.444	0.421	0.445	0.498	0.485	0.813	0.699	0.739	0.619	0.739	0.619
	336	0.474	0.451	0.480	0.457	0.503	0.475	0.492	0.470	0.659	0.603	0.565	0.515	0.487	0.465	0.468	0.472	0.506	0.496	1.181	0.876	0.717	0.613	0.717	0.613
	720	0.463	0.462	0.484	0.479	0.596	0.565	0.532	0.502	0.893	0.736	0.594	0.558	0.516	0.513	0.500	0.493	0.477	0.487	1.182	0.885	0.828	0.678	0.828	0.678
	Avg	0.437	0.435	0.449	0.447	0.488	0.474	0.478	0.466	0.628	0.574	0.541	0.507	0.462	0.458	0.441	0.457	0.476	0.477	0.993	0.788	0.763	0.636	0.763	0.636
ETTth2	96	0.289	0.337	0.301	0.349	0.350	0.403	0.320	0.364	0.681	0.592	0.400	0.440	0.343	0.396	0.347	0.391	0.383	0.424	2.072	1.140	3.171	1.364	1.678	0.950
	192	0.363	0.385	0.382	0.402	0.472	0.475	0.409	0.417	1.837	1.054	0.528	0.509	0.473	0.474	0.430	0.443	0.557	0.511	5.081	1.814	3.222	1.398	1.749	1.009
	336	0.419	0.426	0.430	0.434	0.564	0.528	0.449	0.451	3.000	1.472	0.643	0.571	0.603	0.546	0.469	0.475	0.470	0.481	3.564	1.475	3.306	1.452	1.814	1.030
	720	0.415	0.437	0.447	0.455	0.815	0.654	0.473	0.474	3.024	1.399	0.874	0.679	0.812	0.650	0.473	0.480	0.501	0.515	2.469	1.247	3.599	1.565	2.025	1.126
	Avg	0.371	0.396	0.390	0.410	0.550	0.515	0.413	0.426	2.136	1.130	0.611	0.550	0.558	0.516	0.430	0.447	0.478	0.483	3.296	1.419	3.325	1.445	1.817	1.029
ECL	96	0.144	0.233	0.148	0.239	0.189	0.277	0.171	0.273	0.148	0.248	0.237	0.329	0.210	0.302	0.200	0.315	0.199	0.315	0.252	0.352	0.688	0.621	0.348	0.420
	192	0.159	0.247	0.167	0.258	0.193	0.282	0.188	0.289	0.161	0.263	0.236	0.330	0.210	0.305	0.207	0.322	0.215	0.327	0.266	0.364	0.587	0.582	0.323	0.400
	336	0.172	0.263	0.179	0.272	0.207	0.296	0.208	0.304	0.191	0.289	0.249	0.344	0.223	0.319	0.226	0.340	0.232	0.343	0.292	0.383	0.590	0.588	0.319	0.398
	720	0.204	0.294	0.209	0.298	0.245	0.332	0.289	0.363	0.226	0.314	0.284	0.373	0.258	0.350	0.282	0.379	0.268	0.371	0.287	0.371	0.602	0.601	0.325	0.406
	Avg	0.170	0.259	0.176	0.267	0.209	0.297	0.214	0.307	0.182	0.279	0.251	0.344	0.225	0.319	0.229	0.339	0.228	0.339	0.274	0.367	0.617	0.598	0.329	0.406
Traffic	96	0.391	0.265	0.397	0.272	0.528	0.341	0.609	0.317	0.518	0.269	0.805	0.493	0.697	0.429	0.577	0.362	0.609	0.385	0.686	0.385	1.451	0.744	0.928	0.513
	192	0.410	0.273	0.418	0.279	0.531	0.338	0.621	0.328	0.551	0.285	0.756	0.474	0.647	0.407	0.603	0.372	0.633	0.400	0.679	0.377	0.842	0.622	0.890	0.491
	336	0.424	0.280	0.432	0.286	0.551	0.345	0.641	0.342	0.546	0.293	0.762	0.477	0.653	0.410	0.615	0.378	0.637	0.398	0.663	0.361	0.844	0.620	0.872	0.476
	720	0.460	0.298	0.467	0.305	0.598	0.367	0.671	0.354	0.597	0.323	0.719	0.449	0.694	0.429	0.649	0.403	0.668	0.415	0.693	0.381	0.867	0.624	0.872	0.469
	Avg	0.421	0.279	0.428	0.286	0.552	0.348	0.636	0.335	0.553	0.292	0.760	0.473	0.673	0.419	0.611	0.379	0.637	0.399	0.680	0.376	1.001	0.652	0.890	0.487
Weather	96	0.164	0.202	0.201	0.247	0.184	0.239	0.178	0.226	0.177	0.246	0.202	0.261	0.197	0.259	0.221	0.304	0.284	0.355	0.332	0.383	0.610	0.568	0.246	0.308
	192	0.220	0.253	0.250	0.283	0.223	0.275	0.227	0.266	0.227	0.297	0.242	0.298	0.236	0.294	0.275	0.345	0.313	0.371	0.634	0.539	0.541	0.552	0.279	0.341
	336	0.275	0.294	0.302	0.317	0.272	0.316	0.283	0.305	0.278	0.346	0.287	0.335	0.282	0.332	0.338	0.379	0.359	0.393	0.656	0.579	0.565	0.569	0.320	0.372
	720	0.356	0.347	0.370	0.362	0.340	0.363	0.359	0.355	0.368	0.407	0.351	0.386	0.347	0.384	0.408	0.418	0.440	0.446	0.908	0.706	0.622	0.601	0.378	0.410
	Avg	0.254	0.274	0.281	0.302	0.255	0.299	0.262	0.288	0.262	0.324	0.271	0.320	0.265	0.317	0.311	0.361	0.349	0.391	0.632	0.552	0.584	0.572	0.306	0.357
1 st Count	30	35	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Table 8: Full results on the short-term forecasting task with forecast lengths: yearly, quarterly, and monthly. Avg indicates the results averaged over forecasting lengths.

Models	FreDF			FreTS			iTransformer			Crossformer			DLinear			Fedformer			Autoformer		
	(Ours)			(2023)			(2024)			(2023)			(2023)			(2023)			(2023)		
Metric	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA
Yearly	13.556	3.046	0.798	13.576	3.068	0.801	13.797	3.143	0.818	68.344	17.601	4.305	14.307	3.094	0.827	13.648	3.089	0.806	18.477	4.26	1.101
Quarterly	10.374	1.229	0.919	10.361	1.223	0.916	10.503	1.248	0.932	73.822	13.272	8.191	10.5	1.237	0.928	10.612	1.246	0.936	14.254	1.829	1.314
Monthly	12.999	0.983	0.913	13.088	0.99	0.919	13.227	1.013	0.935	68.67	11.269	7.679	13.362	1.007	0.937	14.181	1.105	1.011	18.421	1.616	1.398
Others	5.294	3.614	1.127	5.563	3.71	1.17	5.101	3.419	1.076	98.68	79.677	22.948	5.12	3.649	1.114	4.823	3.243	1.019	6.772	4.963	1.495
Avg.	12.112	1.648	0.877	12.169	1.66	0.883	12.298	1.68	0.893	71.332	16.626	6.977	12.48	1.674	0.898	12.734	1.702	0.914	16.851	2.443	1.26
1 st Count	3	3	3	1	1	1	0	0	0	0	0	0	0	0	0	1	1	1	0	0	0

Short-term forecast. We provide a detailed comparison for the short-term forecast task in Table 8, with FreTS serving as the base model for FreDF implementation. Similar to the long-term forecast results, FreDF enhances FreTS’s performance in most instances. It is noteworthy that some baseline models, initially designed for long-term forecasting, fail to perform optimally in short-term scenarios (such as Crossformer), despite thorough fine-tuning. Interestingly, FreTS exhibits superior performance over FreDF in quarterly forecast lengths. This observation aligns with the expectation that FreDF is optimized to minimize overall average forecast error on the validation set rather than targeting specific forecast lengths. While it is possible to fine-tune FreDF for each forecast length to cater to the distinct properties and optimal hyperparameter settings of different tasks, this approach was not pursued as the current results adequately demonstrate FreDF’s effectiveness.

Missing data imputation. We provide a thorough comparison on the missing data imputation task in Table 9 with varying missing ratios, where iTransformer is selected as the base model for FreDF implementation. Similar to forecast tasks above, FreDF enhances the imputation performance of iTransformer in all instances, hitting the minimum MSE in 24 out of 30 cases and minimum MAE in 21 out of 30 cases. The efficacy lies in the effective capturing of label autocorrelation among non-missing entries.

Showcases. We provide additional showcases illustrating the improvements in forecast sequences by integrating FreDF in Figure 14 and 15. Overall, FreDF effectively eliminates blurs in the forecast sequences and captures high frequency components in the label sequences. These successes are attributed to the unique capability of FreDF to operate in the frequency domain. In this domain, the challenges of autocorrelation are naturally mitigated, and the expression of high-frequency components becomes more straightforward. These factors underly FreDF’s success in elevating the quality of forecast generation.

E.2 Generalization Studies

In this detailed investigation, we further explore the universality of the Frequency-enhanced Direct Forecast (FreDF) paradigm in improving a range of neural forecasting models across diverse datasets. Our analysis encompasses the impact of FreDF on four prominent models: iTransformer, DLinear, Autoformer, and Transformer. The performance improvements facilitated by FreDF are quantitatively presented in Figure 16 across five distinct datasets, with forecast errors averaged over various prediction lengths and error bars denoting 95% confidence intervals.⁹

FreDF demonstrates a significant ability to elevate the performance of these forecasting models, with Transformer-based models like the Autoformer and Transformer experiencing particularly notable enhancements. A case in point is the ECL dataset, where FreDF enables the Autoformer—a model introduced in 2021—to surpass the performance of DLinear, a state-of-the-art model developed in 2023. This and other examples detailed in Appendix E vividly illustrate FreDF’s effectiveness and general applicability.

The results presented here affirm the broad utility of FreDF in augmenting neural forecast models, suggesting its role as a versatile and universally applicable training methodology in the field of time series forecasting. This evidence

⁹See footnote for error bar methodology.

Table 9: Full results on the missing data imputation task with missing ratios 0.125, 0.25, 0.375, 0.5. The length of history window is set to 96 for all baselines. Avg indicates the results averaged over missing ratios.

Models	FreDF (Ours)		iTransformer (2024)		FreTS (2023)		TimesNet (2023)		Crossformer (2023)		TiDE (2023)		DLinear (2023)		FEDformer (2022)		Autoformer (2021)	
	p_{miss}	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
ETTm1	0.125	0.00153 0.02790		0.00213	0.03307	0.01102	0.07843	0.01152	0.07267	0.14570	0.26727	0.45052	0.45514	0.00148 0.02380	0.68262	0.38111	0.37654	0.35378
	0.25	0.00287 0.03801		0.00402	0.04434	0.01089	0.07753	0.01245	0.07946	0.06801	0.17911	0.41777	0.45884	0.00154 0.02351	0.68235	0.38116	0.37059	0.35261
	0.375	0.00256 0.03669		0.00458	0.04663	0.01100	0.07812	0.01407	0.08673	0.03494	0.12612	0.62935	0.55570	0.00175 0.02385	0.68191	0.38105	0.37877	0.36093
	0.5	0.00152 0.02739		0.00363	0.04359	0.01102	0.07818	0.01676	0.09610	0.02696	0.11517	0.29342	0.39320	0.00192 0.02219	0.68119	0.38085	0.38052	0.36462
	Avg	0.00212 0.03250		0.00359	0.04191	0.01098	0.07807	0.01370	0.08374	0.06891	0.17192	0.44776	0.46572	0.00167 0.02334	0.68202	0.38104	0.37660	0.35798
ETTm2	0.125	0.00363 0.03840	0.00398 0.04034		0.03194	0.13349	0.01189	0.06710	0.35271	0.43486	0.83023	0.62174	0.03822	0.12943	3.10388	1.31356	1.40160	0.80777
	0.25	0.00437 0.04255	0.00431 0.04303		0.03591	0.13655	0.01795	0.08939	0.19400	0.31921	0.81402	0.61100	0.03063	0.11547	3.10364	1.31348	1.41033	0.81363
	0.375	0.00352 0.03823	0.00342 0.03793		0.03250	0.13336	0.02742	0.11499	0.12863	0.25738	1.11225	0.73633	0.01709	0.08822	3.10328	1.31330	1.40812	0.81049
	0.5	0.00137 0.02382	0.00160 0.02538		0.03126	0.13027	0.04053	0.14285	0.07394	0.19618	0.99459	0.70665	0.01025	0.06440	3.10527	1.31389	1.44617	0.81796
	Avg	0.00322 0.03575	0.00333 0.03667		0.03290	0.13342	0.02445	0.10358	0.18732	0.30191	0.93777	0.66893	0.02405	0.09938	3.10402	1.31356	1.41655	0.81246
ETTTh1	0.125	0.00178 0.03059		0.00319	0.04102	0.01400	0.08181	0.00441	0.04403	0.21157	0.34941	0.36363	0.45350	0.00279 0.03617	0.68307	0.38026	0.43136	0.41184
	0.25	0.00218 0.03405		0.00334	0.04205	0.01347	0.08097	0.00320	0.03850	0.15197	0.29466	0.28435	0.40516	0.00236 0.03324	0.68162	0.37973	0.43515	0.41584
	0.375	0.00182 0.03108		0.00280	0.03852	0.01308	0.08017	0.00261	0.03540	0.11596	0.25758	0.21038	0.34029	0.00210 0.03121	0.68181	0.37975	0.44431	0.42505
	0.5	0.00114 0.02414	0.00174 0.03008		0.01276	0.07918	0.00245	0.03472	0.07787	0.20468	0.13344	0.27102	0.00175	0.02844	0.68137	0.37992	0.44312	0.42387
	Avg	0.00173 0.02996		0.00277	0.03792	0.01333	0.08053	0.00317	0.03817	0.13935	0.27658	0.24795	0.36749	0.00225 0.03226	0.68197	0.37992	0.43848	0.41915
ETTTh2	0.125	0.00222 0.03124	0.00473 0.04606		0.04485	0.13849	0.00535	0.04495	0.58587	0.54432	1.15859	0.73871	0.02287	0.10885	3.12756	1.31746	1.45130	0.84467
	0.25	0.00407 0.04258	0.00571 0.05096		0.04647	0.13551	0.00494 0.04476		0.33565	0.41741	0.75643	0.59747	0.02491	0.11511	3.12891	1.31754	1.45386	0.84388
	0.375	0.00306 0.03693	0.00452 0.04519		0.04830	0.13583	0.00512	0.04697	0.28196	0.38453	0.59470	0.52371	0.01944	0.10277	3.12788	1.31728	1.45464	0.84194
	0.5	0.00129 0.02365	0.00249 0.03304		0.04900	0.13469	0.00604	0.05224	0.21866	0.32516	0.35775	0.40497	0.01465	0.08746	3.12882	1.31733	1.45997	0.84644
	Avg	0.00266 0.03360	0.00436 0.04381		0.04715	0.13613	0.00536	0.04723	0.35553	0.41785	0.71687	0.56622	0.02046	0.10355	3.12829	1.31740	1.45494	0.84423
ECL	0.125	0.00029 0.01257	0.00187 0.03191		0.01018	0.08255	0.00466	0.04597	0.25009	0.36799	0.32942	0.42254	0.10658	0.23808	0.45884	0.41005	0.20147	0.29003
	0.25	0.00061 0.01846	0.00216 0.03491		0.01022	0.08269	0.00341	0.03978	0.18890	0.32186	0.28831	0.40031	0.10682	0.23654	0.45887	0.41007	0.20618	0.29771
	0.375	0.00090 0.02242	0.00211 0.03473		0.01022	0.08258	0.00230	0.03296	0.13777	0.27320	0.25310	0.37626	0.10500	0.23415	0.45886	0.41006	0.20998	0.30337
	0.5	0.00103 0.02393	0.00175 0.03177		0.01025	0.08284	0.00171 0.02856		0.09879	0.22980	0.21280	0.34526	0.10362	0.23127	0.45891	0.41011	0.21322	0.30764
	Avg	0.00071 0.01935	0.00197 0.03333		0.01022	0.08266	0.00302	0.03682	0.16889	0.29821	0.27091	0.38609	0.10550	0.23501	0.45887	0.41007	0.20771	0.29969
Weather	0.125	0.00050 0.01259	0.00061 0.01446		0.00661	0.06123	0.00300	0.02110	0.09604	0.20783	0.36982	0.40486	0.00514	0.05275	0.40556	0.42631	0.13538	0.17599
	0.25	0.00067 0.01513	0.00073 0.01715		0.00657	0.06105	0.00214	0.01830	0.04910	0.14269	0.29296	0.36483	0.00476	0.05019	0.40558	0.42635	0.13688	0.18177
	0.375	0.00054 0.01443	0.00067 0.01700		0.00658	0.06113	0.00088	0.00924	0.04304	0.13516	0.17569	0.28913	0.00454	0.04811	0.40550	0.42633	0.13831	0.18700
	0.5	0.00031 0.01107	0.00047 0.01429		0.00650	0.06071	0.00042 0.00463		0.03787	0.12878	0.12578	0.24598	0.00492	0.04961	0.40551	0.42632	0.13850	0.19051
	Avg	0.00051 0.01331	0.00062 0.01573		0.00656	0.06103	0.00161	0.01332	0.05651	0.15362	0.24106	0.32620	0.00484	0.05016	0.40554	0.42633	0.13727	0.18382
1 st Count	24	21	2	1	0	0	0	2	0	0	0	0	4	6	0	0	0	0

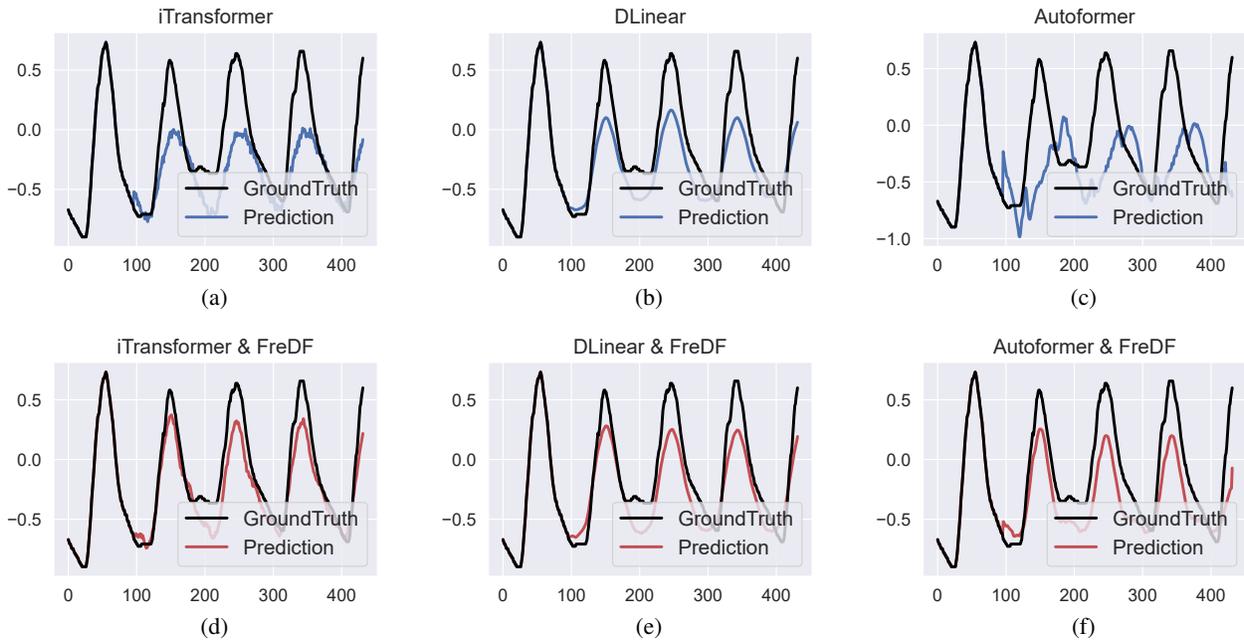


Figure 14: Forecast sequences generated by iTransformer, DLinear and Autoformer with and without FreDF. The prediction length is set to 336 and the experiment is conducted on a snapshot of ETTm2.

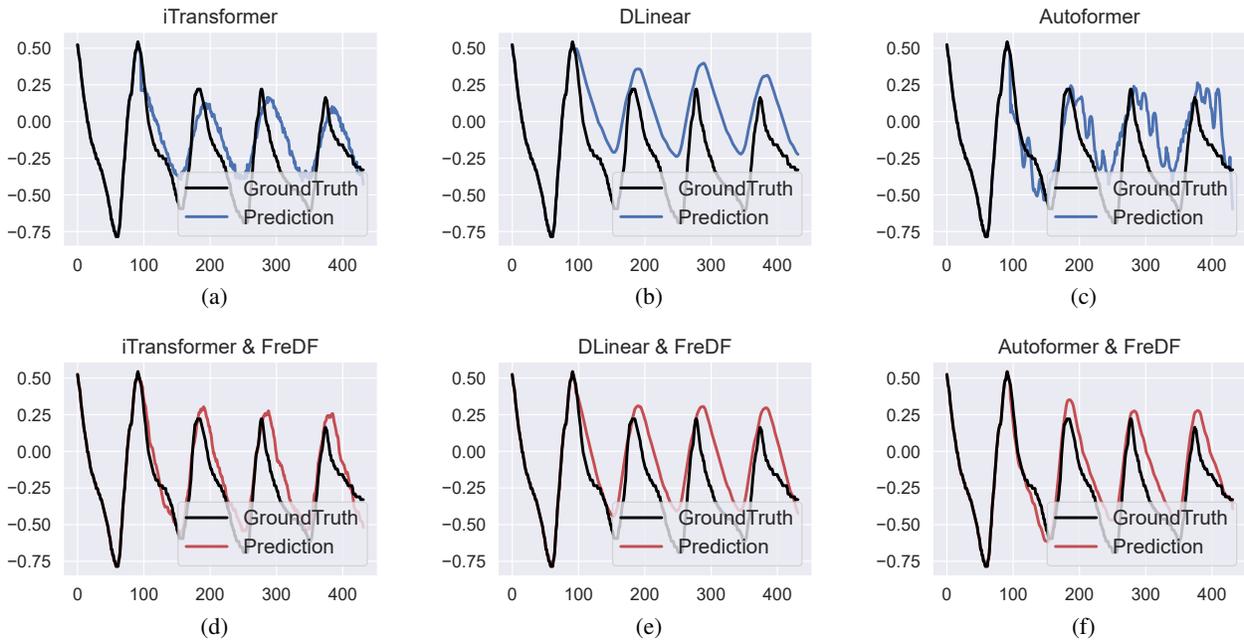


Figure 15: Forecast sequences generated by generated by iTransformer, DLinear and Autoformer with and without FreDF. The prediction length is set to 336 and the experiment is conducted on a snapshot of ETTm2.

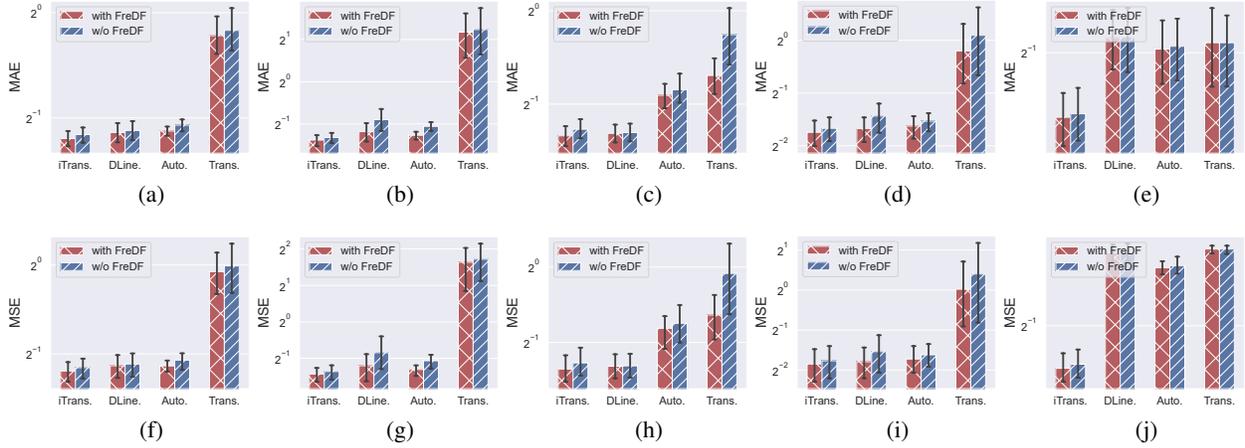


Figure 16: Performance of different forecast models with and without FreDF on the ETTh1 (a, f), ETTh2 (b, g), ETTm1 (c, h), ETTm2 (d, i), and Traffic (e, j) datasets. The forecast errors are averaged over prediction lengths and the error bars represent 95% confidence intervals.

solidifies FreDF’s position as a powerful tool capable of addressing a wide array of forecasting challenges, marking it as a significant contribution to the advancement of forecasting methodologies.

E.3 Hyperparameter Sensitivity

In this section, we investigate the influence of adjusting the frequency loss parameter, α , on the efficacy of the Frequency-enhanced Direct Forecast (FreDF) paradigm. This exploration is conducted across three models: iTransformer, Autoformer, and DLinear, with the respective results depicted in Figures 17, 18, and 19.

A consistent observation across these models is that incrementally increasing α from 0 to 1 generally leads to a decrease in forecast error, although a marginal increase in error is noted as α approaches 1. For example, within the ECL dataset for a prediction length of $T=192$, we witness a reduction in both Mean Absolute Error (MAE) and Mean Squared Error (MSE), from approximately 0.258 and 0.167 down to 0.247 and 0.158, respectively. This pattern of error reduction, observed across various prediction lengths and datasets, affirms the advantages of adopting a frequency domain learning approach.

Notably, the most significant decrease in forecast error often occurs at α values close to 1, such as 0.8 for the ETTh1 dataset, rather than at the maximum value of 1. This finding suggests that integrating supervisory signals from both the time and frequency domains can yield further enhancements in forecasting performance.

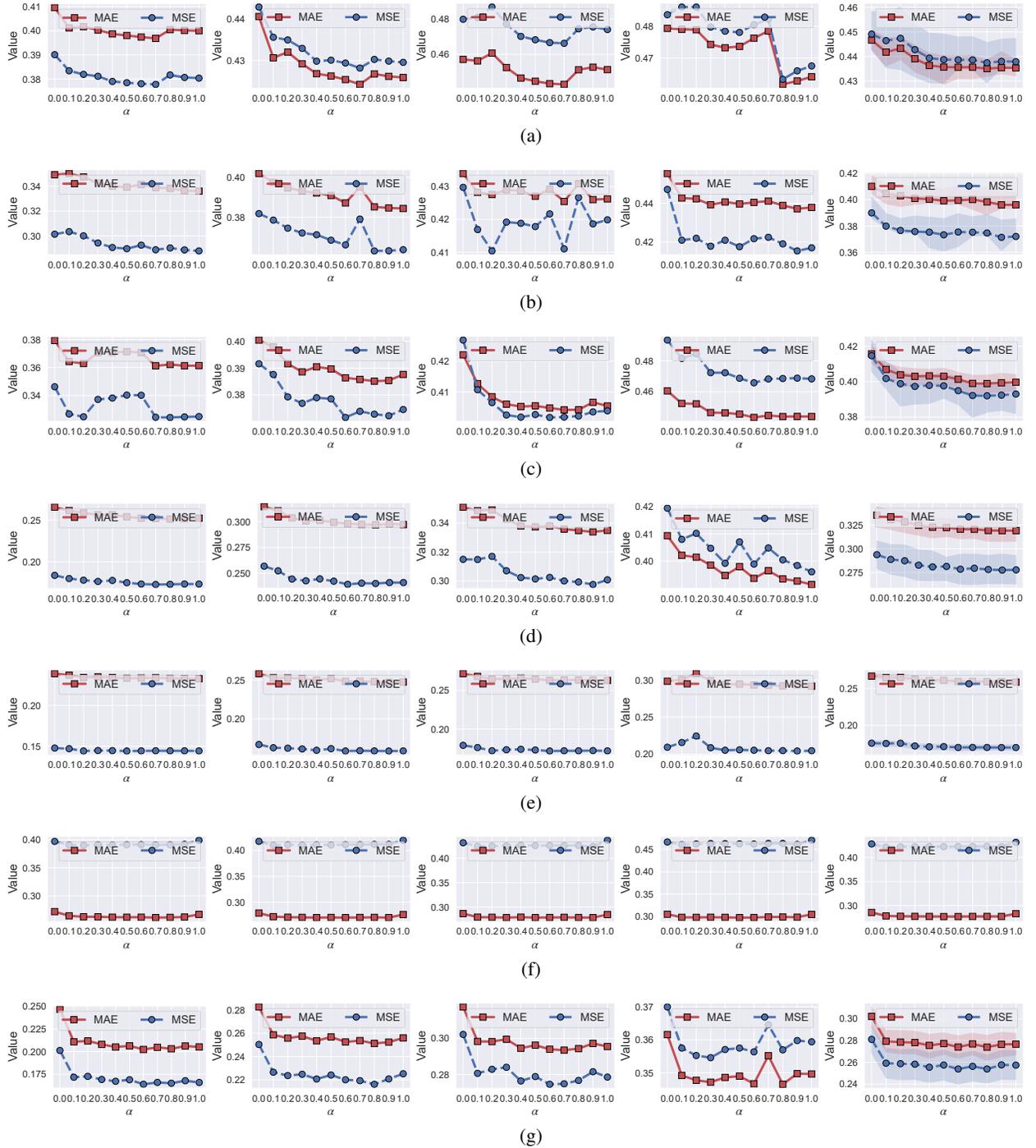


Figure 17: Performance of iTransformer enhanced by FreDF given different relative importance of frequency loss α . These experiments are conducted on ETTh1 (a), ETTh2 (b), ETTm1 (c), ETTm2 (d), ECL (e), Traffic (f) and Weather (g) datasets. Different columns correspond to different forecast lengths T (from left to right: 96, 192, 336, 720, and their average with shaded areas being 50% confidence intervals).

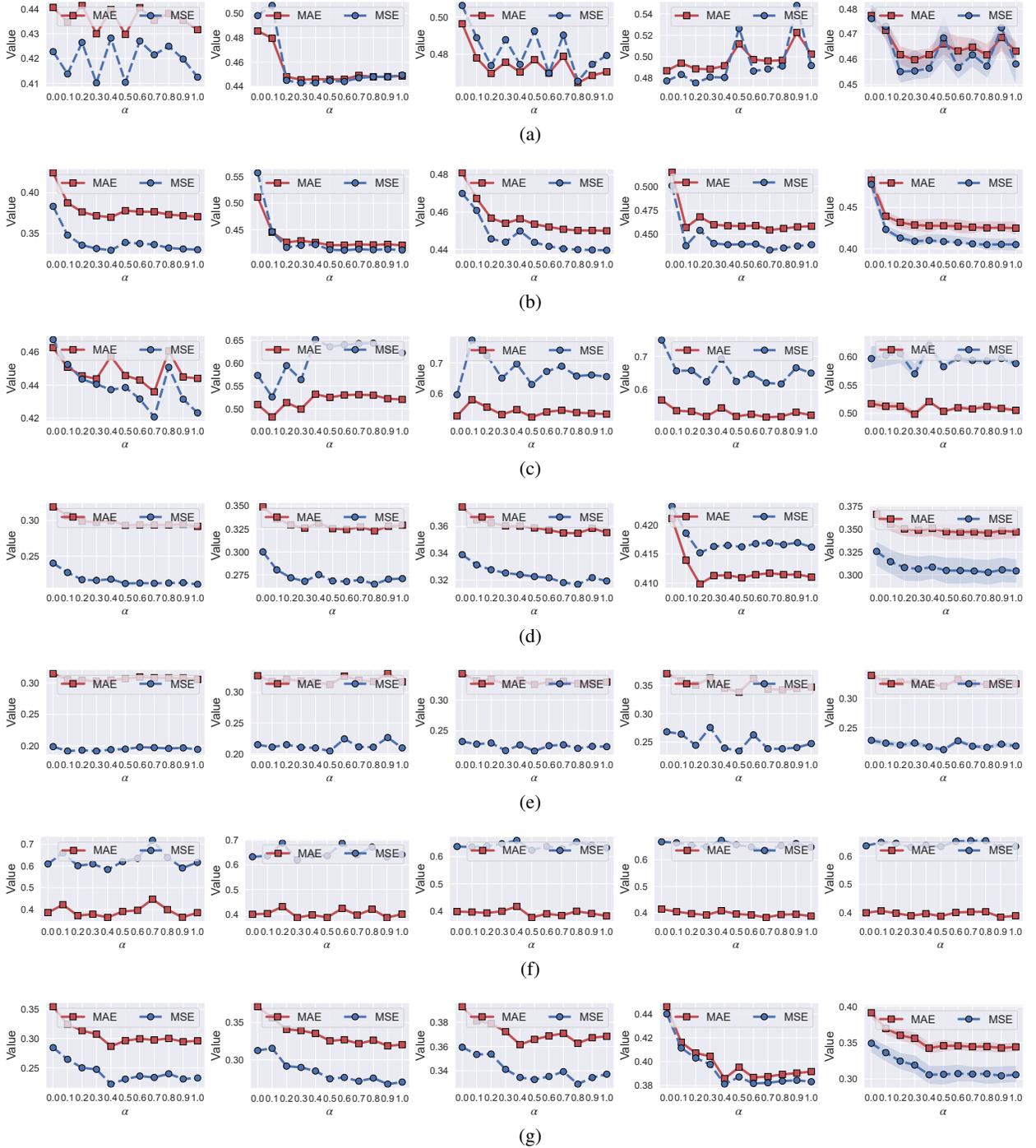


Figure 18: Performance of Autoformer enhanced by FreDF given different relative importance of frequency loss α . These experiments are conducted on ETTh1 (a), ETTh2 (b), ETTm1 (c), ETTm2 (d), ECL (e), Traffic (f) and Weather (g) datasets. Different columns correspond to different forecast lengths T (from left to right: 96, 192, 336, 720, and their average with shaded areas being 50% confidence intervals).

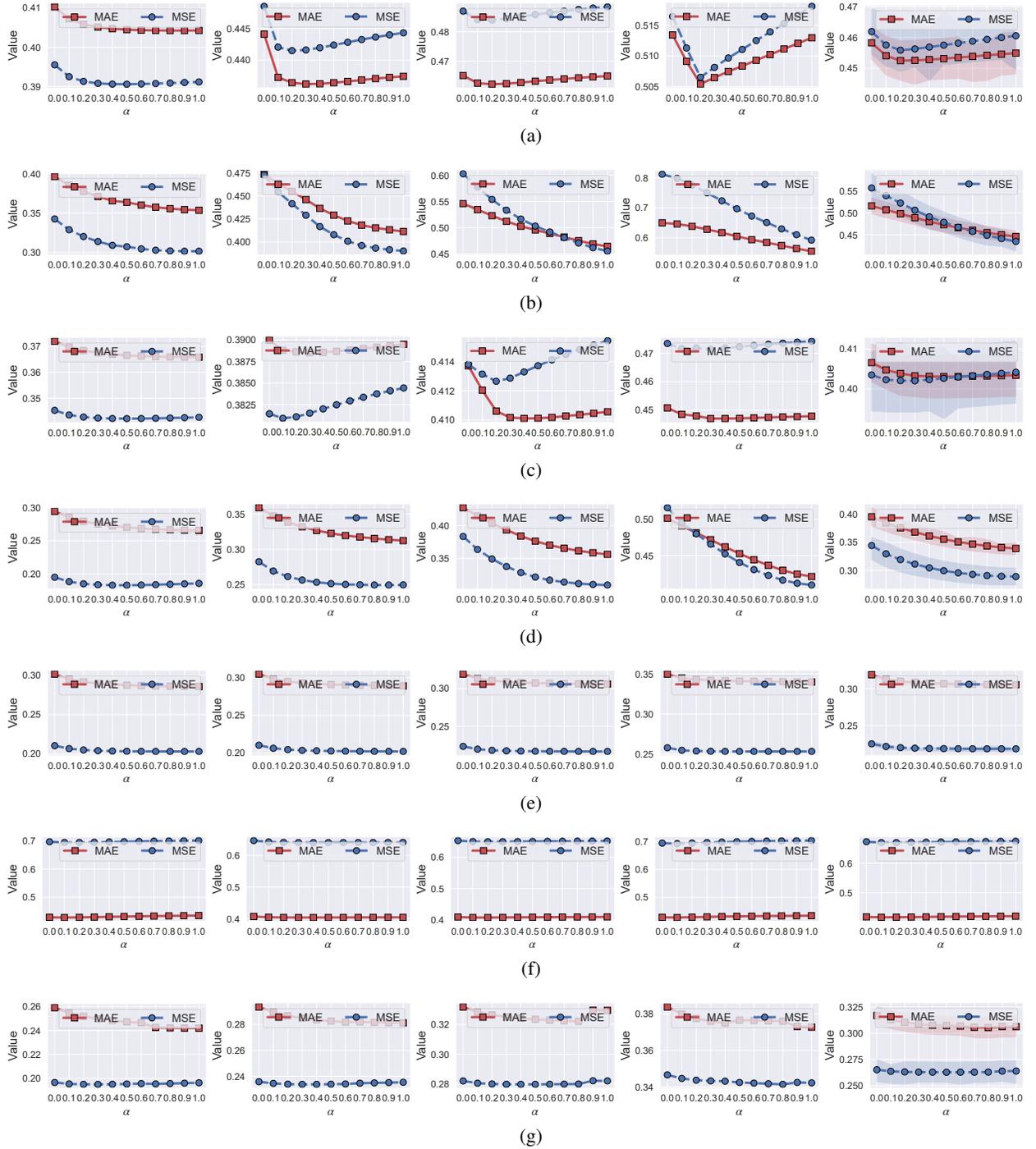


Figure 19: Performance of DLinear enhanced by FreDF given different relative importance of frequency loss α . These experiments are conducted on ETTh1 (a), ETTh2 (b), ETTm1 (c), ETTm2 (d), ECL (e), Traffic (f) and Weather (g) datasets. Different columns correspond to different forecast lengths T (from left to right: 96, 192, 336, 720, and their average with shaded areas being 50% confidence intervals).

Table 10: Full results of the system-level ablation studies. The forecast model is specified as iTransformer. $\mathcal{L}^{(\text{tmp})}$ and $\mathcal{L}^{(\text{freq})}$ indicates whether the temporal forecast loss and frequency forecast loss is incorporated in $\mathcal{L}^{(\text{freq})}$, respectively.

$\mathcal{L}^{(\text{tmp})}$	$\mathcal{L}^{(\text{freq})}$	Data	T=96		T=192		T=336		T=720		Avg	
			MSE	MAE								
✓	✗	ETTm1	0.346	0.379	0.391	0.400	0.426	0.422	0.493	0.460	0.414	0.415
		ETTh1	0.390	0.409	0.442	0.440	0.479	0.457	0.483	0.479	0.449	0.446
		ECL	0.147	0.239	0.166	0.258	0.178	0.271	0.209	0.298	0.175	0.266
		Traffic	<u>0.397</u>	0.271	<u>0.417</u>	0.278	<u>0.431</u>	0.286	<u>0.466</u>	0.305	<u>0.428</u>	0.285
		Weather	0.201	0.246	<u>0.250</u>	0.282	<u>0.302</u>	0.317	0.370	0.361	0.280	0.302
✗	✓	ETTm1	<u>0.324</u>	<u>0.361</u>	<u>0.374</u>	<u>0.387</u>	<u>0.403</u>	<u>0.405</u>	<u>0.468</u>	<u>0.443</u>	<u>0.392</u>	<u>0.399</u>
		ETTh1	<u>0.380</u>	<u>0.399</u>	<u>0.429</u>	<u>0.425</u>	<u>0.474</u>	<u>0.451</u>	<u>0.467</u>	<u>0.464</u>	<u>0.437</u>	<u>0.435</u>
		ECL	<u>0.144</u>	<u>0.232</u>	<u>0.158</u>	<u>0.247</u>	<u>0.171</u>	<u>0.262</u>	<u>0.204</u>	<u>0.291</u>	<u>0.169</u>	<u>0.258</u>
		Traffic	0.399	<u>0.267</u>	0.419	<u>0.276</u>	0.437	<u>0.284</u>	0.470	<u>0.304</u>	0.431	<u>0.283</u>
		Weather	<u>0.165</u>	<u>0.205</u>	<u>0.225</u>	<u>0.255</u>	<u>0.278</u>	<u>0.295</u>	<u>0.359</u>	<u>0.349</u>	<u>0.257</u>	<u>0.276</u>
✓	✓	ETTm1	0.324	0.362	0.372	0.385	0.402	0.404	0.468	0.443	0.391	0.398
		ETTh1	0.381	0.400	0.430	0.426	0.474	0.451	0.463	0.461	0.437	0.435
		ECL	0.144	0.233	0.158	0.247	0.172	0.263	0.204	0.293	0.169	0.259
		Traffic	0.390	0.265	0.410	0.272	0.424	0.280	0.460	0.298	0.421	0.279
		Weather	0.163	0.202	0.220	0.252	0.274	0.293	0.356	0.346	0.253	0.273