

# LABEL CORRELATION BIASES DIRECT TIME SERIES FORECAST

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Time series modeling presents unique challenges due to autocorrelation in both historical data and future sequences. While current research predominantly addresses autocorrelation within historical data, the correlations among future labels are often overlooked. Specifically, modern forecasting models primarily adhere to the direct forecast (DF) paradigm, generating multi-step forecasts independently and disregarding label correlations over time. In this work, we demonstrate that the learning objective of DF is biased in the presence of label correlation. To address this issue, we propose the Frequency-enhanced Direct Forecast (FreDF), which mitigates label correlation by learning to forecast in the frequency domain, thereby reducing estimation bias. Our experiments show that FreDF significantly outperforms existing state-of-the-art methods and is compatible with a variety of forecast models. Code is available at <https://anonymous.4open.science/r/FreDF-0FB1>.

## 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 (Bi et al., 2023), short-term prediction in industrial maintenance (Ma et al., 2023), and data imputation in healthcare (Si et al., 2020). A key challenge in time series modeling, distinguishing it from canonical regression tasks, is the presence of autocorrelation, which refers to the dependence between time steps inherent in *both* the input and label sequences.

To accommodate autocorrelation in input sequences, diverse forecast models have been developed, exemplified by recurrent (Salinas et al., 2020), convolution (Wu et al., 2023) and graph neural networks (Yi et al., 2023a). Recently, Transformer-based models, utilizing self-attention mechanisms to dynamically assess autocorrelation, have gained prominence (Liu et al., 2024; Nie et al., 2023). Concurrently, there is a growing trend of incorporating frequency analysis into forecast models. By representing input sequence in the frequency domain, input autocorrelations can be efficiently accommodated, which improves forecast performance of Transformers (Zhou et al., 2022), GNNs (Yi et al., 2023a) and MLPs (Yi et al., 2023b). These pioneering works highlight the importance of autocorrelation and frequency analysis in advanced time series modeling.

Another critical aspect is the autocorrelation within the label sequence, where each future step is autoregressively dependent on its predecessors. This phenomenon, known as *label autocorrelation*, poses significant challenges that have not been adequately addressed. Specifically, recent forecasting methods predominantly employ the direct forecast (DF) paradigm (Liu et al., 2024; Nie et al., 2023), which generates multi-step predictions simultaneously via a multi-output head (Liu et al., 2022b), optimizing forecast errors across all steps concurrently. However, this approach implicitly assumes *step-wise independence* in the label sequence, overlooking the inherent label autocorrelation present in time series data. We theoretically demonstrate that this oversight results in biased forecasts, revealing a significant flaw in the existing DF paradigm.

To address this issue, we introduce the *Frequency-enhanced Direct Forecast* (FreDF), a straightforward yet effective refinement of the DF paradigm. The central idea is to align the forecasts and label sequences in the frequency domain, where the label correlation is found to be effectively diminished. This method not only resolves the discrepancy between DF assumptions and the actual time-series

characteristics but also retains the advantages of the DF approach, such as sample efficiency and simplicity of implementation. Our main contributions are summarized as follows:

- We uncover label autocorrelation as a critical yet underexplored challenge in modern time series modeling and theoretically justify how it biases the learning objective of the prevalent DF paradigm.
- We propose FreDF, a straightforward yet effective modification to the DF paradigm that learns to forecast in the frequency domain, thereby mitigating label correlation and reducing bias. To our knowledge, this is the first work to leverage frequency analysis to enhance 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 an 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. Define input sequence  $L \in \mathbb{R}^{H \times D}$  and label sequence  $Y \in \mathbb{R}^{T \times D}$  where  $H$  and  $T$  are sequence lengths. At an arbitrary  $n$ -th step, these sequences are observed as  $L = [X(n - H + 1), \dots, X(n)]$  and  $Y = [X(n + 1), \dots, X(n + T)]$ . The goal of time series forecast is identifying a model  $g : \mathbb{R}^{H \times D} \rightarrow \mathbb{R}^{T \times D}$  within a model family  $\mathcal{G}$  (e.g., decision trees, neural networks) that generates the prediction sequence  $\hat{Y} = g(L)$  approximating the label sequence  $Y$ .

There are two critical aspects to accommodate autocorrelation in time series modeling: (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.

### 2.2 MODEL ARCHITECTURES

To exploit autocorrelation in the input sequences, diverse architectures have been developed. Initial statistical methods include VAR (Watson, 1993) and ARIMA (Asteriou & Hall, 2011). Subsequently, neural networks became increasingly prominent for their ability to automate feature interaction and capture nonlinear correlations. Exemplars include RNNs (e.g., DeepAR (Salinas et al., 2020), S4 (Gu et al., 2021)), CNNs (e.g., TimesNet (Wu et al., 2023)), and GNNs (e.g., MTGNN (Mateos et al., 2019)), each designed to effectively encode autocorrelation. Current progress reaches a debate between Transformer-based and MLP-based architectures, each with its advantages and limitations. Transformers (e.g., PatchTST (Nie et al., 2023), iTransformer (Liu et al., 2024), CrossFormer (Zhang & Yan, 2023)) excel in encoding autocorrelation but come with high computational costs, while MLPs (e.g., DLinear (Zeng et al., 2023), TSMixer (Ekambaram et al., 2023)) 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 (Zhou et al., 2022), 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 (Zhou et al., 2022; Wu et al., 2021), MLPs (Yi et al., 2023b) and GNNs (Yi et al., 2023a; Cao et al., 2020), which makes it a versatile plugin in the design of neural networks for time series forecast.

### 2.3 ITERATIVE FORECAST V.S. DIRECT FORECAST

There are two paradigms to generate multi-step forecast: iterative forecast (IF) and direct forecast (DF) (Liu et al., 2022b). The IF paradigm follows the canonical sequence-to-sequence manner, which forecasts one step at a time and uses previous predictions as input for subsequent forecasts. This recursive approach respects label autocorrelation in forecast generation, widely used by early-stage

methods (Lai et al., 2018; Salinas et al., 2020). However, IF suffers from high variance due to error propagation, which significantly impairs performance in long-term forecasts (Taieb & Atiya, 2015). Therefore, modern works Li et al. (2021) advocate the DF paradigm, which generates multi-step forecasts simultaneously using a multi-output head, featured by fast inference, implementation ease and superior accuracy. Currently, DF has been a dominant paradigm, continuing to be employed in modern works (Wu et al., 2023; Liu et al., 2024).

**Significance of this work.** Our work refines the DF paradigm by performing forecasting in the frequency domain<sup>1</sup>. In contrast to recent advancements that incorporate frequency analysis within model architectures to manage *input autocorrelation* (Yi et al., 2023b;a), our approach specifically focuses on refining the loss function to mitigate the bias caused by *label autocorrelation*, which is an unexplored yet significant aspect in modern time series analytics (Li et al., 2021)<sup>2</sup>.

### 3 PROPOSED METHOD

#### 3.1 MOTIVATION

Autocorrelation is a fundamental characteristic of time series data, where each observation is highly dependent on its predecessors (Zeng et al., 2023). This inherent dependency distinguishes time series from other data modalities and poses unique challenges for modeling. To capture autocorrelation, various neural network architectures have been developed (Wu et al., 2021; Liu et al., 2024). These architectures effectively model autocorrelation within the input sequence. However, they fall short when it comes to addressing autocorrelation in the label sequence—the future time steps we aim to predict. Handling autocorrelation in the label sequence is challenging, as it requires the learning objective to implicitly encapsulate label correlations, a task that cannot be achieved merely by modifying neural architectures.

Modern time series forecasting models are primarily trained under the multitask learning manner, known as the direct forecasting (DF) paradigm. Specifically, the DF paradigm employs a multi-output model  $g_\theta : \mathbb{R}^{H \times D} \rightarrow \mathbb{R}^{T \times D}$  to generate  $T$ -step forecasts  $\hat{Y} = g_\theta(L)$ . The model parameters  $\theta$  are optimized by minimizing the temporal loss:

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

In this learning objective, the temporal loss at each forecast step is computed independently, treating each future time step as a separate task. While this method has shown empirical effectiveness, it overlooks the autocorrelation present within the label sequence  $Y$ . Specifically, the label sequence is autoregressively generated, with  $Y_{t+1}$  being highly dependent on  $Y_t$ , as illustrated by the blue arrows in Figure 1(a). In contrast, the learning objective in (1) assumes that each step in the label sequence can be independently modeled, as indicated by the black arrows in Figure 1(a). This misalignment between the model’s assumptions and the data’s characteristics introduces bias into the learning objective of the DF paradigm, as demonstrated in Theorem 3.1.

**Theorem 3.1** (Bias of DF). *Given input sequence  $L$  and label sequence  $Y$ , the learning objective (1) of the DF paradigm is biased against the practical negative-log-likelihood (NLL), expressed as:*

$$\text{Bias} = \sum_{i=1}^T \frac{1}{2\sigma^2} (Y_i - \hat{Y}_i)^2 - \sum_{i=1}^T \frac{1}{2\sigma^2(1 - \rho_i^2)} \left( Y_i - \left( \hat{Y}_i + \sum_{j=1}^{i-1} \rho_{ij}(Y_j - \hat{Y}_j) \right) \right)^2, \quad (2)$$

where  $\hat{Y}_i$  indicates the prediction at the  $i$ -th step,  $\rho_{ij}$  denotes the partial correlation between  $Y_i$  and  $Y_j$  given  $L$ ,  $\rho_i^2 = \sum_{j=1}^{i-1} \rho_{ij}^2$ .

<sup>1</sup>Given the inferior performance of the IF paradigm (Li et al., 2021), this paper advocates adapting the DF paradigm to handle label autocorrelation, rather than revisiting IF to directly model label autocorrelation.

<sup>2</sup>Compared to prevailing works that applied FFT for image reconstruction (Jiang et al., 2021; Wang et al., 2023b; Xie et al., 2023), FreDF distinguishes itself in two aspects: (1) it innovatively applies FFT to enhance the DF paradigm, specifically tailored for time series analysis, and (2) it introduces a novel perspective on managing label autocorrelation, which elucidates the rationale of FFT to enhance forecasting performance.

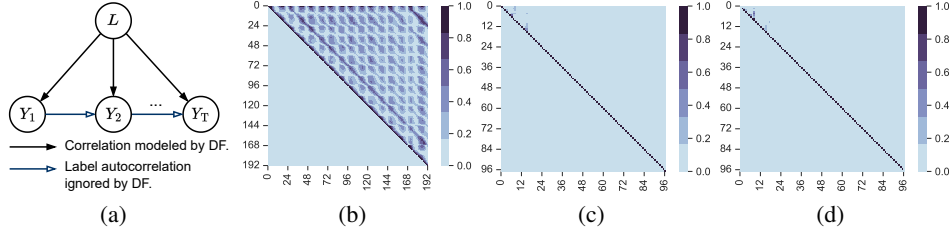


Figure 1: Visualizing label correlation in time series forecasting. (a) shows the generation process of time series with dependencies depicted as arrows. (b) shows the label correlation in the time domain, where each element  $\rho_{i,j}$  indicates the partial correlation between  $Y_i$  and  $Y_j$  given  $L$ . (c-d) shows the label correlation in the frequency domain, where each element  $\rho_{i,j}$  indicates the partial correlation between  $F_i$  and  $F_j$  given  $L$ , shown with the real (c) and imaginary part (d). Due to the symmetry inherent in FFT, the prediction length in the frequency domain is halved.

According to Theorem 3.1, the presence of label correlations  $\rho_{ij}$  causes the loss to be biased against the NLL of the real data. Notably, this bias diminishes to zero when the labels are uncorrelated ( $\rho_{ij} = 0$ ). Therefore, label correlation is a crucial aspect for training time series forecast models.

### 3.2 REDUCE LABEL CORRELATION WITH FOURIER TRANSFORM

As established in Theorem 3.1, the bias in the learning objective decreases as label correlations diminish. To achieve this reduction, a promising strategy is transforming the label sequence into a representation where these correlations are minimized. The Discrete Fourier Transform (DFT), defined in Definition 3.2, offers an intuitive and effective approach by projecting the sequence onto a set of orthogonal sine and cosine bases. In this transformed space, the label sequence is described as a linear combination of predefined temporal patterns that are orthogonal, which effectively bypasses the autocorrelation in the time domain. The efficacy of this transformation in reducing label correlation is formalized in Theorem 3.3, where different frequency components become decorrelated. Consequently, the reduced  $\rho_{i \neq j}$  lowers the bias against the NLL, which benefits the training of time series forecast models.

**Definition 3.2** (Discrete Fourier Transform, DFT). *The normalized DFT of a sequence  $Y = [Y_0, \dots, Y_{T-1}]$  is defined as its projection onto a set of orthogonal Fourier bases with different frequencies. The projection on the basis associated with frequency  $k$  is computed as*

$$F_k = \sum_{t=0}^{T-1} Y_t \exp \left( -j \left( \frac{2\pi k}{T} \right) t \right) / \sqrt{T},$$

where  $j$  is the imaginary unit,  $\exp(\cdot)$  is the Fourier basis orthogonal for different  $k$  values. DFT refers to the set of projections  $F = [F_1, \dots, F_{T-1}]$ , denoted as  $F = \mathcal{F}(Y)$ , which can be computed via the fast Fourier transform (FFT) algorithm with complexity  $\mathcal{O}(H \log H)$ .

**Theorem 3.3** (Decorrelation between frequency components). *Let  $Y$  be a zero-mean, discrete-time, wide-sense stationary random process of length  $T$ . As  $T \rightarrow \infty$ , the DFT coefficients become asymptotically uncorrelated at different frequencies:*

$$\lim_{T \rightarrow \infty} \mathbb{E}[F_k F_{k'}^*] = \begin{cases} S_Y(f_k), & \text{if } k = k', \\ 0, & \text{if } k \neq k', \end{cases}$$

where  $f_k = \frac{k}{T}$  and  $S_Y(f)$  is the power spectral density of  $Y$ .

**Case study.** To validate our theoretical claims, we conducted a case study on the Weather dataset, illustrated in Figure 1. The main observations are summarized as follows:

- **Evidence of Label Autocorrelation:** We quantified the partial correlations between different steps  $Y_i$  and  $Y_j$  of the label sequence  $Y$ , conditioned on the input  $L$ . The results revealed that a significant number of non-diagonal elements exhibit substantial values, with approximately 37.5% exceeding



0.3. This indicates that different time steps in  $Y$  are correlated even after accounting for  $L$ , confirming the presence of label autocorrelation. Moreover, the correlation patterns display regular variations, evidenced by alternating light and dark regions in the correlation matrix, suggesting a periodic nature in the series. Additional implementation details, empirical evidence, and formal analysis are provided in Appendix A. The existence of label autocorrelation contributes to the bias in the DF learning objective, as established in Theorem 3.1.

- **Effectiveness of Frequency Domain Transformation:** Figures 1 (c-d) visualize the partial correlations between different frequency components of the transformed label sequence  $F$ . The majority of non-diagonal elements show negligible values, with only about 3.6% exceeding 0.1. This demonstrates that transforming the label sequence to the frequency domain significantly reduces the partial correlations between different components, corroborating Theorem 3.3. The reduction in label correlations  $\rho_{i \neq j}$  leads to a decrease in the bias identified in Theorem 3.1, underscoring the potential of forecasting in the frequency domain for more accurate and unbiased predictions.

### 3.3 MODEL IMPLEMENTATION

In this section, we construct FreDF, a simple yet effective enhancement to the current DF training paradigm. The key is to align the forecasts and label sequences in the frequency domain to mitigate the bias caused by label autocorrelation.

According to the workflow in Figure 2, the historical sequence  $L$  is fed into the model to generate  $T$ -step forecasts, denoted as  $\hat{Y} = g(L)$ . The forecast error in the time domain,  $\mathcal{L}^{(\text{tmp})}$ , is computed according to (1). Subsequently, both the predicted and actual label sequences are transformed into the frequency domain using DFT, and the forecast error in the frequency domain is calculated as:

$$\mathcal{L}^{(\text{freq})} := \left| \mathcal{F}(\hat{Y}) - \mathcal{F}(Y) \right|, \quad (3)$$

where  $|\cdot|_1$  denotes the element-wise  $\ell_1$  norm, summing the absolute values of all elements in the matrix. Since FFT is differentiable (Wu et al., 2021; Zhou et al., 2022), the frequency-domain loss  $\mathcal{L}^{(\text{freq})}$  can be optimized using standard stochastic gradient descent methods. We advocate the use of the  $\ell_1$  loss in the frequency domain instead of the squared loss due to the numerical characteristics of the transformed label sequence. Specifically, different frequency components often exhibit vastly varying magnitudes; lower frequencies possess significantly higher amplitudes compared to higher frequencies, making the squared loss prone to instability. By employing the  $\ell_1$  loss, we ensure a more balanced and stable optimization process.

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:

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

Aligning forecasts and label sequence in the frequency domain, FreDF reduces the bias produced by label correlation while preserving the benefits of DF, such as efficient inference and multi-task learning capabilities. An important feature of FreDF is its model and transformation agnosticism. It is compatible with various forecasting models  $g$  (e.g., Transformers and MLPs). This flexibility significantly broadens the potential application scope of FreDF across different time series forecasting scenarios.

## 4 EXPERIMENTS

To demonstrate the efficacy of FreDF, six aspects are empirically investigated:

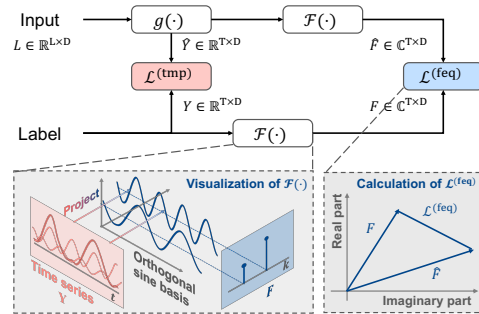


Figure 2: The workflow of FreDF. Key operations in the time and frequency domains are highlighted in red and blue, respectively.

Table 1: Long-term forecasting performance.

Models	FreDF (Ours)		iTransformer (2024)		FreTS (2023)		TimesNet (2023)		MICN (2023)		TiDE (2023)		DLinear (2023)		FEDformer (2022)		Autoformer (2021)		Transformer (2017)		TCN (2017)	
Metrics	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	<b>0.392</b>	<b>0.399</b>	0.415	0.416	0.407	0.415	0.413	0.418	<u>0.399</u>	0.423	0.419	0.419	0.404	<u>0.407</u>	0.440	0.451	0.596	0.517	0.943	0.733	0.891	0.632
ETTm2	<b>0.278</b>	<b>0.319</b>	<u>0.294</u>	0.335	0.335	0.379	0.297	<u>0.332</u>	0.300	0.356	0.358	0.404	0.344	0.396	0.302	0.348	0.326	0.366	1.322	0.814	3.411	1.432
ETTh1	<b>0.437</b>	<b>0.435</b>	0.449	<u>0.447</u>	0.488	0.474	0.478	0.466	0.525	0.515	0.628	0.574	0.462	0.458	<u>0.441</u>	0.457	0.476	0.477	0.993	0.788	0.763	0.636
ETTh2	<b>0.371</b>	<b>0.396</b>	<u>0.390</u>	<u>0.410</u>	0.550	0.515	0.413	0.426	0.624	0.549	0.611	0.550	0.558	0.516	0.430	0.447	0.478	0.483	3.296	1.419	3.325	1.445
ECL	<b>0.170</b>	<b>0.259</b>	<u>0.176</u>	<u>0.267</u>	0.209	0.297	0.214	0.307	0.187	0.297	0.251	0.344	0.225	0.319	0.229	0.339	0.228	0.339	0.274	0.367	0.617	0.598
Traffic	<b>0.421</b>	<b>0.279</b>	<u>0.428</u>	<u>0.286</u>	0.552	0.348	0.535	0.309	0.636	0.335	0.760	0.473	0.673	0.419	0.611	0.379	0.637	0.399	0.680	0.376	1.001	0.652
Weather	<b>0.254</b>	<b>0.274</b>	0.281	0.302	<u>0.255</u>	0.299	0.262	<u>0.288</u>	0.261	0.319	0.271	0.320	0.265	0.317	0.311	0.361	0.349	0.391	0.632	0.552	0.584	0.572
PEMS03	<u>0.113</u>	<u>0.219</u>	0.116	0.226	0.146	0.257	0.118	0.223	<b>0.099</b>	<b>0.214</b>	0.316	0.370	0.233	0.344	0.174	0.302	0.501	0.513	0.126	0.233	0.666	0.634
PEMS08	<b>0.141</b>	<b>0.238</b>	0.159	0.258	0.174	0.277	<u>0.154</u>	<u>0.245</u>	0.717	0.459	0.319	0.378	0.294	0.377	0.232	0.322	0.630	0.572	0.249	0.266	0.713	0.629

Note: We fix the input length as 96 following the established benchmarks (Liu et al., 2024; Wu et al., 2023). **Bold** typeface highlights the top performance for each metric, while underlined text denotes the second-best results. The results are averaged over prediction lengths (96, 192, 336 and 720), with full results in Table 5.

- Performance:** *Does FreDF work?* Section 4.2 compares FreDF against state-of-the-art baselines using public datasets. The long-term forecasting task is investigated in Section 4.2 and the short-term forecasting and imputation tasks are explored in Appendix E.1.
- Mechanism:** *How does it work?* Section 4.3 offers an ablative study to dissect the the contributions of FreDF’s individual components, elucidating their roles in enhancing forecasting accuracy.
- Generality:** *Does it support other forecasting models?* Section 4.4 verifies the adaptability of FreDF across different forecasting models, with additional results documented in Appendix E.5.
- Flexibility:** *Does it support alternatives to FFT?* Section 4.4 replaces FFT with other transformations to showcase its flexibility of implementation.
- Sensitivity:** *Does it necessitate careful finetuning?* Section 4.5 presents a sensitivity analysis of the hyperparameter  $\alpha$ , where FreDF maintains efficacy across a broad range of parameter values.
- Efficiency:** *Is FreDF effective given limited samples?* Section 4.6 offers a learning curve analysis, where FreDF achieves comparable performance with limited samples to that obtained using substantially more time-domain labels, indicating an advantageous sample efficiency.

#### 4.1 SETUP

**Datasets.** The datasets for long-term forecast and imputation include ETT (4 subsets), ECL, Traffic, Weather and PEMS following Wu et al. (2021) and Liu et al. (2024). The dataset for short-term forecast is M4 following Wu et al. (2023). 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 (Vaswani et al., 2017), Autoformer (Wu et al., 2021), FEDformer (Zhou et al., 2022), iTransformer (Liu et al., 2024); (2) MLP-based methods: DLinear (Zeng et al., 2023), TiDE (Das et al., 2023), FreTS (Yi et al., 2023b); (3) other notable models: TimesNet (Wu et al., 2023), MICN (Wang et al., 2023a), TCN (Bai et al., 2018). Notably, iTransformer (Liu et al., 2024) is the state-of-the-art baseline released in ICLR-24.

**Implementation.** The baseline models are reproduced using the scripts sourced from TimesNet (Wu et al., 2023). They are trained with Adam (Kingma & Ba, 2015) optimizer to minimize the MSE loss. When integrating FreDF to enhance an established model, we respect the associated hyperparameter settings in the public benchmark (Wu et al., 2023), merely tuning  $\alpha$  and learning rate conservatively. Experiments are conducted on Intel(R) Xeon(R) Platinum 8383C CPUs NVIDIA RTX 3090 GPUs. More implementation details are provided in Appendix D.2.

#### 4.2 OVERALL PERFORMANCE

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. Overall, FreDF improves the performance

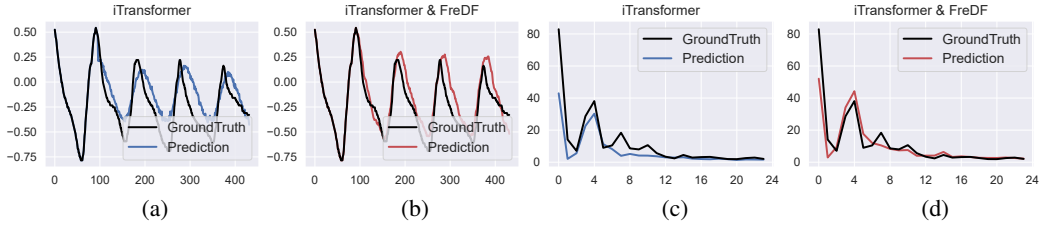


Figure 3: Visualization of forecast sequence generated with and without FreDF in the time (a-b) and frequency (c-d) domain.

Table 2: Ablation study results.

Model	$\mathcal{L}^{(tmp)}$	$\mathcal{L}^{(freq)}$	Data	T=96		T=192		T=336		T=720		Avg	
				MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
DF	✓	✗	ETTh1	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
			Weather	0.201	0.246	0.250	0.282	0.302	0.317	0.370	0.361	0.280	0.302
FreDF <sup>†</sup>	✗	✓	ETTh1	<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>
			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>
FreDF	✓	✓	ETTh1	<b>0.324</b>	<b>0.362</b>	<b>0.372</b>	<b>0.385</b>	<b>0.402</b>	<b>0.404</b>	<b>0.468</b>	<b>0.443</b>	<b>0.391</b>	<b>0.398</b>
			ETTh1	<b>0.381</b>	<b>0.400</b>	<b>0.430</b>	<b>0.426</b>	<b>0.474</b>	<b>0.451</b>	<b>0.463</b>	<b>0.461</b>	<b>0.437</b>	<b>0.435</b>
			ECL	<b>0.144</b>	<b>0.233</b>	<b>0.158</b>	<b>0.247</b>	<b>0.172</b>	<b>0.263</b>	<b>0.204</b>	<b>0.293</b>	<b>0.169</b>	<b>0.259</b>
			Weather	<b>0.163</b>	<b>0.202</b>	<b>0.220</b>	<b>0.252</b>	<b>0.274</b>	<b>0.293</b>	<b>0.356</b>	<b>0.346</b>	<b>0.253</b>	<b>0.273</b>

of iTransformer substantially. For instance, on the ETTh1 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.

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 handling label autocorrelation and FreDF.

**Showcases.** We visualize the forecast sequences to highlight the improvements of FreDF in forecast quality. A ETTh2 snapshot with T=336 is depicted in Figure 3. 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 sequential 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.

Table 3: Varying FFT implementation results.

Model	ETTh1				ETTh1				ECL			
	MSE	$\Delta$	MAE	$\Delta$	MSE	$\Delta$	MAE	$\Delta$	MSE	$\Delta$	MAE	$\Delta$
iTransformer	0.449	-	0.447	-	0.415	-	0.416	-	0.176	-	0.267	-
+ FreDF-T	0.437	↓ 2.63%	0.435	↓ 2.62%	0.392	↓ 5.49%	0.399	↓ 4.01%	0.170	↓ 3.41%	0.259	↓ 2.77%
+ FreDF-D	0.445	↓ 0.92%	0.440	↓ 1.42%	0.395	↓ 4.77%	0.398	↓ 4.33%	0.171	↓ 2.51%	0.260	↓ 2.52%
+ FreDF-2	0.432	↓ 3.94%	0.431	↓ 3.57%	0.392	↓ 5.60%	0.399	↓ 4.05%	0.166	↓ 5.32%	0.256	↓ 4.20%

Note:  $\Delta$  denotes the relative error reduction compared to iTransformer with DF paradigm.

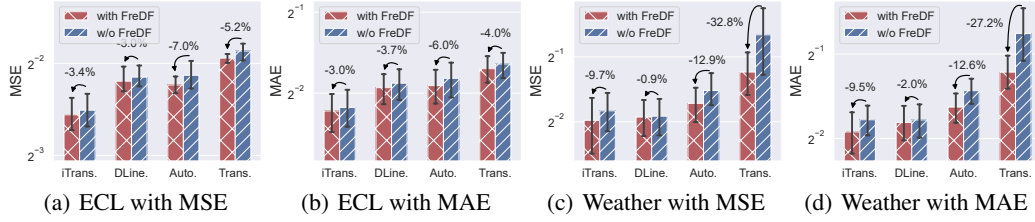


Figure 4: Benefit of incorporating FreDF in varying models, shown with colored bars for means over prediction lengths (96, 192, 336, 720) and error bars for 99.9% confidence intervals.

### 4.3 ABLATION STUDIES

In this section, we dissect the contributions of the temporal and frequency loss for enhancing forecast performance. The results are detailed in Table 2, where iTransformer is employed as the forecast model. Overall, the frequency loss consistently improves performance compared to the temporal loss. The rationale is that label autocorrelation can be effectively managed in the frequency domain, aligning better with the conditional independence assumption inherent in DF. Moreover, learning to forecast in both domains generally showcase improvement compared to relying solely on one domain. However, the improvement over  $\mathcal{L}^{(\text{freq})}$  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.

### 4.4 GENERALIZATION STUDIES

In this section, we investigate the utility of FreDF with different forecast models and domain transformation strategies, to showcase the generality of FreDF. In the bar-plots, the forecast errors are averaged over prediction lengths (96, 192, 336, 720), with error bars as 95% confidence intervals.

**Varying forecast models.** We explore the versatility of FreDF in augmenting representative neural forecasting models: iTransformer, DLinear, Autoformer, and Transformer. FreDF demonstrates significant enhancements across these models compared to the traditional DF paradigm, as illustrated in Figure 4. Notably, Transformer-based models such as the Autoformer and Transformer substantially benefit from the integration of FreDF. On the ECL dataset, for instance, the Autoformer (developed in 2021) enhanced by FreDF outperforms DLinear (developed in 2023). More evidence of FreDF’s versatility is provided in Appendix E. These results confirm FreDF’s potential as a plugin-and-play strategy to enhance various time-series forecasting models.

**Varying FFT implementations.** We note that label correlation exists between not only different steps, but also variables in multivariate forecasting. Therefore, we implement FFT along the time (FreDF-T) and variable dimension (FreDF-D) to handle the corresponding correlations, with the outcomes illustrated in Table 3. 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, FreDF-T slightly outperforms FreDF-D, 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 axes (FreDF-2), which accommodates the correlations between both time steps and variables simultaneously and further improves performance.

**Varying transformations.** Motivated by the fact that FFT can be viewed as projections onto sine polynomials, we extend the implementation of FreDF by replacing FFT with projections onto other established polynomials. Each polynomial set is adept at capturing specific data patterns, such as trends and periodicity, which are challenging to learn in the time domain. The results are summarized in Figure 5. Notably, projections onto Legendre and Fourier bases demonstrate superior performance. This superiority is attributed to the orthogonality between polynomials, a feature not guaranteed by others as analyzed in Appendix C. It underscores orthogonality when selecting polynomials for implementing FreDF, which is pivotal for eliminating autocorrelations.

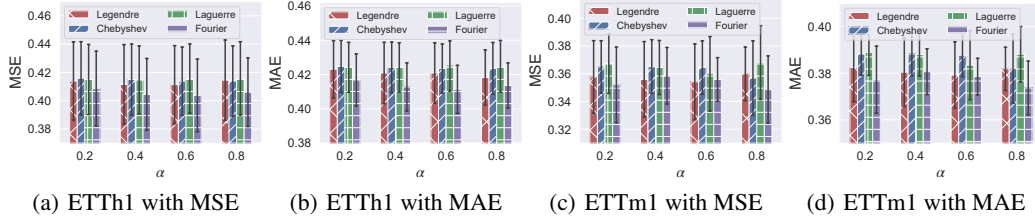


Figure 5: Varying projection bases results, shown with colored bars for means over prediction lengths (96, 192, 336, 720) and error bars for 99.9% confidence intervals.

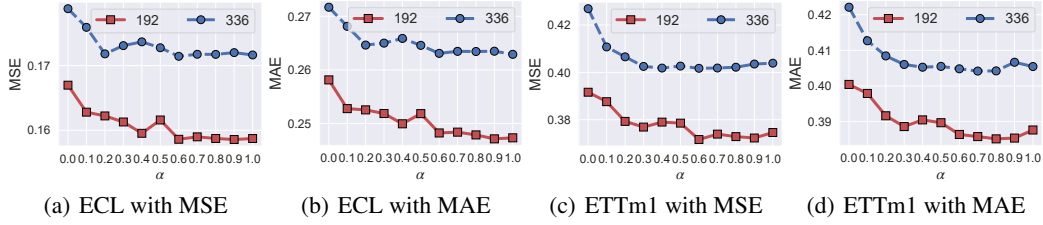


Figure 6: Varying strength of frequency loss ( $\alpha$ ) results, shown with colored lines for  $T=192, 336$ .

#### 4.5 HYPERPARAMETER SENSITIVITY

The key hyperparameter of FreDF is the frequency loss strength  $\alpha$ . The performance given different  $\alpha$  is summarized in Figure 6. Overall, increasing  $\alpha$  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. Notably, the optimal reduction in forecast error typically occurs at  $\alpha$  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. The claims above can be supported by more

#### 4.6 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 7. 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 learning in the time domain 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.

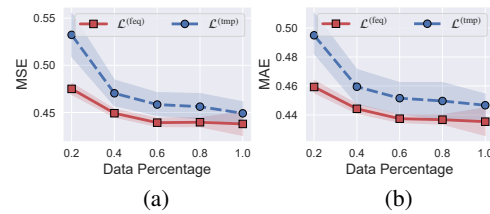


Figure 7: Learning curve on ETTm1 dataset.



## 5 CONCLUSION

In this study, we investigated the challenge of label correlation in time series modeling, which biases the learning objective of the DF paradigm away from the true likelihood of time series data. To mitigate this issue, we developed FreDF, which reduces label correlation by transforming the label sequence into the frequency domain, thereby diminishing the bias in the learning objective. Our experiments demonstrate that FreDF not only enhances forecasting accuracy but also exhibits strong adaptability across various tasks and forecasting 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 independent component analysis can produce orthogonal bases considering 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.

## REFERENCES

- Dimitros Asteriou and Stephen G Hall. Arima models and the box-jenkins methodology. *Appl. Econ.*, 2(2):265–286, 2011.
- 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.
- 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.
- Michela Bia, Martin Huber, and Lukáš Lafférs. Double machine learning for sample selection models. *Journal of Business & Economic Statistics*, 42(3):958–969, 2024.
- 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 *Proc. Adv. Neural Inf. Process. Syst.*, volume 33, pp. 17766–17778, 2020.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. Double/debiased machine learning for treatment and structural parameters: Double/debiased machine learning. *The Econometrics Journal*, 21(1), 2018.
- 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.
- Vijay Ekambaram, Arindam Jati, Nam Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. Tsmixer: Lightweight mlp-mixer model for multivariate time series forecasting. In *Proc. ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, pp. 459–469, 2023.
- Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *Proc. Int. Conf. Learn. Represent.*, 2021.
- Liming Jiang, Bo Dai, Wayne Wu, and Chen Change Loy. Focal frequency loss for image reconstruction and synthesis. In *Proc. IEEE Int. Conf. Comput. Vis.*, pp. 13899–13909. IEEE, 2021.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proc. Int. Conf. Learn. Represent.*, 2015.
- Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *SIGIR*, 2018.
- Jianxin Li, Xiong Hui, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proc. AAAI Conf. Artif. Intell.*, 2021.

- 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 *Proc. Adv. Neural Inf. Process. Syst.*, 2022a.
- Shiyu Liu, Rohan Ghosh, and Mehul Motani. Towards better long-range time series forecasting using generative forecasting. *CoRR*, abs/2212.06142, 2022b.
- 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 *Proc. Int. Conf. Learn. Represent.*, 2024.
- 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.
- 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.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *Proc. Int. Conf. Learn. Represent.*, 2023.
- David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *Int. J. Forecast.*, 36(3):1181–1191, 2020.
- 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.
- 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.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proc. Adv. Neural Inf. Process. Syst.*, 2017.
- Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. Micn: Multi-scale local and global context modeling for long-term series forecasting. In *Proc. Int. Conf. Learn. Represent.*, 2023a.
- Wenxuan Wang, Jing Wang, Chen Chen, Jianbo Jiao, Lichao Sun, Yuanxiu Cai, Shanshan Song, and Jiangyun Li. Fremae: Fourier transform meets masked autoencoders for medical image segmentation. *CoRR*, abs/2304.10864, 2023b.
- Mark W. Watson. Vector autoregressions and cointegration. *Working Paper Series, Macroeconomic Issues*, 4, 1993.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with Auto-Correlation for long-term series forecasting. In *Proc. Adv. Neural Inf. Process. Syst.*, 2021.
- Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *Proc. Int. Conf. Learn. Represent.*, 2023.
- Jiahao Xie, Wei Li, Xiaohang Zhan, Ziwei Liu, Yew-Soon Ong, and Chen Change Loy. Masked frequency modeling for self-supervised visual pre-training. In *Proc. Int. Conf. Learn. Represent. OpenReview.net*, 2023.
- Kun Yi, Qi Zhang, Wei Fan, Hui He, Liang Hu, Pengyang Wang, Ning An, Longbing Cao, and Zhendong Niu. Fourierrnn: Rethinking multivariate time series forecasting from a pure graph perspective. In *Proc. Adv. Neural Inf. Process. Syst.*, 2023a.
- 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 *Proc. Adv. Neural Inf. Process. Syst.*, 2023b.

Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proc. AAAI Conf. Artif. Intell.*, 2023.

Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *Proc. Int. Conf. Learn. Represent.*, 2023.

Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *Proc. Int. Conf. Mach. Learn.*, 2022.

## CONTENTS

<b>A Overview of DML for Partial Correlation Estimation</b>	<b>14</b>
A.1 Motivation . . . . .	14
A.2 Method . . . . .	14
A.3 Experimental Settings . . . . .	14
A.4 More Experimental Results . . . . .	15
<b>B Theoretical Justification</b>	<b>15</b>
<b>C Generalized Transformation onto Different Bases</b>	<b>20</b>
<b>D Reproduction Details</b>	<b>22</b>
D.1 Dataset Descriptions . . . . .	22
D.2 Implementation Details . . . . .	22
<b>E More Experimental Results</b>	<b>23</b>
E.1 Overall Performance . . . . .	23
E.2 Running cost analysis . . . . .	25
E.3 Random Seed Sensitivity . . . . .	29
E.4 Amplitude v.s. Phase Alignment . . . . .	29
E.5 Generalization Studies . . . . .	29
E.6 Hyperparameter Sensitivity . . . . .	30
<b>F Broader Impact</b>	<b>30</b>

## A OVERVIEW OF DML FOR PARTIAL CORRELATION ESTIMATION

### A.1 MOTIVATION

In this section, we introduce the rationale for employing double machine learning (DML) to quantify the partial correlations. Our focus is on the autocorrelation represented by  $Y_t \rightarrow Y_{t'}$  where  $0 \leq t < t' < T$ . However, the fork structure  $Y_t \leftarrow L(n) \rightarrow Y_{t'}$  creates a pseudo correlation between  $Y_{t'}$  and  $Y_t$ . In this case, the autocorrelation  $Y_t \rightarrow Y_{t'}$  is influenced 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 this influence and quantify partial correlation, it is essential to employ methods that excel in distinguishing direct relationships from spurious ones. DML is chosen for calculating partial correlation for three key reasons (Bia et al., 2024; Chernozhukov et al., 2018): (1) its model-agnostic nature, which does not depend on specific machine learning model specifications; (2) 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 partial correlation. 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 control variable that needs to be accounted for. The implementation of DML is depicted in Figure 8 (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 control variables ( $\mathcal{X}$ ). To this end, we first use two machine learning models, namely  $\phi$  and  $\psi$ , to predict the outcome and the treatment based on  $\mathcal{X}$ . These predictions aim to capture the components in  $\mathcal{Y}$  and  $\mathcal{T}$  that are influenced by  $\mathcal{X}$ . 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  $\epsilon$  is the error term;  $\beta$  is the model coefficient that can be identified via ordinary least squares. The  $\beta$  can be identified in a supervised learning manner, with objective to minimize the MSE of the prediction and real values. The identified  $\beta$  quantifies the partial correlation between the treatment and the outcome, having accounted for the influence of  $\mathcal{X}$ .

By regressing the orthogonalized outcome on the orthogonalized treatment, DML captures the direct effect of the treatment on the outcome without the interference from control variables, as depicted in Figure 8 (c). That is, DML isolates the desired partial correlation  $\mathcal{T} \rightarrow \mathcal{Y}$  from the influencing 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  $\phi$  and  $\psi$ , we opt for a linear regression model optimized using ordinary least squares for its efficiency<sup>3</sup>. Following Appendix A.1, we treat the history sequence

<sup>3</sup>The 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.



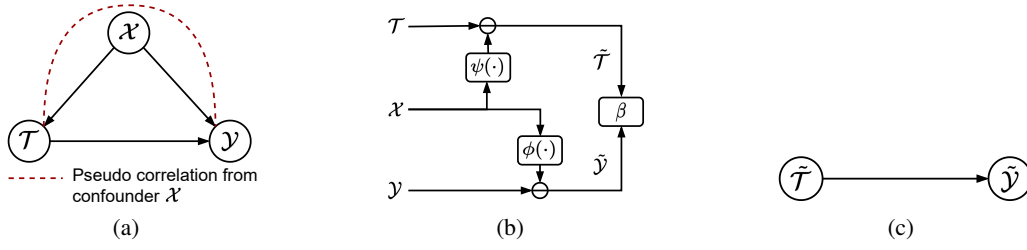


Figure 8: Visualization of partial correlation and DML approach for partial correlation quantification. (a) The correlation graph where the pseudo correlation is caused by the fork structure  $\mathcal{T} \leftarrow \mathcal{X} \rightarrow \mathcal{Y}$ . (b) The implementation of DML, where  $\beta$  is the identified strength of the partial correlation  $\mathcal{T} \rightarrow \mathcal{Y}$ . (c) The partial correlation identified by DML.

$L$  as the control variable 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 dataset<sup>4</sup>. This focus makes  $Y$  a scalar value within the real number space rather than a D-dimensional vector in this experiment.

**Specifications for identifying time-domain partial correlation.** To assess the partial correlation  $Y_t \rightarrow Y_{t'}$ , 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  $\beta$  derived from the DML model is interpreted as the strength of the partial correlation  $Y_t \rightarrow Y_{t'}$ .

**Specifications for identifying frequency-domain partial correlation.** To quantify the partial correlation  $F_k \rightarrow F_{k'}$ , 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  $\beta$  derived from the DML model is interpreted as the strength of the partial correlation  $F_k \rightarrow F_{k'}$ . A notable complexity arises due to  $F_k$  being a complex number. Since DML and similar analytical 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 partial correlation 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 9, with prediction length set to 192. Subsequently, we present the results given varying prediction lengths: 48, 96, 192, 336 in Figure 10, 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 9 (a-c) and Figure 10 (a-d) demonstrate significant values, which affirms the presence of label autocorrelation in the time domain. In contrast, the non-diagonal elements in Figure 9 (d-i) and Figure 10 (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** (Bias of DF, simplified). *Given a input sequence  $L$  and a univariate label sequence  $Y = [Y_1, Y_2]$  (the prediction length is set to 2 for simplicity), the learning objective (1) of DF*

<sup>4</sup>This focus is aligned with the study's objective of analyzing autocorrelations instead of inter-feature correlations, which simplifies the interpretation of results.

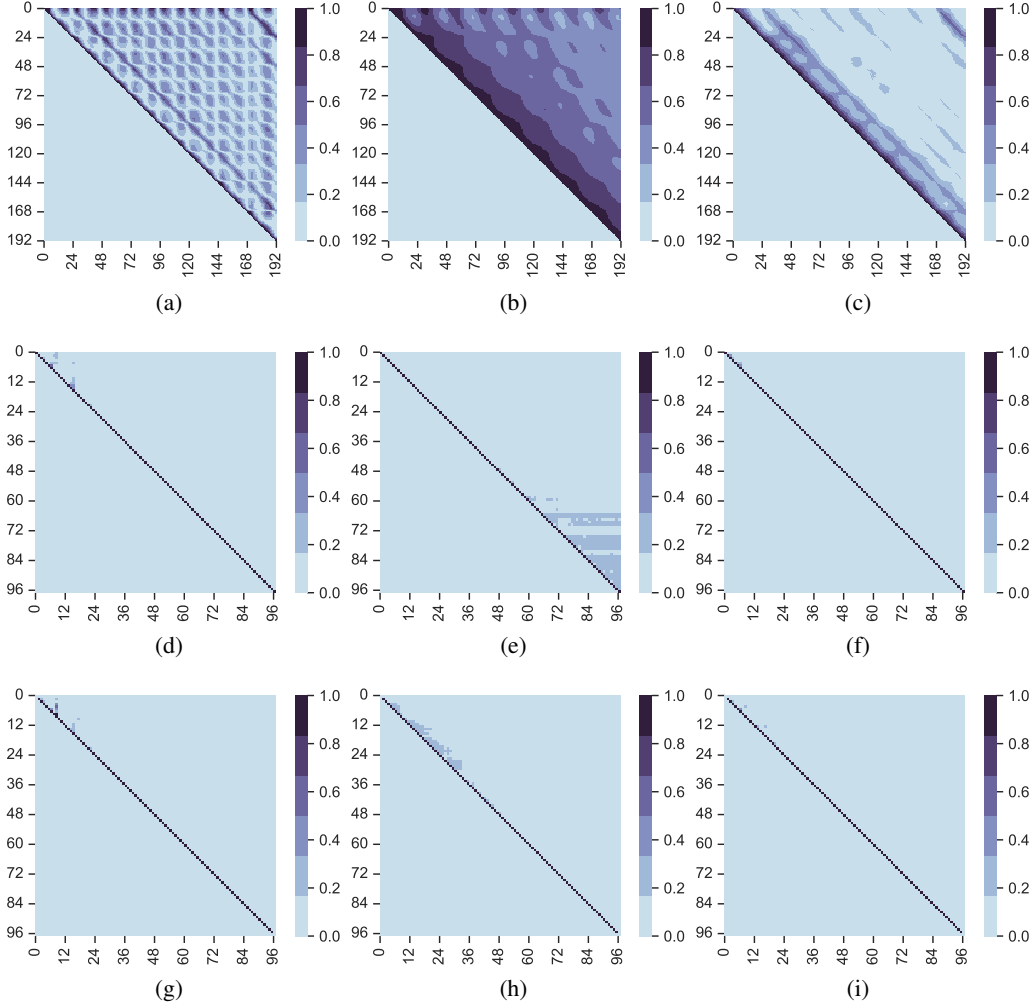


Figure 9: More comprehensive visualizations of label autocorrelation in different domains and datasets, with columns representing different datasets: Traffic, ETTh1, and ECL, from left to right. Panels (a-c) show the label correlation in the time domain, where each element  $\rho_{i,j}$  indicates the partial correlation between  $Y_i$  and  $Y_j$  given  $L$ . Panels (d-i) show the label correlation in the frequency domain, where each element  $\rho_{i,j}$  indicates the partial correlation between  $F_i$  and  $F_j$  given  $L$ , shown with the real (d-f) and imaginary part (g-i). Due to the symmetry inherent in FFT, the prediction length in the frequency domain is halved.

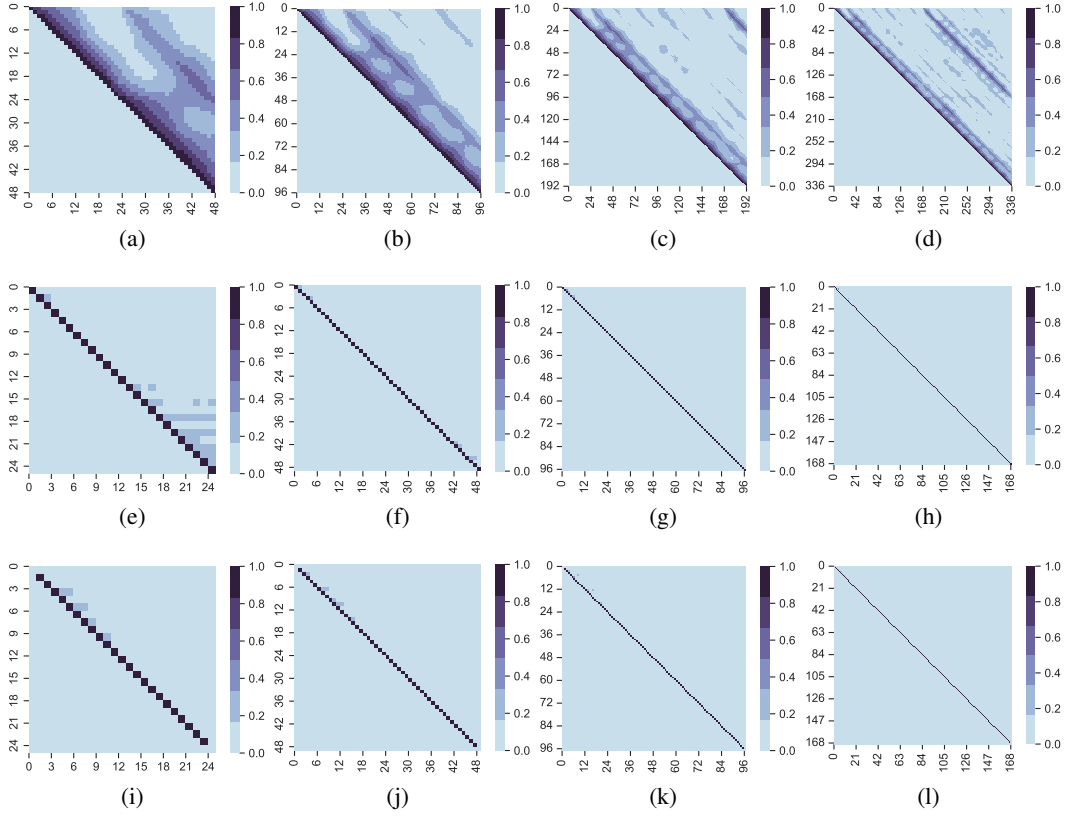


Figure 10: More comprehensive visualizations of label autocorrelation in different domains and label lengths, with columns representing label lengths  $H=48, 96, 192, 336$  from left to right. Panels (a-d) show the label correlation in the time domain, where each element  $\rho_{i,j}$  indicates the partial correlation between  $Y_i$  and  $Y_j$  given  $L$ . Panels (e-l) show the label correlation in the frequency domain, where each element  $\rho_{i,j}$  indicates the partial correlation between  $F_i$  and  $F_j$  given  $L$ , shown with the real (e-h) and imaginary part (i-l).

paradigm is biased, and the bias is quantified as:

$$\text{Bias} = \frac{1}{2\sigma^2}(Y_2 - \hat{Y}_2)^2 - \frac{1}{2\sigma^2(1 - \rho^2)}(Y_2 - (\hat{Y}_2 + \rho(Y_1 - \hat{Y}_1)))^2, \quad (7)$$

where  $\hat{Y}_i$  indicates the prediction at the  $i$ -th step and  $\rho$  denotes the partial correlation between  $Y_1$  and  $Y_2$  given  $L$ .

*Proof.* Aligning with the maximum likelihood analysis, we assume the label sequence obeys a normal distribution with mean  $\mu = [\hat{Y}_1, \hat{Y}_2]$  and covariance  $\zeta = [[\sigma^2, \rho\sigma^2], [\rho\sigma^2, \sigma_2^2]]$ . The negative log-likelihood (NLL) of  $Y$  given historical sequence  $L$  can be expressed as

$$\begin{aligned} -\log p(Y|L) &= -\log p(Y_1|L) - \log p(Y_2|L, Y_1) \\ &= -\log\left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(Y_1 - \hat{Y}_1)^2}{2\sigma^2}\right)\right) \\ &\quad - \log\left(\frac{1}{\sqrt{2\pi(1 - \rho^2)}\sigma} \exp\left(-\frac{(Y_2 - (\hat{Y}_2 + \rho(Y_1 - \hat{Y}_1)))^2}{2\sigma^2(1 - \rho^2)}\right)\right) \end{aligned}$$

Removing coefficients unrelated to  $g$ , the practical NLL that contributes the gradients to update  $g$  is

$$\text{NLL} := \frac{1}{2\sigma^2}(Y_1 - \hat{Y}_1)^2 + \frac{1}{2\sigma^2(1 - \rho^2)}(Y_2 - (\hat{Y}_2 + \rho(Y_1 - \hat{Y}_1)))^2$$

If the independence assumption of different time step holds (i.e.,  $Y_1$  and  $Y_2$  are conditionally independent given  $L$ ), we have  $\rho = 0$  and  $p(Y_2|L, Y_1) = p(Y_2|L)$ . In this case, the MSE loss that is employed by DF mirrors the practical NLL (the term  $\sigma$  is often set to 1 when implementing MSE):

$$\text{MSE} = \frac{1}{2\sigma^2}(Y_1 - \hat{Y}_1)^2 + \frac{1}{2\sigma^2}(Y_2 - \hat{Y}_2)^2$$

If the independence assumption does not hold, i.e.,  $\rho \neq 0$ , the MSE loss in the time domain is biased to the practical NLL. The bias is quantified as:

$$\text{Bias} = \frac{1}{2\sigma^2}(Y_2 - \hat{Y}_2)^2 - \frac{1}{2\sigma^2(1 - \rho^2)}(Y_2 - (\hat{Y}_2 + \rho(Y_1 - \hat{Y}_1)))^2$$

So far we have specified the bias introduced by label autocorrelation, which makes the MSE loss in the time domain fail to reflect the practical NLL and therefore misleads the update of forecast model  $g$  under DF paradigm.  $\square$

**Theorem B.2** (Bias of DF). *Given a input sequence  $L$  and a univariate label sequence  $Y$ , the learning objective (1) of the DF paradigm is biased against the practical NLL, expressed as:*

$$\text{Bias} = \sum_{i=1}^T \frac{1}{2\sigma^2}(Y_i - \hat{Y}_i)^2 - \sum_{i=1}^T \frac{1}{2\sigma^2(1 - \rho_i^2)} \left( Y_i - \left( \hat{Y}_i + \sum_{j=1}^{i-1} \rho_{ij}(Y_j - \hat{Y}_j) \right) \right)^2, \quad (8)$$

where  $\hat{Y}_i$  indicates the prediction at the  $i$ -th step,  $\rho_{ij}$  denotes the partial correlation between  $Y_i$  and  $Y_j$  given  $L$ ,  $\rho_i^2 = \sum_{j=1}^{i-1} \rho_{ij}^2$ .

*Proof.* Assume that the label sequence  $Y$  conditioned on the input sequence  $L$  follows a multivariate normal distribution with mean vector  $\mu = [\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_T]$  and covariance matrix  $\Sigma$ , where the diagonal entries  $\Sigma_{ii} = \sigma^2$  and the off-diagonal entries are  $\Sigma_{ij} = \rho_{ij}\sigma^2$  for  $i \neq j$ . Here,  $\rho_{ij}$  denotes the partial correlation between  $Y_i$  and  $Y_j$  given  $L$ . Under these assumptions, the NLL of the label sequence  $Y$  given  $L$  can be decomposed into a sum of conditional NLLs due to the properties of the multivariate normal distribution:

$$-\log p(Y | L) = -\sum_{i=1}^T \log p(Y_i | L, Y_1, Y_2, \dots, Y_{i-1}),$$

where each conditional probability  $p(Y_i \mid L, Y_1, \dots, Y_{i-1})$  is Gaussian with mean  $\hat{Y}_i + \sum_{j=1}^{i-1} \rho_{ij}(Y_j - \hat{Y}_j)$  and variance  $\sigma^2(1 - \rho_i^2)$ ,  $\rho_i^2 = \sum_{j=1}^{i-1} \rho_{ij}^2$ . Thus, the NLL can be expressed as

$$-\log p(Y \mid L) = \sum_{i=1}^T \left( \frac{1}{2} \log(2\pi\sigma^2(1 - \rho_i^2)) + \frac{1}{2\sigma^2(1 - \rho_i^2)} \left( Y_i - \left( \hat{Y}_i + \sum_{j=1}^{i-1} \rho_{ij}(Y_j - \hat{Y}_j) \right) \right)^2 \right).$$

For the purpose of gradient-based optimization, constant terms independent of the model predictions  $\hat{Y}_i$  can be omitted. Therefore, the practical NLL contributing to the gradients is given by

$$\text{NLL} = \sum_{i=1}^T \frac{1}{2\sigma^2(1 - \rho_i^2)} \left( Y_i - \left( \hat{Y}_i + \sum_{j=1}^{i-1} \rho_{ij}(Y_j - \hat{Y}_j) \right) \right)^2,$$

which immediately follows from discarding the constants irrelevant with the gradient of  $\hat{Y}$ .

On the other hand, the DF paradigm typically employs the MSE loss, expressed as

$$\text{MSE} = \sum_{i=1}^T \frac{1}{2\sigma^2} (Y_i - \hat{Y}_i)^2.$$

which deviates from the practical NLL. The bias is expressed as:

$$\text{Bias} = \text{MSE} - \text{NLL} = \sum_{i=1}^T \frac{1}{2\sigma^2} (Y_i - \hat{Y}_i)^2 - \sum_{i=1}^T \frac{1}{2\sigma^2(1 - \rho_i^2)} \left( Y_i - \hat{Y}_i + \sum_{j=1}^{i-1} \rho_{ij}(Y_j - \hat{Y}_j) \right)^2.$$

When there exists label autocorrelation, i.e.,  $\rho_{ij} \neq 0$ , the bias above exists. *In the special case (e.g., calculating the loss in the frequency domain) where the label autocorrelation is diminished, i.e.,  $\rho_{ij} \rightarrow 0$ , the bias approaches zero almost surely.*

□

**Corollary B.3** (Bias of DF, multivariate). *Given a input sequence  $L$  and a multivariate label sequence  $Y \in \mathbb{R}^{T \times D}$ , suppose  $Z \in \mathbb{R}^{T \times D}$  be the flattened version of  $Y$  obtained by concatenating the rows, the learning objective (1) of the DF paradigm is biased against the practical NLL, expressed as:*

$$\text{Bias} = \sum_{i=1}^{T \times D} \frac{1}{2\sigma^2} (Z_i - \hat{Z}_i)^2 - \sum_{i=1}^{T \times D} \frac{1}{2\sigma^2(1 - \rho_i^2)} \left( Z_i - \left( \hat{Z}_i + \sum_{j=1}^{i-1} \rho_{ij}(Z_j - \hat{Z}_j) \right) \right)^2, \quad (9)$$

where  $\hat{Z}_i$  indicates the prediction of  $Z_i$ ,  $\rho_{ij}$  denotes the partial correlation between  $Z_i$  and  $Z_j$  given  $L$ ,  $\rho_i^2 = \sum_{j=1}^{i-1} \rho_{ij}^2$ .

*Proof.* This corollary immediately follows from Theorem B.2, by viewing the multivariate label sequence  $Z$  as an augmented univariate sequence. □

**Theorem B.4** (Decorrelation between frequency components). *Let  $Y$  be a zero-mean, discrete-time, wide-sense stationary random process of length  $T$ . As  $T \rightarrow \infty$ , the DFT coefficients become asymptotically uncorrelated at different frequencies:*

$$\lim_{T \rightarrow \infty} \mathbb{E}[F_k F_{k'}^*] = \begin{cases} S_Y(f_k), & \text{if } k = k', \\ 0, & \text{if } k \neq k', \end{cases}$$

where  $f_k = \frac{k}{T}$  and  $S_Y(f)$  is the power spectral density of  $Y$ .

*Proof.* Recalling that the normalized DFT coefficients  $F_k$  are defined as  $F_k = 1/\sqrt{T} \sum_{t=0}^{T-1} Y_t e^{-j2\pi kt/T}$ ,  $k = 0, 1, \dots, T-1$ . On this basis, the expected value of the



product  $F_k F_{k'}^*$  can be expressed as:

$$\begin{aligned}\mathbb{E}[F_k F_{k'}^*] &= \mathbb{E} \left[ \sum_{t=0}^{T-1} Y_t e^{-j2\pi kt/T} \cdot \sum_{t'=0}^{T-1} Y_{t'}^* e^{j2\pi k' t'/T} \right] / T \\ &= \sum_{t=0}^{T-1} \sum_{t'=0}^{T-1} R_Y[t - t'] e^{-j2\pi kt/T} e^{j2\pi k' t'/T} / T,\end{aligned}\tag{10}$$

which immediately follows from interchanging the order of summation and expectation, and using the autocorrelation function  $R_Y[t - t'] = \mathbb{E}[Y_t Y_{t'}^*]$ . Let  $\tau = t - t'$ ; then  $t' = t - \tau$ . Substituting and simplifying the exponentials gives:

$$\begin{aligned}\mathbb{E}[F_k F_{k'}^*] &= \sum_{t=0}^{T-1} \sum_{\tau=-t}^{T-1-t} R_Y[\tau] e^{-j2\pi(k t/T - k'(t-\tau)/T)} / T \\ &= \sum_{\tau=-(T-1)}^{T-1} R_Y[\tau] e^{j2\pi k' \tau/T} \left( \sum_{t=\max(0, \tau)}^{\min(T-1, T-1+\tau)} e^{j2\pi(k-k')t/T} / T \right).\end{aligned}$$

where the inner sum over  $t$  involves a sum of complex exponentials. When  $k \neq k'$ , the exponentials oscillate rapidly as  $T$  increases, causing the inner term to diminish due to destructive interference. That is, for  $k \neq k'$ :

$$\lim_{T \rightarrow \infty} \mathbb{E}[F_k F_{k'}^*] = 0.$$

When  $k = k'$ , the exponential term becomes unity, and the inner sum simplifies to:

$$\lim_{T \rightarrow \infty} \sum_{t=\max(0, \tau)}^{\min(T-1, T-1+\tau)} 1/T = \lim_{T \rightarrow \infty} 1 - |\tau|/T = 1.$$

which immediately follows by  $\mathbb{E}[F_k F_{k'}^*] = S_Y(f_k)$ , where  $S_Y$  is the power spectral density of  $Y$  that can be calculated as the DFT of  $R_Y$ . The proof is therefore completed.  $\square$

**Theorem B.5** (Bias of DF, extended). *Given input sequence  $L$  and label sequence  $Y$ , the learning objective (1) of the DF paradigm is biased against the practical NLL, expressed as:*

$$\text{Bias} = \sum_{i=1}^T \frac{1}{2\sigma^2} (Y_i - \hat{Y}_i)^2 - \sum_{i=1}^T \frac{1}{2\sigma^2(1 - \rho_i^2)} \left( Y_i - \left( \hat{Y}_i + \sum_{j=1}^{i-1} \rho_{ij} (Y_j - \hat{Y}_j) \right) \right)^2, \tag{11}$$

where  $\hat{Y}_i$  indicates the prediction at the  $i$ -th step,  $\rho_{ij}$  denotes the partial correlation between  $Y_i$  and  $Y_j$  given  $L$ ,  $\rho_i^2 = \sum_{j=1}^{i-1} \rho_{ij}^2$ .

## C GENERALIZED 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:

$$\begin{aligned}f_k(t) &= \exp(-j(2\pi/H)kt), \\ F_k &= \int_{-\infty}^{\infty} x(t) f_k(t) dt\end{aligned}\tag{12}$$

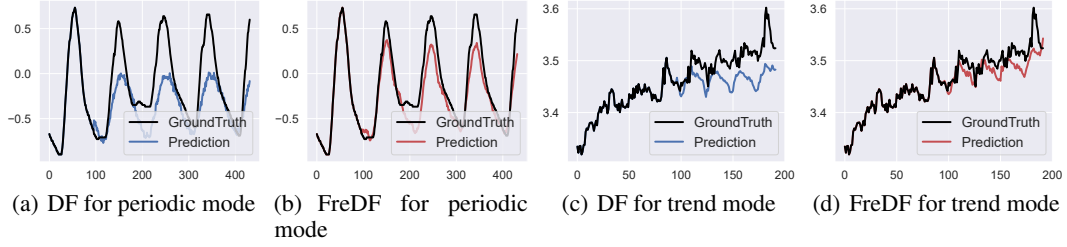


Figure 11: Forecast sequences generated by iTransformer on the snapshots where periodicity and trend are dominant. The bases are selected as Fourier polynomial (b) for periodic mode and Legendre (d) polynomial for trend mode.

**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 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:

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

$$F_k = \int_{-1}^1 x(t) f_k(t) dt$$
(13)

**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.

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

$$F_k = \int_{-1}^1 \frac{x(t) f_k(t)}{\sqrt{1-t^2}} dt$$
(14)

**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:

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

$$F_k = \int_0^\infty \frac{x(t) f_k(t)}{\exp(t)} dt$$
(15)

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 11. 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.

Table 4: 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 aligned with the recent works (Wu et al., 2023; Liu et al., 2024).

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
ETTm1	7	96, 192, 336, 720	34465/11521/11521	15min	Health
ETTm2	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
PEMS03	358	12, 24, 36, 48	15617/5135/5135	5min	Transportation
PEMS08	170	12, 24, 36, 48	10690/3548/265	5min	Transportation

## D REPRODUCTION DETAILS

**All scripts have been incorporated in our anonymous repository, to reproduce the overall performance results, ablation results and hyper-parameter sensitivity results.** They will be released publicly with pretrained checkpoints, hyperparameter settings and logging files.

### 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:

- ETT (Li et al., 2021) 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 (Wu et al., 2021) includes 21 meteorological variables gathered every 10 minutes throughout 2020 from the Weather Station of the Max Planck Biogeochemistry Institute.
- ECL (Electricity Consumption Load) (Wu et al., 2021) presents hourly electricity consumption data for 321 clients.
- Traffic (Wu et al., 2021) features hourly road occupancy rates from 862 sensors in the San Francisco Bay area freeways, spanning from January 2015 to December 2016.
- PEMS (Liu et al., 2022a) contains the public traffic network data in California collected by 5-minute windows. Two public subsets (PEMS03, PEMS08) are adopted in this work.

Data processing and the division into training, validation, and testing sets adhere to the protocol established by TimesNet (Wu et al., 2023). 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 4.

### D.2 IMPLEMENTATION DETAILS

The baseline models in this study were meticulously reproduced using training scripts obtained from the TimesNet Repository (Wu et al., 2023) after reproducibility verification. Models were trained employing the Adam optimizer (Kingma & Ba, 2015), 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.

In experiments integrating FreDF for existing forecast models, we closely adhered to the original hyperparameter settings 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. Although finetuning for each prediction length separately can further boost the performance, we omitted it since the efficacy of FreDF does not rely on dedicated 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 5. 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 31 out of 45 cases and MAE in 40 out of 45 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).

**Case Study with PatchTST and Varying Historical Length.** PatchTST (Nie et al., 2023) is a powerful baseline, whose predictive performance strongly correlates with historical length. To explore this, we open a thread here to study the performance of iTransformer and PatchTST under different historical lengths, and to observe whether FreDF could improve both. The results in Table 6 with the Weather dataset indicate a consistent improvement with the integration of FreDF for both iTransformer and PatchTST. It’s noteworthy that under our experimental conditions, PatchTST with  $H = 336$  achieved results comparable to the original “PatchTST/42” in Nie et al. (2023), and FreDF still managed to reduce the MSE by 0.002, demonstrating its robustness across historical lengths.

**Short-term forecast.** We provide a detailed comparison for the short-term forecast task in Table 7, 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. 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 investigate missing data imputation task. iTransformer, identified as the best baseline for imputation tasks, is selected as the testbed for FreDF. iTransformer is selected as the base model for FreDF implementation. 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 empirical results in Table 8 demonstrate the efficacy of FreDF in this task: it improves the performance of iTransformer significantly, outperforming most competitive methods, hitting the minimum MSE in 23 out of 30 cases and minimum MAE in 19 out of 30 cases. 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.

Table 5: Full results on the long-term forecasting task. The length of history window is set to 96 for all baselines. Avg indicates the results averaged over forecasting lengths: T=96, 192, 336 and 720.

Models		FreDF (Ours)		iTransformer (2024)		FreTS (2023)		TimesNet (2023)		MICN (2023)		TiDE (2023)		DLinear (2023)		FEDformer (2022)		Autoformer (2021)		Transformer (2017)		TCN (2017)	
Metrics		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	96	0.324	0.362	0.346	0.379	0.339	0.374	0.338	0.379	0.318	0.366	0.364	0.387	0.345	0.372	0.389	0.427	0.468	0.463	0.591	0.549	0.887	0.613
	192	0.373	0.385	0.392	0.400	0.382	0.397	0.389	0.400	0.364	0.396	0.398	0.404	0.381	0.390	0.402	0.431	0.573	0.509	0.704	0.629	0.877	0.626
	336	0.402	0.404	0.427	0.422	0.421	0.426	0.429	0.428	0.398	0.428	0.428	0.425	0.414	0.414	0.438	0.451	0.596	0.527	1.171	0.861	0.890	0.636
	720	0.469	0.444	0.494	0.461	0.485	0.462	0.495	0.464	0.514	0.501	0.487	0.461	0.473	0.451	0.529	0.498	0.749	0.569	1.307	0.893	0.911	0.653
	Avg	0.392	0.399	0.415	0.416	0.407	0.415	0.413	0.418	0.399	0.423	0.419	0.419	0.404	0.407	0.440	0.451	0.596	0.517	0.943	0.733	0.891	0.632
ETTh2	96	0.173	0.252	0.184	0.266	0.190	0.282	0.185	0.264	0.178	0.275	0.207	0.305	0.195	0.294	0.194	0.284	0.240	0.319	0.317	0.408	3.125	1.345
	192	0.241	0.298	0.257	0.315	0.260	0.329	0.254	0.307	0.240	0.317	0.290	0.364	0.283	0.359	0.264	0.324	0.300	0.349	1.069	0.758	3.130	1.350
	336	0.298	0.334	0.315	0.351	0.373	0.405	0.314	0.345	0.299	0.354	0.377	0.422	0.384	0.427	0.319	0.359	0.339	0.375	1.325	0.869	3.185	1.375
	720	0.398	0.393	0.419	0.409	0.517	0.499	0.434	0.413	0.482	0.479	0.558	0.524	0.516	0.502	0.430	0.424	0.423	0.421	2.576	1.223	4.203	1.658
	Avg	0.278	0.319	0.294	0.335	0.335	0.379	0.297	0.332	0.300	0.356	0.358	0.404	0.344	0.396	0.302	0.348	0.326	0.366	1.322	0.814	3.411	1.432
ETTh1	96	0.382	0.400	0.390	0.410	0.399	0.412	0.422	0.433	0.383	0.418	0.479	0.464	0.396	0.410	0.377	0.418	0.423	0.441	0.796	0.691	0.767	0.633
	192	0.430	0.427	0.443	0.441	0.453	0.443	0.465	0.457	0.500	0.491	0.521	0.503	0.449	0.444	0.421	0.445	0.498	0.485	0.813	0.699	0.739	0.619
	336	0.474	0.451	0.480	0.457	0.503	0.475	0.492	0.470	0.546	0.530	0.659	0.603	0.487	0.465	0.468	0.472	0.506	0.496	1.181	0.876	0.717	0.613
	720	0.463	0.462	0.484	0.479	0.596	0.565	0.532	0.502	0.671	0.620	0.893	0.736	0.516	0.513	0.500	0.493	0.477	0.487	1.182	0.885	0.828	0.678
	Avg	0.437	0.435	0.449	0.447	0.488	0.474	0.478	0.466	0.525	0.515	0.628	0.574	0.462	0.458	0.441	0.457	0.476	0.477	0.993	0.788	0.763	0.636
ETTh2	96	0.289	0.337	0.301	0.349	0.350	0.403	0.320	0.364	0.361	0.404	0.400	0.440	0.343	0.396	0.347	0.391	0.383	0.424	2.072	1.140	3.171	1.364
	192	0.363	0.385	0.382	0.402	0.472	0.475	0.409	0.417	0.495	0.490	0.528	0.509	0.473	0.474	0.430	0.443	0.557	0.511	5.081	1.814	3.222	1.398
	336	0.419	0.426	0.430	0.434	0.564	0.528	0.449	0.451	0.671	0.588	0.643	0.571	0.603	0.546	0.469	0.475	0.470	0.481	3.564	1.475	3.306	1.452
	720	0.415	0.437	0.447	0.455	0.815	0.654	0.473	0.474	0.968	0.712	0.874	0.679	0.812	0.650	0.473	0.480	0.501	0.515	2.469	1.247	3.599	1.565
	Avg	0.371	0.396	0.390	0.410	0.550	0.515	0.413	0.426	0.624	0.549	0.611	0.550	0.558	0.516	0.430	0.447	0.478	0.483	3.296	1.419	3.325	1.445
ECL	96	0.144	0.233	0.148	0.239	0.189	0.277	0.171	0.273	0.168	0.280	0.237	0.329	0.210	0.302	0.200	0.315	0.199	0.315	0.252	0.352	0.688	0.621
	192	0.159	0.247	0.167	0.258	0.193	0.282	0.188	0.289	0.177	0.289	0.236	0.330	0.210	0.305	0.207	0.322	0.215	0.327	0.266	0.364	0.587	0.582
	336	0.172	0.263	0.179	0.272	0.207	0.296	0.208	0.304	0.185	0.296	0.249	0.344	0.223	0.319	0.226	0.340	0.232	0.343	0.292	0.383	0.590	0.588
	720	0.204	0.294	0.209	0.298	0.245	0.332	0.289	0.363	0.218	0.323	0.284	0.373	0.258	0.350	0.282	0.379	0.268	0.371	0.287	0.371	0.602	0.601
	Avg	0.170	0.259	0.176	0.267	0.209	0.297	0.214	0.307	0.187	0.297	0.251	0.344	0.225	0.319	0.229	0.339	0.228	0.339	0.274	0.367	0.617	0.598
Traffic	96	0.391	0.265	0.397	0.272	0.528	0.341	0.504	0.298	0.609	0.317	0.805	0.493	0.697	0.429	0.577	0.362	0.609	0.385	0.686	0.385	1.451	0.744
	192	0.410	0.273	0.418	0.279	0.531	0.338	0.526	0.305	0.621	0.328	0.756	0.474	0.647	0.407	0.603	0.372	0.633	0.400	0.679	0.377	0.842	0.622
	336	0.424	0.280	0.432	0.286	0.551	0.345	0.540	0.310	0.641	0.342	0.762	0.477	0.653	0.410	0.615	0.378	0.637	0.398	0.663	0.361	0.844	0.620
	720	0.460	0.298	0.467	0.305	0.598	0.367	0.570	0.324	0.671	0.354	0.719	0.449	0.694	0.429	0.649	0.403	0.668	0.415	0.693	0.381	0.867	0.624
	Avg	0.421	0.279	0.428	0.286	0.552	0.348	0.535	0.309	0.636	0.335	0.760	0.473	0.673	0.419	0.611	0.379	0.637	0.399	0.680	0.376	1.001	0.652
Weather	96	0.164	0.202	0.201	0.247	0.184	0.239	0.178	0.226	0.182	0.250	0.202	0.261	0.197	0.259	0.221	0.304	0.284	0.355	0.332	0.383	0.610	0.568
	192	0.220	0.253	0.250	0.283	0.223	0.275	0.227	0.266	0.234	0.301	0.242	0.298	0.236	0.294	0.275	0.345	0.313	0.371	0.634	0.539	0.541	0.552
	336	0.275	0.294	0.302	0.317	0.272	0.316	0.283	0.303	0.268	0.325	0.287	0.335	0.282	0.332	0.338	0.379	0.359	0.393	0.656	0.579	0.565	0.569
	720	0.356	0.347	0.370	0.362	0.340	0.363	0.359	0.355	0.361	0.399	0.351	0.386	0.347	0.384	0.408	0.418	0.440	0.446	0.908	0.706	0.622	0.601
	Avg	0.254	0.274	0.281	0.302	0.255	0.299	0.262	0.288	0.261	0.319	0.271	0.320	0.265	0.317	0.311	0.361	0.349	0.391	0.632	0.552	0.584	0.572
PEMS03	12	0.068	0.172	0.069	0.175	0.083	0.194	0.082	0.188	0.087	0.203	0.117	0.225	0.122	0.245	0.123	0.248	0.239	0.365	0.107	0.209	0.632	0.606
	24	0.096	0.205	0.098	0.210	0.127	0.241	0.110	0.216	0.086	0.198	0.233	0.320	0.202	0.320	0.160	0.287	0.492	0.506	0.121	0.227	0.655	0.626
	36	0.128	0.240	0.131	0.243	0.169	0.281	0.133	0.236	0.105	0.220	0.380	0.422	0.275	0.382	0.191	0.321	0.399	0.459	0.133	0.243	0.678	0.644
	48	0.161	0.269	0.164	0.275	0.204	0.311	0.146	0.251	0.120	0.235	0.536	0.511	0.335	0.429	0.223	0.350	0.875	0.723	0.144	0.253	0.699	0.659
	Avg	0.113	0.219	0.116	0.226	0.146	0.257	0.118	0.223	0.099	0.214	0.316	0.370	0.233	0.344	0.174	0.302	0.501	0.513	0.126	0.233	0.666	0.634
PEMS08	12	0.080	0.182	0.085	0.189	0.095	0.204	0.110	0.209	2.193	0.871	0.121	0.231	0.152	0.274	0.175	0.275	0.446	0.483	0.213	0.236	0.680	0.607
	24	0.118	0.220	0.131	0.236	0.150	0.259	0.142	0.239	0.235	0.339	0.232	0.326	0.245	0.350	0.211	0.305	0.488	0.509	0.238	0.256	0.701	0.622
	36	0.161	0.258	0.182	0.282	0.202	0.305	0.167	0.258	0.197	0.300	0.379	0.428	0.344	0.417	0.250	0.338	0.532	0.513	0.263	0.277	0.727	0.637
	48	0.206	0.293	0.236	0.323	0.250	0.341	0.195	0.274	0.242	0.324	0.543	0.527	0.437	0.469	0.293	0.371	1.052	0.781	0.283	0.295	0.746	0.648
	Avg	0.141	0.238	0.159	0.258	0.174	0.277	0.154	0.245	0.717	0.459	0.319	0.378	0.294	0.377	0.232	0.322	0.630	0.572	0.249	0.266	0.713	0.629
1 <sup>st</sup> Count	31	40	0	0	1	0	1	1	1	10	4	0	0	0	0	3	0	0	0	0	0	0	0



Table 6: Evaluation results on varying history window length with Weather dataset.

Models		FreDF		iTransformer		FreDF		PatchTST		
Metrics		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
History Window Length	96	96	0.164	0.202	0.201	0.247	0.174	0.217	0.200	0.244
		192	0.220	0.253	0.250	0.283	0.230	0.266	0.234	0.268
		336	0.275	0.294	0.302	0.317	0.279	0.301	0.311	0.321
		720	0.356	0.347	0.370	0.362	0.355	0.351	0.365	0.353
		Avg	0.254	0.274	0.281	0.302	0.259	0.284	0.278	0.297
	192	96	0.164	0.207	0.184	0.235	0.158	0.205	0.167	0.213
		192	0.211	0.250	0.236	0.277	0.200	0.241	0.204	0.244
		336	0.262	0.290	0.268	0.296	0.259	0.287	0.266	0.291
		720	0.341	0.343	0.342	0.345	0.330	0.334	0.333	0.337
		Avg	0.244	0.272	0.258	0.288	0.237	0.267	0.242	0.271
	336	96	0.159	0.204	0.164	0.215	0.150	0.200	0.153	0.203
		192	0.204	0.248	0.211	0.256	0.193	0.240	0.194	0.240
		336	0.253	0.288	0.260	0.292	0.245	0.280	0.247	0.282
		720	0.325	0.336	0.327	0.339	0.320	0.332	0.321	0.336
		Avg	0.235	0.269	0.241	0.276	0.227	0.263	0.229	0.265
	720	96	0.164	0.215	0.172	0.228	0.144	0.194	0.191	0.246
		192	0.209	0.257	0.218	0.265	0.190	0.242	0.192	0.241
		336	0.251	0.291	0.273	0.306	0.243	0.283	0.241	0.285
		720	0.318	0.342	0.340	0.353	0.310	0.330	0.311	0.331
		Avg	0.236	0.276	0.251	0.288	0.222	0.262	0.234	0.276

Table 7: Full results on the short-term forecasting task. Avg indicates the results averaged over forecasting lengths: yearly, quarterly, and monthly.

Models	FreDF (Ours)			FreTS (2023)			iTransformer (2024)			MICN (2023)			DLinear (2023)			Fedformer (2023)			Autoformer (2023)		
Metric	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA
Yearly	<b>13.556</b>	<b>3.046</b>	<b>0.798</b>	<b>13.576</b>	<b>3.068</b>	<b>0.801</b>	13.797	3.143	0.818	14.594	3.392	0.873	14.307	3.094	0.827	13.648	3.089	0.806	18.477	4.26	1.101
Quarterly	<b>10.374</b>	<b>1.229</b>	<b>0.919</b>	<b>10.361</b>	<b>1.223</b>	<b>0.916</b>	10.503	1.248	0.932	11.417	1.385	1.023	10.500	1.237	0.928	10.612	1.246	0.936	14.254	1.829	1.314
Monthly	<b>12.999</b>	<b>0.983</b>	<b>0.913</b>	<b>13.088</b>	<b>0.99</b>	<b>0.919</b>	13.227	1.013	0.935	13.834	1.080	0.987	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	<b>5.101</b>	<b>3.419</b>	<b>1.076</b>	6.137	4.201	1.308	5.12	3.649	1.114	<b>4.823</b>	<b>3.243</b>	<b>1.019</b>	6.772	4.963	1.495
Avg.	<b>12.112</b>	<b>1.648</b>	<b>0.877</b>	<b>12.169</b>	<b>1.66</b>	<b>0.883</b>	12.298	1.68	0.893	13.044	1.841	0.962	12.48	1.674	0.898	12.734	1.702	0.914	16.851	2.443	1.26
1 <sup>st</sup> Count	<b>3</b>	<b>3</b>	<b>3</b>	<b>1</b>	<b>1</b>	<b>1</b>	0	0	0	0	0	0	0	0	0	<b>1</b>	<b>1</b>	<b>1</b>	0	0	0

**Showcases.** We provide additional showcases illustrating the improvements in forecast sequences by integrating FreDF in Figure 12 and 14. 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 RUNNING COST ANALYSIS

In this section, we analyze the computational complexity and running cost of FreDF through empirical investigation. The core computation of FreDF involves calculating the FFT of both predicted and label sequences, followed by calculating their point-wise MAE loss. Therefore, the overall complexity of FreDF is governed primarily by the FFT operation, which is  $\mathcal{O}(T \log T)$ , where  $T$  is the length of predicted sequence. Figure 16 illustrates the empirical running costs of FreDF for varying sequence lengths, capturing both the Forward Pass (calculating FFT for the predicted sequence) and Backward Pass (computing the frequency loss and gradients w.r.t. the prediction sequence). Our results confirm that the additional computational duration imposed by FreDF is about 1ms for a prediction sequence with  $T < 720$ . Importantly, the frequency loss computation is not required during inference. Therefore, FreDF does not compromise the model’s efficiency during either training or inference stages.

Table 8: Full results on the missing data imputation task. The length of history window is set to 96 for all baselines. Avg indicates the results averaged over missing ratios: 0.125, 0.25, 0.375, 0.5.

Models	FreDF (Ours)			iTransformer (2024)		FreTS (2023)		TimesNet (2023)		MICN (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	MSE	MAE
ETm1	0.125	0.00153	0.02790	0.00213	0.03307	0.01102	0.07843	0.01152	0.07267	0.00236	0.03371	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	0.00284	0.03691	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	0.00323	0.03900	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	0.00352	0.04028	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	0.00299	0.03747	0.44776	0.46572	0.00167	0.02334	0.68202	0.38104	0.37660	0.35798
ETm2	0.125	0.00363	0.03840	0.00398	0.04034	0.03194	0.13349	0.01189	0.06710	0.00219	0.03345	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	0.00331	0.04100	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	0.00431	0.04598	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	0.00505	0.04918	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	0.00371	0.04240	0.93777	0.66893	0.02405	0.09938	3.10402	1.31356	1.41655	0.81246
ETTh1	0.125	0.00178	0.03059	0.00319	0.04102	0.01400	0.08181	0.00441	0.04403	0.00432	0.04655	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	0.00454	0.04769	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	0.00454	0.04730	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	0.00437	0.04594	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	0.00444	0.04687	0.24795	0.36749	0.00225	0.03226	0.68197	0.37992	0.43848	0.41915
ETTh2	0.125	0.00222	0.03124	0.00473	0.04606	0.04485	0.13849	0.00535	0.04495	0.00334	0.04202	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	0.00457	0.04950	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	0.00535	0.05363	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	0.00584	0.05547	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	0.00477	0.05016	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	0.003678	0.14078	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	0.004106	0.14847	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	0.004373	0.15224	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	0.004520	0.15380	0.21280	0.34526	0.10362	0.23127	0.45891	0.41011	0.21322	0.30764
	Avg	0.00071	0.01331	0.00197	0.03333	0.01022	0.08266	0.00302	0.03682	0.004169	0.14882	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	0.00317	0.03646	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	0.00325	0.03900	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	0.00326	0.03997	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	0.00309	0.03929	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	0.00320	0.03868	0.24106	0.32620	0.00484	0.05016	0.40554	0.42633	0.13727	0.18382
1 <sup>st</sup> Count		23	19	1	1	0	0	0	2	2	2	0	0	4	6	0	0	0	0

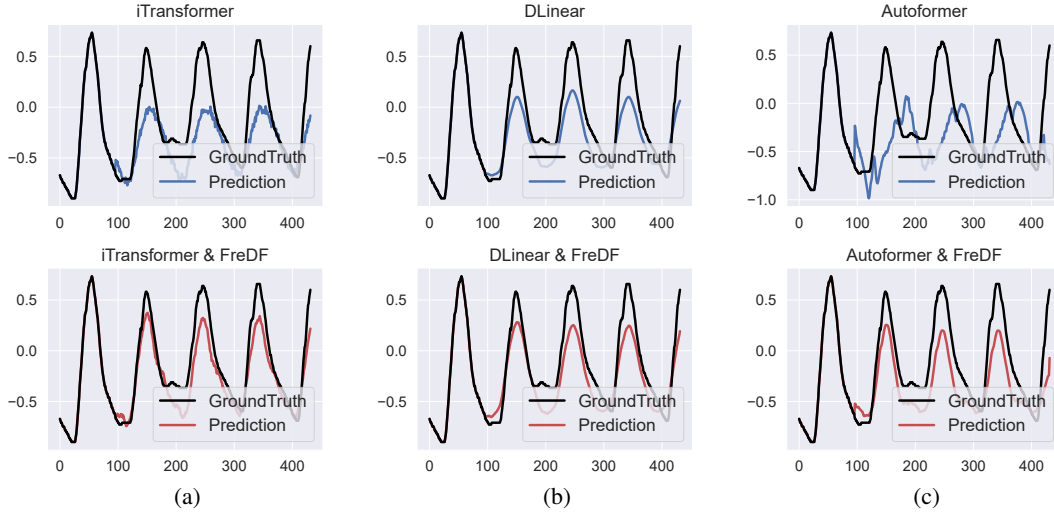


Figure 12: Forecast sequences generated by iTransformer, Dlinear and Autoformer with and without FreDF. The prediction length is set to 336 and the experiment is conducted on a snapshot of ETm2.

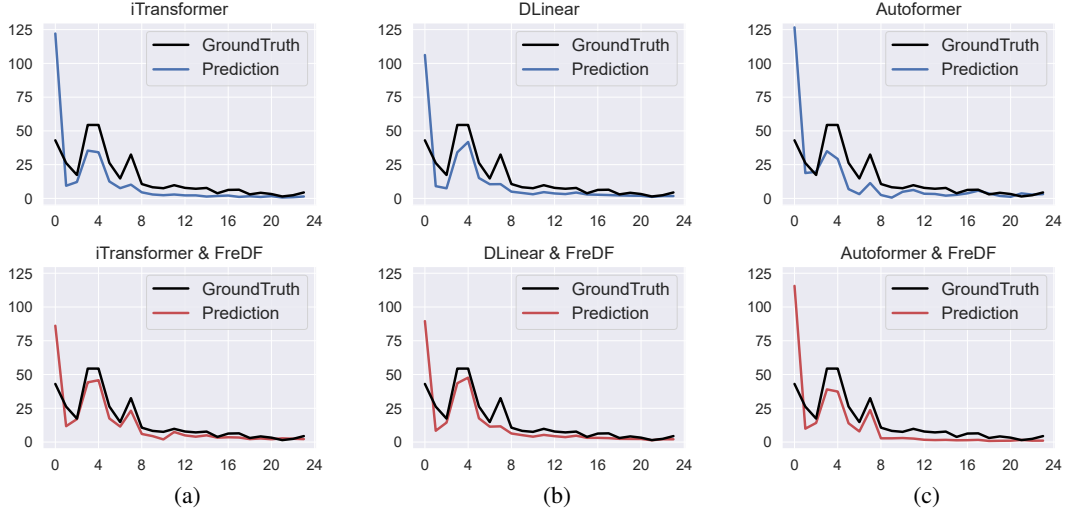


Figure 13: Spectrum of forecast sequences generated by iTransformer, Dlinear and Autoformer with and without FreDF, where only the first 24 frequencies of the spectrum are selected. The prediction length is set to 336 and the experiment is conducted on a snapshot of ETTm2.

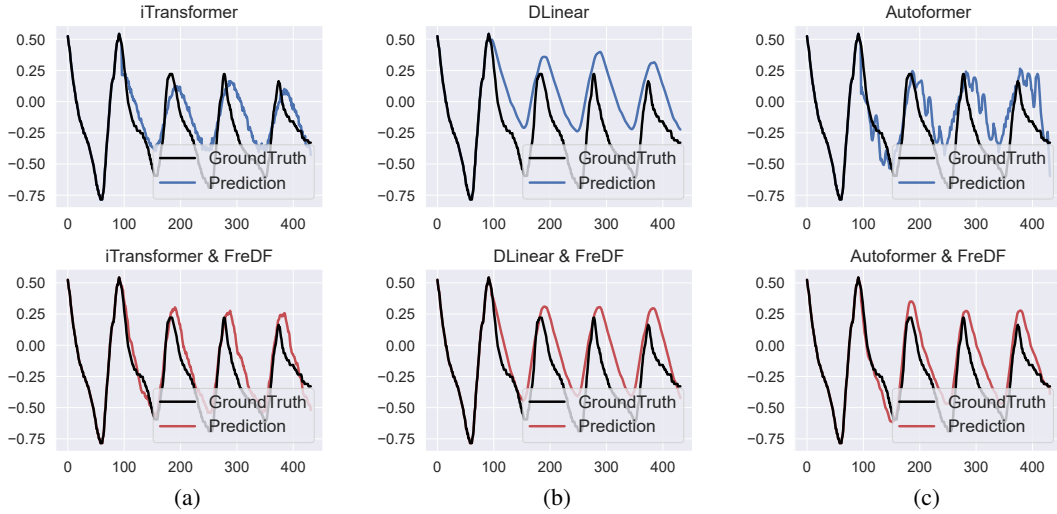


Figure 14: Forecast sequences generated by generated by iTransformer, Dlinear and Autoformer with and without FreDF. The prediction length is set to 336 and the experiment is conducted on a snapshot of ETTm2.

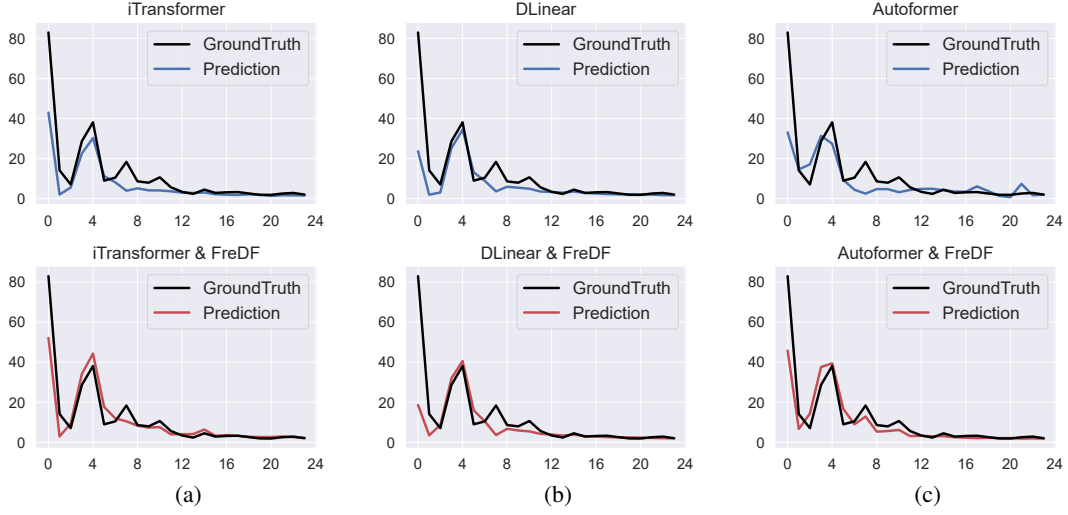


Figure 15: Spectrum of forecast sequences generated by iTransformer, Dlinear and Autoformer with and without FreDF, where only the first 24 frequencies of the spectrum are selected. The prediction length is set to 336 and the experiment is conducted on a snapshot of ETTm2.

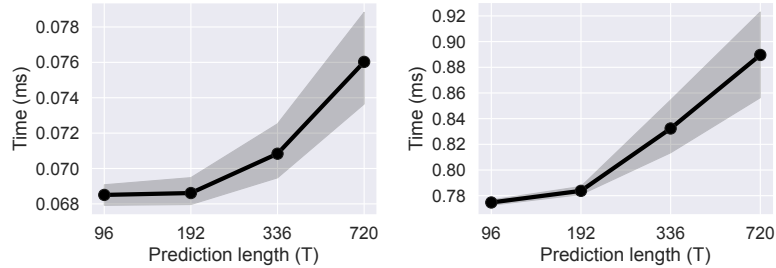


Figure 16: Running time of the frequency loss in the forward pass (left panel) and backward pass (right panel), shown with dashed lines for average and shaded areas for 99.9% confidence intervals.

Table 9: Experimental results (mean $\pm$ std) with varying seeds (2020, 2021, 2022, 2023, 2024).

Dataset	ETTh1				Weather			
Models	FreDF		iTransformer		FreDF		iTransformer	
Metrics	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
96	0.377 $\pm$ 0.001	0.396 $\pm$ 0.001	0.391 $\pm$ 0.001	0.409 $\pm$ 0.001	0.168 $\pm$ 0.003	0.205 $\pm$ 0.003	0.203 $\pm$ 0.002	0.246 $\pm$ 0.002
192	0.428 $\pm$ 0.001	0.424 $\pm$ 0.001	0.446 $\pm$ 0.002	0.441 $\pm$ 0.002	0.220 $\pm$ 0.001	0.254 $\pm$ 0.001	0.249 $\pm$ 0.001	0.281 $\pm$ 0.001
336	0.466 $\pm$ 0.001	0.442 $\pm$ 0.001	0.484 $\pm$ 0.005	0.460 $\pm$ 0.003	0.281 $\pm$ 0.002	0.298 $\pm$ 0.002	0.299 $\pm$ 0.002	0.315 $\pm$ 0.002
720	0.468 $\pm$ 0.005	0.465 $\pm$ 0.003	0.499 $\pm$ 0.015	0.489 $\pm$ 0.010	0.364 $\pm$ 0.008	0.354 $\pm$ 0.006	0.371 $\pm$ 0.001	0.361 $\pm$ 0.001
Avg	0.435 $\pm$ 0.002	0.432 $\pm$ 0.002	0.455 $\pm$ 0.006	0.450 $\pm$ 0.004	0.258 $\pm$ 0.004	0.278 $\pm$ 0.003	0.280 $\pm$ 0.001	0.301 $\pm$ 0.002

Table 10: Performance comparison of aligning amplitude and phase characteristics.

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

### E.3 RANDOM SEED SENSITIVITY

In this section, we investigate the sensitivity of the results to the specification of random seeds. To this end, we report the mean and standard deviation of the results by conducting experiments using 5 random seeds (2020, 2021, 2022, 2023, 2024). We investigate (1) iTransformer; (2) FreDF, which is applied to refine iTransformer. The results in Table 9 indicate minimal performance variation, with standard deviations below 0.005 in 7 out of 8 Avg cases. This showcases the insensitivity of models in time series forecast to different random seed specifications.

### E.4 AMPLITUDE V.S. PHASE ALIGNMENT

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 signal processing, both are foundational for representing the dynamics of signals. We dissect their contributions in Table 10 with results averaged over forecast lengths. Overall, both characteristics are essential for FreDF. In particular, phase alignment emerges as particularly crucial: aligning amplitude characteristics without phase alignment results in poor performance. This outcome is reasonable, since minor deviations of phase characteristics could correspond significant discrepancies in the time domain.

### E.5 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 17 across five distinct datasets. The forecast errors are averaged over prediction lengths (96, 192, 336, 720), with error bars as 95% confidence intervals.

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 solidifies FreDF’s position as a powerful tool capable of addressing

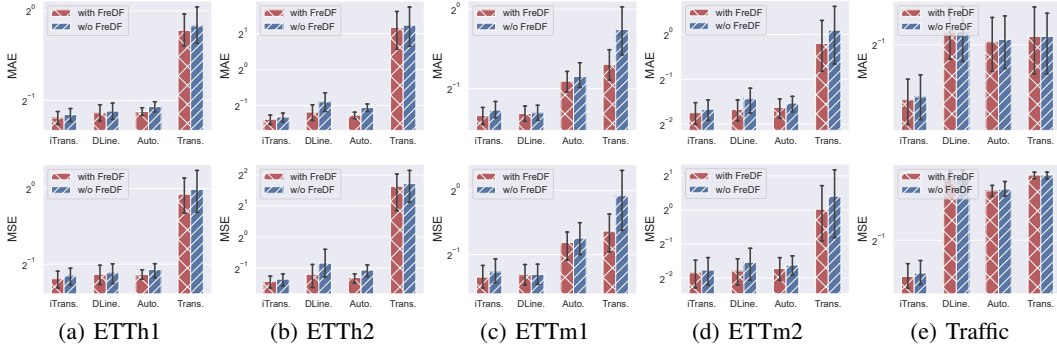


Figure 17: Performance of different forecast models with and without FreDF. The forecast errors are averaged over prediction lengths and the error bars represent 95% confidence intervals.

a wide array of forecasting challenges, marking it as a significant contribution to the advancement of forecasting methodologies.

## E.6 HYPERPARAMETER SENSITIVITY

In this section, we investigate the influence of adjusting the frequency loss parameter,  $\alpha$ , 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 18, 19, and 20.

A consistent observation across these models is that incrementally increasing  $\alpha$  from 0 to 1 generally leads to a decrease in forecast error, although a marginal increase in error is noted as  $\alpha$  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  $\alpha$  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.

## F 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.

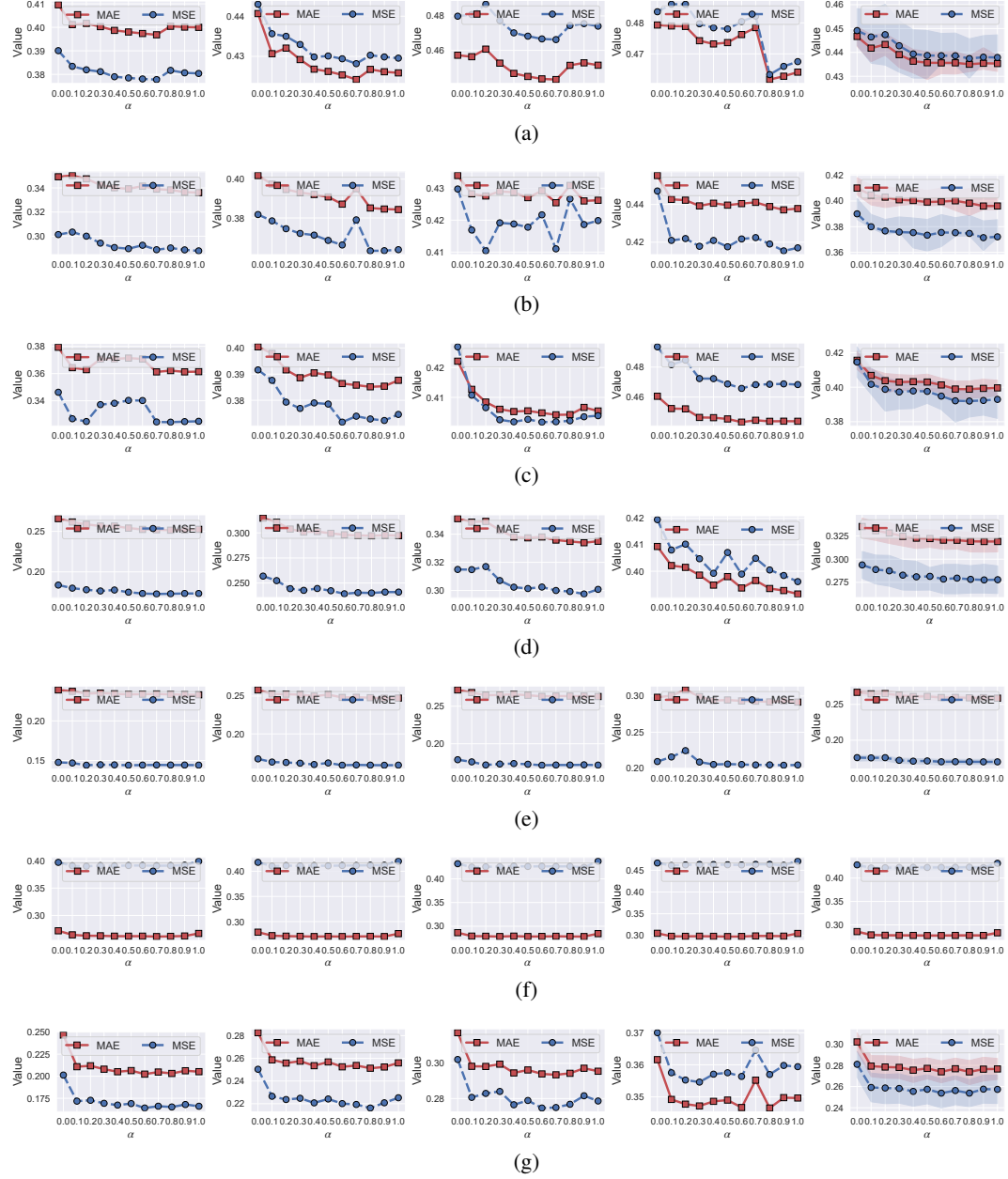


Figure 18: Performance of iTransformer enhanced by FreDF given different relative importance of frequency loss  $\alpha$ . These experiments are conducted on ETTh1 (a), ETTh2 (b), ETTm1 (c), ETTm2 (d), ECL (e), Traffic (f) and Weather (g) datasets. Different columns correspond to different forecast lengths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



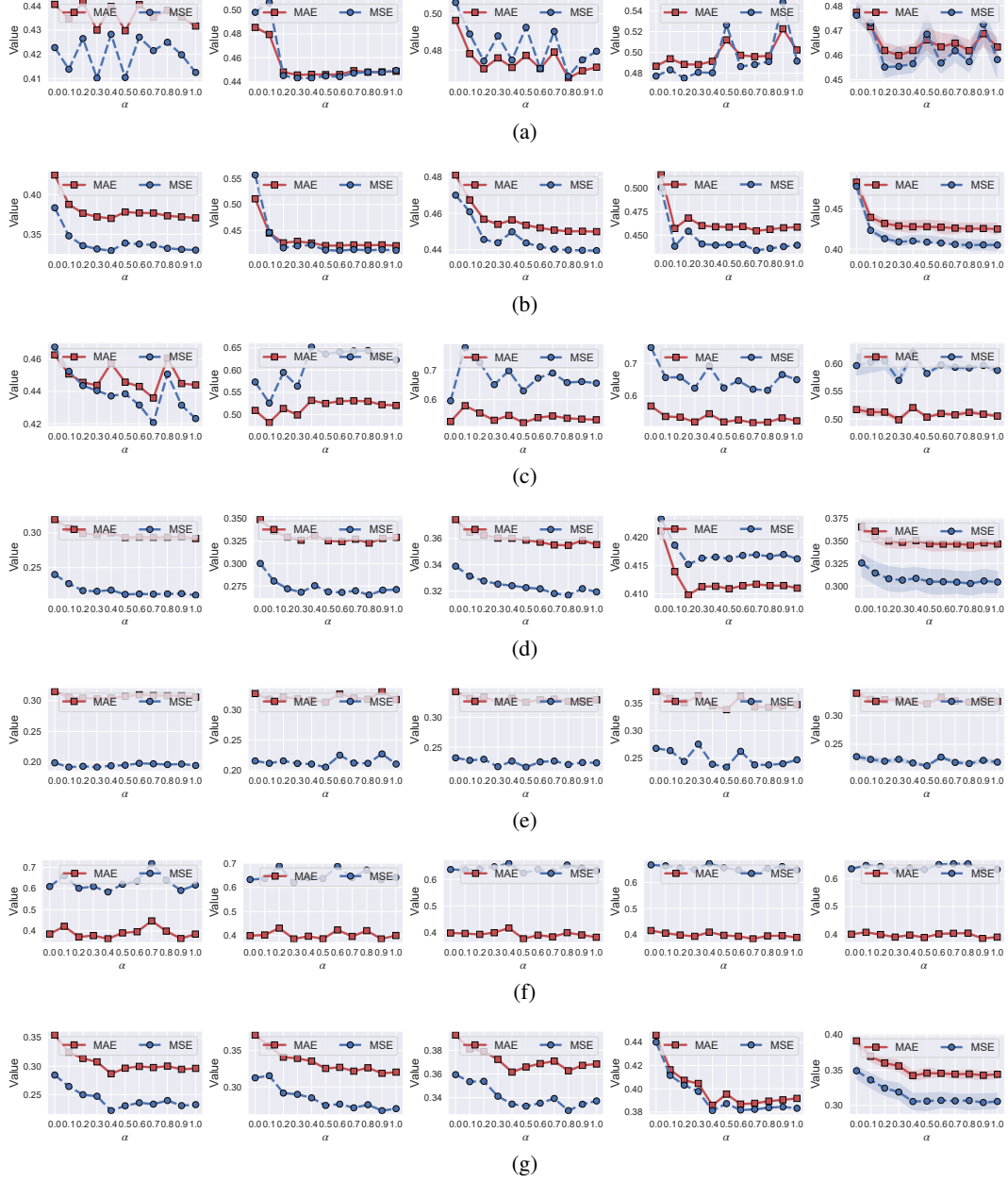


Figure 19: Performance of Autoformer enhanced by FreDF given different relative importance of frequency loss  $\alpha$ . These experiments are conducted on ETTh1 (a), ETTh2 (b), ETTm1 (c), ETTm2 (d), ECL (e), Traffic (f) and Weather (g) datasets. Different columns correspond to different forecast lengths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



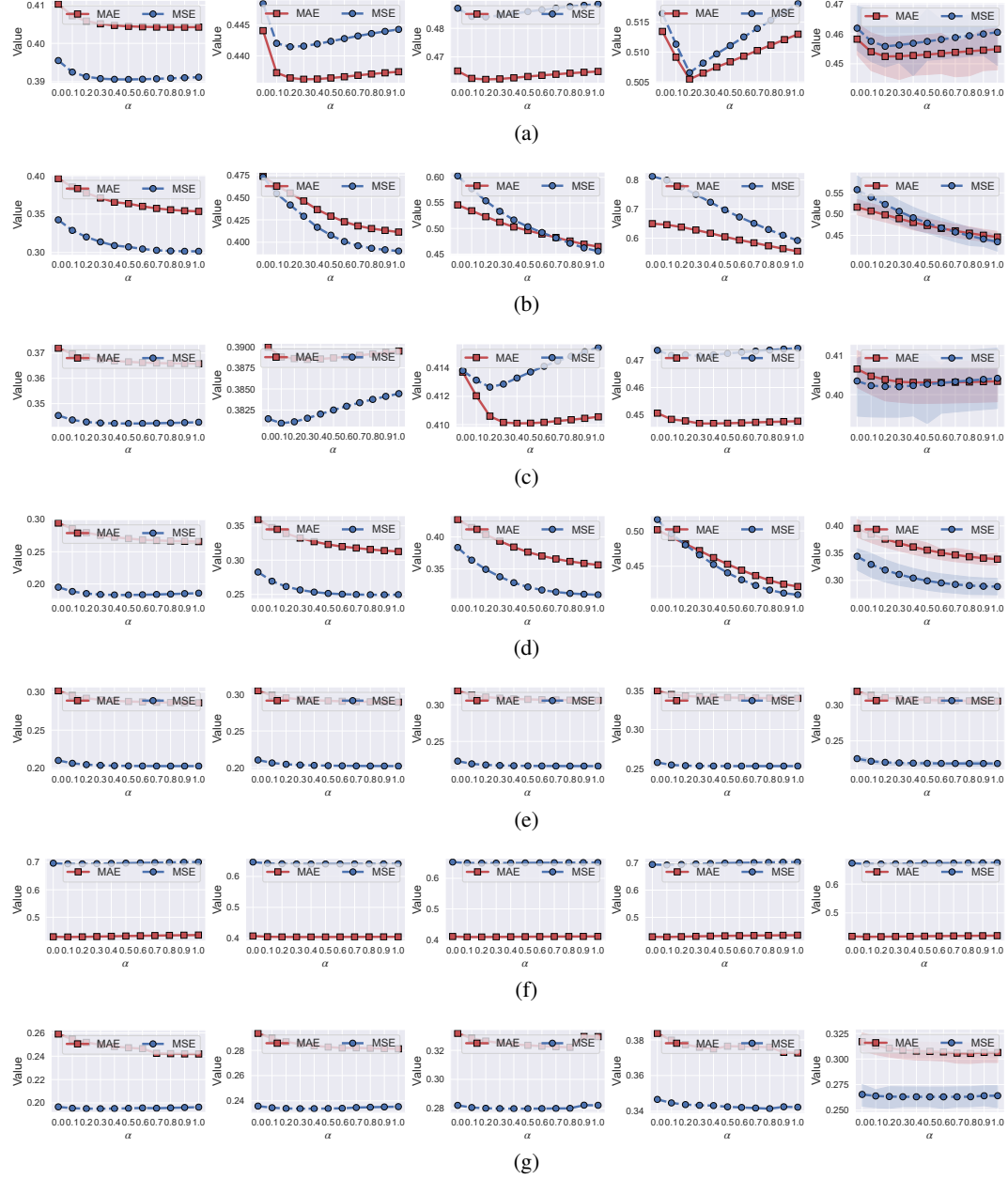


Figure 20: Performance of DLinear enhanced by FreDF given different relative importance of frequency loss  $\alpha$ . These experiments are conducted on ETTh1 (a), ETTh2 (b), ETTm1 (c), ETTm2 (d), ECL (e), Traffic (f) and Weather (g) datasets. Different columns correspond to different forecast lengths  $T$  (from left to right: 96, 192, 336, 720, and their average with shaded areas being 50% confidence intervals).