# A DYNAMIC MULTISCALE ANTI-ALIASING NETWORK FOR TIME SERIES FORECASTING

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

035

037

038

040

041

042

043

044

046 047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

Real-world time series inherently exhibit complex temporal patterns. Within chaotic systems, significant mixing and entanglement occur between different time-varying modes. Given that time series exhibit distinctly different patterns at various sampling scales, downsampling to extract multiscale features is a common approach. However, conventional downsampling causes high-frequency components in the original signal, those exceeding the new Nyquist frequency, to undergo spectral folding. This erroneously introduces spurious low-frequency patterns, perceived as low-frequency noise, thereby leading to the *aliasing problem*. To address this problem, we propose a Decomposition-Prevention-Fusion architecture framework called **DMANet**, which introduces the Dynamic Multiscale Anti-Aliasing Network. Specifically, DMANet comprises two key components: Multiscale Convolutional Downsampling, designed to capture temporal dependencies and inter-channel interactions, and an Anti-Aliasing Operation, which includes Pre-Sampling Anti-Aliasing Filtering and Post-Sampling Interpolation. These designs guarantee the fidelity of multiscale features before and after downsampling. We show that by mitigating the risk of aliasing, our proposed simple convolutional downsampling architecture achieves performance competitive with common baselines and larger Transformer-based models prevalent in existing studies across multiple benchmark datasets. Our codes are available at https://anonymous.4open.science/r/DMANet-ED7A.

# 1 Introduction

Time series analysis is widely applied in various fields such as health Morid et al. (2023), economics Sezer et al. (2020), transportation Shu et al. (2021), and weather Volkovs et al. (2024). With the widespread adoption of physical and virtual sensors, vast amounts of time series data are continuously generated, offering unprecedented opportunities for in-depth analysis and modeling. In contrast to image, video, and text data, which often possess defined syntax or intuitive patterns, time series data consist of scalar values continuously recorded at each time point. Semantic information in time series data is mainly derived from temporal changes Wu et al. (2023).

The complexity and non-stationarity inherent in real-world systems mean that observed time series often exhibit intricate temporal patterns (e.g., ascents, descents, fluctuations, sudden drifts). These patterns interact, and such interactions become particularly pronounced in chaotic systems, where significant overlap and aliasing can occur Wu et al. (2024). This challenge intensifies when distinct temporal patterns emerge at multiple scales, resulting in the entanglement of various temporal variations Shang et al. (2024) Kou et al. (2025). Therefore, time series analysis must carefully consider the intricate interactions and dynamic relationships among temporal patterns.

To address the complex time-varying entanglement in time series, an increasing number of studies focus on leveraging prior knowledge to decompose time series into more interpretable and simpler components that provide a basis for forecasting. For example, models such as Autoformer Wu et al. (2021) and Dlinear Zeng et al. (2023) decompose series into seasonal and trend components. Times-Net Wu et al. (2023) and Peri-midformer Wu et al. (2024) leverage the periodicity of time series by dividing long sequences into shorter segments based on period length, enabling separate modeling of inter- and intra-period dependencies. Beyond time-domain decomposition, frequency-domain analysis offers a valuable complementary perspective to understand temporal entanglement. Techniques

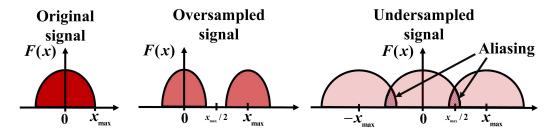


Figure 1: The illustration demonstrates the occurrence of aliasing when sampling a time series signal from the perspective of its frequency spectrum. **Left:** The frequency spectrum of a signal with maximal frequency  $x_{\text{max}}$ . **Center:** After sampling at a sufficiently high rate, replicated spectra do not overlap means that no aliasing occurs. **Right:** After undersampling, spectral replicas overlap, causing aliasing due to mixed frequency components. More details in Appendix.A.3

like the Fourier transform allow signals to be decomposed into orthogonal frequency components, where low frequencies might represent long-term periodic variations and high frequencies capture abrupt events, revealing intrinsic patterns often obscured in the time domain.

However, as time series exhibit distinct temporal patterns at varying sampling scales Wang et al. (2024a), future variations are jointly determined by the interplay of multiple scales Hu et al. (2025) Liu et al. (2025). Despite the effectiveness of the aforementioned methods in decomposing specific aspects, modeling complex time-varying entanglement remains a critical challenge. Increasingly, multiscale decomposition approaches, exemplified by TimeMixer Wang et al. (2024a), aim to model multiscale variations by decomposing them into different temporal granularities. These methods often select downsampling operations, progressively reducing temporal resolution using techniques such as strided convolutions or pooling layers to expand the models' receptive field and capture dependencies across different scales.

However, existing downsampling processes are susceptible to critical *aliasing risks* as shown in Figure.1 (see Appendix.A for a detailed explanation). When downsampling operators such as strided convolutions or pooling are used, high-frequency components of the original signal that exceed the new Nyquist frequency undergo spectral folding Shannon (1949); Nyquist (1928). If undersampled, these folded components are incorrectly represented as spurious low-frequency patterns, compromising the precision and reliability of the extracted multiscale features Chen et al. (2024a). In high-sensitivity domains such as industrial fault diagnosis Ahmed et al. (2022), such distortions and the introduction of incorrect frequency can hinder diagnostic capabilities for domain experts.

Motivated by these observations, we posit that directly addressing the aliasing problem inherent in downsampling processes, particularly within convolutional architectures, represents a key breakthrough for constructing reliable multiscale time series models capable of effectively modeling timevarying entanglement. Technically, we introduce a novel multiscale convolutional downsampling framework centered around a Decomposition-Prevention-Fusion architecture, designed to mitigate aliasing during the downsampling process. Our contributions can be summarized as follows:

- We reexamine the multiscale downsampling framework for time series from a synergistic time-frequency perspective, proposing a Decomposition-Prevention-Fusion architecture that effectively disentangles time-series features to address the challenges posed by complex time-varying entanglement.
- We introduce novel mechanisms for pre-emptive prevention and post-hoc suppression of aliasing explicitly within the multiscale decomposition process, thereby further leveraging the potential of convolutional downsampling for time-series analysis.
- Through extensive experiments, we demonstrate that our proposed method achieves stateof-the-art performance with a parameter-efficient design across multiple benchmarks.

# 2 Related Work

**Frequency-aware Models.** In time series analysis, the frequency domain can effectively capture periodic information that is difficult to represent in the time domain, thus becoming an important

complement to time domain modeling. Some methods aim to enhance time domain operations by incorporating frequency domain features as auxiliary information. For example, FEDformer performs attention weight aggregation in the frequency domain Zhou et al. (2022b). Film separates the signal from the noise in historical information through Fourier filtering Zhou et al. (2022a). Meanwhile, approaches are proposed which replace time domain input with frequency domain representations directly. FITS Xu et al. (2024) and FreTS Yi et al. (2024b) use frequency-domain MLPs for prediction, significantly reducing computational complexity. FreDF Wang et al. (2025) introduces an additional loss function in the frequency domain to supervise the alignment of the model's spectrum with the real values. However, the effectiveness of frequency domain methods is constrained by the spectrum utilization bottleneck. FilterNet Yi et al. (2024a) through simple filters demonstrates that traditional feature selection strategies in the frequency domain, such as top-K or random-K, may lead to the loss of key frequency band information. Although Fredformer Piao et al. (2024b) and proposes a frequency band equal learning mechanism and CFPT Kou et al. (2025) introduced a dual-branch architecture featuring a cross-frequency interaction module, it still does not address the issue of modeling dynamic interactions between frequency bands.

**Decomposition-based Models.** Real-world time series are often composed of various underlying patterns. To take advantage of the features of different patterns, recent methods tend to decompose the sequence into multiple subcomponents, including trend-seasonal decomposition, multiperiod decomposition, and multiscale decomposition Huang et al. (2025). Methods such as Autoformer Wu et al. (2021) and DLinear Zeng et al. (2023) use moving averages to decouple seasonal and trend components, followed by modeling with attention mechanisms or MLP layers. TimesNet Wu et al. (2023) and PDF Dai et al. (2024) utilize Fourier analysis to decouple the sequence into multiple subperiodic sequences based on computational periods. FRENet Zhang et al. (2024) introduces a frequency-based rotation network that can capture the features of dynamically complex periods. Furthermore, TimeMixer Wang et al. (2024a) uses past decomposable mixes for multiscale representation learning and future multi-prediction mixes to enhance forecasting with complementary skills. TimeStacker Liu et al. (2025) progressively stacks features from patches of varying sizes and employs a frequency-based self-attention mechanism. However, the information fidelity of multiscale decomposition is facing challenges. Downsampling operations may lead to the loss of fine-grained features due to spectral aliasing. To address the limitations, this paper proposes a multiscale decomposition framework based on frequency domain adaptive filtering, which automatically suppresses aliasing noise through frequency band masking, ensuring the integrity of multiscale feature transfer.

#### 3 Model Framework

#### 3.1 Overall Architecture

In this section, we explain the workflow of DMANet based on a single sample for clarity. The overall architecture adopts a Decomposition-Prevention-Fusion paradigm shown in Figure.2. The input sequence is initially normalized and projected into the latent space. Then, a hierarchical extractor progressively decomposes the sequence into multiscale representations through stacked depth-wise and point-wise convolutions, with gradual downsampling of temporal resolution to capture both intra-variable and inter-variable interactions, respectively. Adhering to the prevention design principle, anti-aliasing filters are performed before each downsampling step. Next, in the upsampling phase, learnable spectral filters are adopted to expand the channel, while zero-padding expands the temporal length, effectively suppressing aliasing distortions. In the fusion stage, multiscale features are integrated using Softmax, while residual connections are made between stacked encoder blocks. Finally, the hierarchical representations are decoded. This architecture fully leverages the potential of convolutional downsampling through joint time-frequency operations.

# 3.2 NORMALIZATION AND EMBEDDING

First, we apply RevIN to the input data  $X \in \mathbb{R}^{C \times L}$  to reduce the discrepancy between the training and testing data distributions Kim et al. (2022). Following this, an initial linear embedding layer re-encodes the normalized series into a latent space that is more suitable for pattern extraction and anti-aliasing. During this encoding process, we add a learnable positional encoding to preserve crucial temporal context by providing an absolute positional reference. The resulting embedded

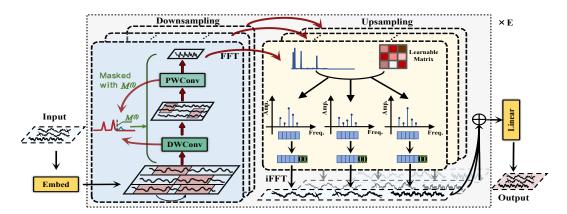


Figure 2: The overall architecture of DMANet.

representation is denoted as  $X^{'} = \text{Linear}(\text{RevIN}(X)) + W$ ,  $X^{'} \in \mathbb{R}^{C \times T}$ , where T represents the dimension of the embedded representation and W represents the positional encoding. This X' serves as the input to the subsequent multiscale extracting layers. The detailed rationale for this embedding-first approach is provided in the Appendix.G.1.

#### 3.3 Multiscale Convolutional Downsampling

To capture features at varying temporal resolutions, we process the embedded X' through a hierarchy of H downsampling layers. Unlike methods relying solely on pooling Wang et al. (2024a), we employ convolutions for efficient multiscale feature extraction. This process generates a set of downsampled feature maps  $X_{\text{down}} = \{x^{(0)}, x^{(1)}, x^{(2)}, \dots, x^{(H)}\}$ , where  $x^{(0)}$  is the initial input X' and  $x^{(l)} \in \mathbb{R}^{C_l \times T_l}$ . The temporal dimension decreases at each layer:  $T_0 = T$  and  $T_l = \lfloor T_{l-1}/s \rfloor$  for  $l \geq 1$ , with s being the fixed downsampling stride. The number of channels  $C_l$  can also vary across layers.

**Depth-wise Convolution.** In the l-th layer, the input  $x^{(l-1)} \in \mathbb{R}^{C_{l-1} \times T_{l-1}}$  first undergoes a depthwise convolution (DWConv). This operation applies distinct filters to each input channel, focusing on modeling temporal dependencies within channels without cross-channel interference:

$$f^{(l)} = \text{DWConv}(x^{(l-1)}; \text{ stride} = s, \text{ groups} = C_{l-1}) \in \mathbb{R}^{C_{l-1} \times T_l}. \tag{1}$$

**Point-wise Convolution.** Following the depth-wise convolution, a point-wise convolution (PW-Conv, that is, a  $1 \times 1$  convolution) performs a linear transformation across channels. This enhances inter-channel communication and maps the feature from  $C_{l-1}$  channels to  $C_l$  channels:

$$x^{(l)} = \text{PWConv}(f^{(l)}) \in \mathbb{R}^{C_l \times T_l}.$$
 (2)

When iterating this process up to the H-th layer, we can obtain the produced multiscale feature set  $X_{\mathrm{down}} = \{x^{(0)}, x^{(1)}, x^{(2)}, \dots, x^{(H)}\}$ . This design efficiently separates the learning of temporal patterns and channel interactions while significantly reducing parameters and computation, providing rich hierarchical information for subsequent interpolation and fusion.

## 3.4 ANTI-ALIASING OPERATION

In Section.3.3, the depth-wise convolution and point-wise convolution are proposed. To reduce the negative effect of aliasing, a Pre-Sampling Filtering and Post-Sampling Interpolation should be performed before feeding  $x^{(l-1)}$  and after acquiring  $x^{(l)}$ , respectively.

**Pre-Sampling Filtering.** Downsampling inevitably introduces the risk of aliasing, where high-frequency components fold into lower frequency bands, potentially corrupting the signal or losing critical information, especially with larger strides s. Inspired by the work in computer visionGrabinski et al. (2022a)Chen et al. (2024a), we introduce the concept of Equivalent Sampling

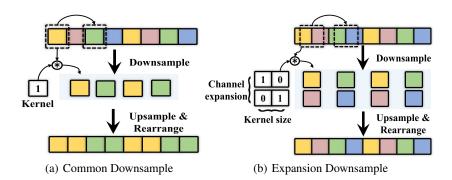


Figure 3: Illustration of how a larger kernel and channel expansion preserve sampling information during downsampling. **Left:** Pointwise downsampling with a stride of 2 and no channel expansion results in a sampling rate of 1/2, discarding half the input information. This leads to aliasing, where high-frequency content is misrepresented as low frequencies. **Right:** Downsampling with 2×1 identity kernels and 2× channel expansion ensures that every input element is sampled and preserved in separate channels. This approach maintains the effective sampling rate at 1,as all pixels are sampled.

Rate (ESR) to dynamically compute the appropriate Nyquist frequency for anti-aliasing filtering during downsampling. We provide a detailed proof for ESR in Appendix.B. As shown in Figure.3, the size of the convolutional kernel and the transformation of channels play a role in determining the sampling ability. Concretely, the ESR at the l-th layer can be calculated with the following strategy:

$$ESR^{(l)} = \frac{\min\left(K, C_l/C_{l-1}\right)}{s},\tag{3}$$

where K is the kernel size of the depth-wise convolution,  $C_{l-1}$  and  $C_l$  are the input and output channels for the layer's point-wise convolution, and s is the stride. Before applying the down-sampling convolution, we use FFT  $(\mathcal{F})$  to transform the input  $x^{(l-1)}$  into the frequency domain:  $\mathcal{X}^{(l-1)} = \mathcal{F}(x^{(l-1)})$ . The layer-specific Nyquist frequency is  $f_{\text{Nyquist}}^{(l)} = \text{ESR}^{(l)}/2$ . Based on  $f_{\text{Nyquist}}^{(l)}$ , we construct a low-pass frequency mask  $\mathbf{M}^{(l)}$ :

$$\mathbf{M}^{(l)}[i] = \begin{cases} 1, & f_i \le f_{\text{Nyquist}}^{(l)} \\ 0, & f_i > f_{\text{Nyquist}}^{(l)} \end{cases}, \tag{4}$$

where  $f_i$  is the *i*-th frequency compotent in  $\mathcal{X}^{(l-1)}$ . Then, we apply the mask in an element-wise strategy  $(\odot)$ , and transform back using IFFT  $(\mathcal{F}^{-1})$ :

$$\tilde{\mathcal{X}}^{(l-1)} = \mathcal{X}^{(l-1)} \odot \mathbf{M}^{(l)}, \ x_{\text{filtered}}^{(l-1)} = \mathcal{F}^{-1}(\tilde{\mathcal{X}}^{(l-1)}). \tag{5}$$

This filtered signal  $x_{\text{filtered}}^{(l-1)}$ , which serves as the real  $x^{(l-1)}$ , is then fed into the depth-wise convolution, effectively suppressing high frequencies prone to aliasing during downsampling.

**Post-Sampling Interpolation.** After obtaining the downsampled feature  $x^{(l)} \in \mathbb{R}^{C_l \times T_l}$ , we propose an anti-aliasing interpolation method and frequency domain channel expansion to restore the temporal resolution to a target length T and expand channels to the model dimension C.

First, the downsampled feature is transformed to the frequency domain using  $\mathcal{F}$ . Inspired by frequency domain filtering strategies Yi et al. (2024a), we introduce our designed learnable complex-valued filters  $\mathcal{H}_{\phi}^{(l)} \in \mathbb{C}^{C \times C_l \times F_l}$ , which are also transformed to the frequency domain. We then perform channel expansion through a weighted sum in the frequency domain, effectively implementing channel mixing. This operation computes each output channel c as a learned combination of all  $C_l$  input channels at each frequency f, fusing cross-channel information while preserving the spectral structure:

$$\mathcal{X}^{(l)} = \mathcal{F}(x^{(l)}) \in \mathbb{C}^{C_l \times F_l}, \ F_l = \lfloor T_l/2 \rfloor + 1, \tag{6}$$

$$\mathcal{S}^{(l)}[c,f] = \sum_{k=1}^{C_l} \mathcal{X}^{(l)}[k,f] \odot \mathcal{H}_{\phi}^{(l)}[c,k,f], \text{ for } c \in [1,C] \text{ and } f \in [1,F_l],$$
 (7)

where  $\mathcal{S}^{(l)} \in \mathbb{C}^{C \times F_l}$ . To restore the sequence length to T, we calculate the corresponding target frequency count  $F = \lfloor T/2 \rfloor + 1$ . We apply zero-padding (Pad) in the frequency domain to extend the spectrum  $\mathcal{S}^{(l)}$  from  $F_l$  components to F components. This zero-padding primarily serves to interpolate the signal in the time domain while implicitly acting as a low-pass filter, further mitigating potential aliasing introduced during the process. Finally,  $\mathcal{F}^{-1}$  transforms the padded spectrum back to the time domain, producing the interpolated feature map for level l:

$$\tilde{\mathcal{S}}^{(l)} = \operatorname{Pad}^{(l)}(\mathcal{S}^{(l)}) \in \mathbb{C}^{C \times F}, \ Y^{(l)} = \mathcal{F}^{-1}(\tilde{\mathcal{S}}^{(l)}, n = T) \in \mathbb{R}^{C \times T}.$$
(8)

Executing channel expansion before zero-padding ensures that the spectral information for each expanded channel is complete before interpolation, avoiding potential spectral leakage or distortion and contributing to high-fidelity reconstruction in multi-channel time series tasks.

#### 3.5 Multiscale Feature Fusion and Decoding Output

**Feature Fusion.** The above operations are processed in a single encoder block, and the interpolation step generates a set of feature maps  $\{Y^{(1)},Y^{(2)},\ldots,Y^{(H)}\}$ , each residing at the target resolution T, but derived from different temporal scales. To integrate these multiscale information, we employ adaptive weighting. A learnable weight vector  $w \in \mathbb{R}^H$  is introduced, and its Softmax normalization yields attention scores  $\alpha_p$  for each scale. The final output  $\hat{Y}^e$  of the e-th encoder block is the weighted sum of these multiscale features:

$$\alpha_p = \frac{\exp(w_p)}{\sum_{k=1}^H \exp(w_k)}, \hat{Y}^e = \sum_{p=1}^H \alpha_p \cdot Y^{(p)} \in \mathbb{R}^{C \times T}.$$
 (9)

Our model stacks E such multiscale encoder blocks. To facilitate the training of this deep architecture and preserve information flow, we incorporate residual connections around each encoder block, where  $\hat{Y}^{(0)} = X'$  means the initial embedded representation, and  $\hat{Y}^e$  is the output of the e-th encoder layer:

$$\hat{Y}^e = \text{MultiScaleEncoder}(\hat{Y}^{e-1}) + \hat{Y}^{e-1}, \text{ for } e = 1, 2, 3, \dots, E.$$
 (10)

**Decoding Output.** The output  $\hat{Y}^E$  from the final encoder block, representing rich multiscale features, first passes through a Layer Normalization step,  $Y^* = \text{LayerNorm}(\hat{Y}^E)$ . Then a simple Feed-Forward Network (FFN) decoder projects these features into the future prediction horizon  $L_{\text{next}}$ . Then we apply the inverse RevIN transformation (iRevIN) to obtain the final forecast  $\hat{X}_{O}$ :

$$\hat{X}_o = iRevIN(FFN(Y^*)) \in \mathbb{R}^{C \times L_{next}}.$$
(11)

#### 4 EXPERIMENTS

#### 4.1 EXPERIMENTAL SETTINGS

**Datasets.** We conduct experiments on many real-world public datasets for long-term forecasting: ETT (four subsets), Weather, ECL, Solar-Energy, PEMS (four subsets). For short-term forecasting, we adopt the ILI, COVID-19, NASDAQ, Wiki, SP500, DowJones, CarSales, Power, Website, Unemp. These datasets are standard benchmarks Wang et al. (2024a) Liu et al. (2024) Yue et al. (2025). Details and statistics of these multivariate time series datasets are summarized in Appendix.C.

**Baseline.** Our primary analysis focuses on long-term forecasting with a 96-step lookback window (Tables 9, 10 and Table.13). In addition to this main task, we also conducted evaluations on univariate long-term forecasting (Table 14), short-term forecasting (Table 15), and long-term forecasting with an extended 720-step lookback window (Tables 11 and 12). Across these diverse settings, we chose a comprehensive set of recent state-of-the-art models to serve as baselines. This includes MLP-based models (SOFTS Han et al. (2024), TimeMixer Wang et al. (2024a), DLinear Zeng et al. (2023)), CNN-based models (TVNet Li et al. (2025), ModernTCN Donghao & Xue (2024), PDF Dai et al. (2024)), frequency-based models (TimeStacker Liu et al. (2025), FreDF Wang et al. (2025), etc.), Transformer-based models (TimeXer Wang et al. (2024b), iTransformer Liu et al. (2024), etc.), and recent architectures based on Mamba (TimePro Ma et al. (2025)), KAN (TimeKAN Huang et al. (2025)) and Retrieval-Augmented (RAFT Han et al. (2025)), among others. Detailed descriptions are provided in Appendix C.

Table 1: Long-term forecasting results (L=96). All results are averaged across four forecasting horizon:  $T \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively. See Appendix.D (Table.9 and Table.10) for full results.

Models		ANet irs		former 124		Mixer 24a		rNet 24a		ormer 24a	FI 20		Fre 202		Time 20		Fre 20	DF 25	SOI 20			eXer 24b
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	0.428	0.429	0.454	0.447	0.447	0.440	0.440	0.432	0.445	0.432	0.447	0.448	0.488	0.474	0.438	0.438	0.437	0.435	0.449	0.442	0.437	0.437
ETTh2	0.361	0.388	0.383	0.407	0.364	0.395	0.378	0.397	0.367	0.396	0.383	0.408	0.550	0.515	0.377	0.403	0.371	0.396	0.373	0.400	0.368	0.396
ETTm1	0.373	0.385	0.407	0.410	0.381	0.395	0.384	0.398	0.393	0.403	0.387	0.408	0.407	0.415	0.391	0.400	0.392	0.399	0.393	0.403	0.382	0.397
ETTm2	0.268	0.310	0.288	0.332	0.275	0.323	0.276	0.322	0.279	0.324	0.286	0.328	0.335	0.379	0.281	0.326	0.278	0.319	0.287	0.330	0.274	0.322
Weather	0.236	0.262	0.258	0.279	0.240	0.271	0.248	0.278	0.246	0.272	0.249	0.276	0.255	0.363	0.251	0.276	0.254	0.274	0.255	0.278	0.241	0.271
Electricity	0.170	0.264	0.178	0.270	0.182	0.272	0.201	0.285	0.175	0.269	0.217	0.295	0.202	0.290	0.169	0.262	0.170	0.259	0.174	0.264	0.171	0.270
Solar-Energy	0.227	0.249	0.233	0.262	0.216	0.280	0.263	0.286	0.232	0.274	0.397	0.398	0.283	0.338	0.232	0.266	0.279	0.292	0.229	0.256	0.237	0.302

**Implementation Details.** The experiments in this paper were conducted using an NVIDIA GeForce RTX 3090 24GB GPU. Inspired by FreDF Wang et al. (2025), we uses the Mean Absolute Error (MAE) in the frequency domain. For details on the hyperparameter settings of the models presented in Appendix.C.

#### 4.2 MAIN RESULTS

Long-term Forecasting. The long-term forecasting results, reported in Table.1 (more results in Appendix.D), demonstrate that DMANet consistently achieves optimal or near-optimal performance across all datasets. Its performance is comparable to TimeMixer, highlighting the general effectiveness of time-series decomposition architectures. However, a key distinction lies in their downsampling mechanisms: while TimeMixer's reliance on average pooling is susceptible to information loss, DMANet's spectral preservation mechanism effectively suppresses aliasing artifacts, enabling a more faithful layer-wise learning of multi-granularity representations. Conversely, when compared to models with channel-wise self-attention like iTransformer, DMANet's reliance on simpler convolutional operations for dependency modeling suggests a potential area for future optimization, particularly on high-dimensional datasets. Furthermore, guided by the principles of scaling laws in Time Series Forecasting (TSF), we extended the lookback window *L* to 720 in Table.2 (full results can be found in Table.11 and Table.12). In this long-context setting, DMANet exhibits robust noise resilience, maintaining state-of-the-art performance and surpassing other convolutional counterparts like ModernTCN and TVNet. This result further validates DMANet's superior adaptability in capturing multi-scale temporal dependencies, even with extended input lengths.

Table 2: Long-term forecasting results (L=720). For baseline, the input length L is searched from  $\{192, 336, 512, 720\}$ , while DMANet is fixed 720. The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively. See Appendix.D (Table.11 and Table.12) for full results.

Models	DMANet	iTrans	former	Time	Mixer	Mode	rnTCN	TV	Net	TSL.	ANet	PI	)F	Patcl	nTST	FI	TS	Time	sNet	DLi	inear
Metric	MSE MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	0.338 0.369	0.361	0.390	0.356	0.380	0.351	0.381	0.348	0.379	0.348	0.383	0.342	0.376	0.349	0.381	0.357	0.377	0.408	0.415	0.356	0.378
ETTm2	0.248 0.307	7   0.269	0.327	0.257	0.318	0.253	0.314	0.256	0.316	0.256	0.316	0.250	0.313	0.256	0.314	0.254	0.313	0.292	0.331	0.259	0.324
Weather	0.218 0.252	2   0.232	0.270	0.226	0.264	0.224	0.264	0.221	0.261	0.325	0.337	0.227	0.263	0.224	0.261	0.244	0.280	0.255	0.282	0.242	0.293
Electricity	0.154 0.252	2   0.163	0.258	0.169	0.265	0.156	0.253	0.165	0.254	0.165	0.257	0.160	0.253	0.171	0.270	0.169	0.265	0.190	0.290	0.167	0.264

**Short-term Forecasting.** The short-term forecasting results, presented in Table.3, validate the superiority of DMANet in handling non-stationary time series. Across a diverse set of challenging datasets such as ILI, COVID-19, and DowJones, DMANet consistently achieves the best performance. It significantly outperforms other methods, including strong frequency-domain baselines like Fredformer and FilterNet. These results underscore DMANet's exceptional capability in short-term and non-stationary forecasting, attributable to its synergistic design: the convolutional architecture excels at preserving local features, while the anti-aliasing structure effectively mitigates disruptive high-frequency noise.

#### 4.3 ABLATION STUDY

In this section, we investigate key components of DMANet, including our novel Anti-aliasing Filter, the Convolutional Downsampling and Frequency Upsampling Mechanisms, and the Basic Settings.

Table 3: Short-term forecasting results. The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively. See Appendix.D (Table.15) for full results and setting details.

Models	DMA	ANet	Timel	Mixer	Filter	rNet	FI	rs .	DLiı	near	Fredfo	ormer	Patch	TST
Metric	MSE	MAE												
ILI	1.763	0.824	2.020	0.878	2.073	0.885	4.130	1.465	3.083	1.217	1.947	0.899	2.128	0.885
COVID-19	1.910	0.670	2.234	0.782	2.088	0.780	2.875	0.979	3.483	1.102	1.902	0.765	2.221	0.820
NASDAQ	0.177	0.273	0.186	0.281	0.197	0.289	0.210	0.302	0.228	0.331	0.194	0.285	0.198	0.286
Wiki	6.506	0.393	6.572	0.409	6.572	0.411	8.515	0.553	6.634	0.481	6.705	0.406	6.523	0.404
SP500	0.225	0.329	0.241	0.353	0.254	0.365	0.291	0.412	0.277	0.391	0.261	0.378	0.246	0.361
DowJones	11.957	0.850	13.948	0.877	13.439	0.873	13.755	0.893	12.688	0.857	12.992	0.858	12.916	0.862
CarSales	0.338	0.333	0.338	0.336	0.336	0.335	0.379	0.365	0.373	0.368	0.340	0.338	0.338	0.335
Power	1.373	0.899	1.484	0.937	1.614	0.986	1.711	1.028	1.549	0.972	1.588	0.981	1.650	0.998
Website	0.137	0.252	0.143	0.261	0.136	0.255	0.278	0.383	0.204	0.319	0.135	0.254	0.141	0.259
Unemp	0.064	0.146	0.094	0.183	0.079	0.166	0.308	0.394	0.154	0.292	0.075	0.163	0.078	0.160

Table 4: Ablation study of DMANet. All results are averaged across four different forecasting horizon. The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

Cata	agorie	s	1	Do	wnsamj	oling Re	place			Up	samplii	ng Repl	ace			Basic S	Settings	;
Cases	DN	/IANet	Linear	r Down	Self-A	ttention	Standa	ırd Conv	Line	ar Up	Inter	olate	Trans	Conv	w/o F	leVIN	MSE	Loss
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Electricity	0.17	2 0.265	0.176	0.271	0.180	0.276	0.179	0.274	0.174	0.268	0.184	0.274	0.179	0.271	0.215	0.316	0.173	0.268
ETTm1	0.37	3 0.385	0.379	0.389	0.379	0.389	0.377	0.388	0.377	0.388	0.377	0.388	0.378	0.388	0.423	0.445	0.381	0.393
Unemp	0.06	4 0.140	0.081	0.171	0.076	0.161	0.075	0.163	0.077	0.164	0.073	0.161	0.068	0.155	0.759	0.414	0.076	0.166
NASDAQ	0.17	7 0.273	0.190	0.283	0.184	0.279	0.185	0.281	0.183	0.279	0.182	0.277	0.184	0.278	2.337	1.132	0.195	0.288

**Basic Settings.** Ablation analysis showed that removing DMANet's ReVIN significantly hurts performance by failing to mitigate distribution shift. The frequency-domain MAE loss is also preferred over MSE for anti-aliasing due to enabling direct frequency adjustment.

Convolution Downsampling. To evaluate the effectiveness of convolutional downsampling, we experimented with the following alternative strategies: (1) LinearDown: two separate linear for downsampling; (2) Standard Conv: standard convolution with stride; (3) Self-attention: employing self-attention to capture temporal dependencies, combined with average pooling along the temporal and convolutional downsampling along the channel. The results are summarized in Table.4.

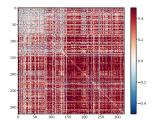
Overall, the combination of depth-wise convolution and point-wise convolution demonstrates the best performance. Notably, replacing convolution with linear or with self-attention followed by average pooling results in a performance drop, which highlights the capability of depthwise convolution in learning temporal dependencies. In addition, using standard convolution alone leads to a substantial increase in parameter count and a worsening of most metrics, suggesting that focusing solely on depthwise convolution to extract temporal dependencies is a more reasonable design.

Table 5: Ablation study on different filter designs. Performance is compared against our DMANet (utilizing the ESR filter), heuristic filters (Max, Random), and classical filters (Ideal, Chebyshev, Gaussian, Butterworth). All results are averaged over four horizons. Details are in Appendix.C.5.

Models	DMA	Net	Ma	ax	Rano	dom	Ide	al	Cheby	yshev	Gaus	sian	Butter	worth	w/o l	ESR
Metric	MSE	MAE														
ILI	1.763	0.824	1.957	0.855	2.043	0.872	1.994	0.849	1.990	0.855	1.974	0.862	1.940	0.849	1.830	0.835
Unemp	0.064	0.146	0.073	0.157	0.074	0.159	0.073	0.159	0.075	0.164	0.071	0.154	0.072	0.158	0.075	0.161
DowJones	11.957	0.850	12.300	0.852	12.261	0.851	12.397	0.855	12.382	0.855	12.351	0.854	12.402	0.855	12.379	0.853
ETTm1	0.373	0.385	0.376	0.387	0.376	0.387	0.375	0.387	0.375	0.387	0.374	0.385	0.375	0.387	0.376	0.387
PEMS08	0.090	0.198	0.117	0.218	0.113	0.210	0.109	0.206	0.110	0.207	0.108	0.205	0.109	0.206	0.110	0.207

Frequency Upsamping. To validate the unique advantages of frequency-domain upsampling, we conducted experiments comparing our approach with three alternative upsampling methods that do not explicitly target frequency information: (1) Linear Up: two separate linear for upsampling; (2) Interpolate: simply interpolation along the temporal and channel; (3) Trans Conv: utilizes transposed convolution mirroring the structure of the downsampling counterpart. As shown clearly in Table.4, replacing our frequency-domain upsampling with any of these alternatives resulted in a significant performance degradation. This indicates that these methods fail to effectively preserve or reconstruct the crucial frequency components of time series during the upsampling process.

In contrast, our strategy first performs expansion in the channel dimension, followed by high-frequency truncation in the frequency domain. This carefully designed approach ensures structural integrity and independence of each channel in the frequency domain and completely avoids the spectral leakage and the aliasing problem inherent to interpolation-



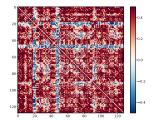


Figure 4: Visualization of dependency differences, comparing feature representations with and without the anti-aliasing filter. Red indicates an increase, while blue indicates a decrease. (a) Left: Channel-wise dependency. (b) Right: Temporal dependency differences.

based methods. As a result, our method demonstrates outstanding performance in the high-fidelity reconstruction of time series.

Anti-aliasing Filter. To validate our ESR-based filter across a wide range of real-world scenarios, we benchmarked DMANet against two categories of alternatives on six diverse datasets: simple heuristic filters and established classical filters. The results in Table.5 show that on the ILI, Unemp, Dowjone and PEMS08 datasets, DMANet outperforms all competing filters. This efficacy is rooted in our model's unique architecture-aware design. Unlike static filters, DMANet dynamically adapts to the varying signal processing capabilities of each layer in its multi-scale pipeline. It automatically calculates a precise, optimal cutoff frequency for every downsampling operation, ensuring anti-aliasing remains theoretically grounded throughout the network. This exposes the core limitations of alternatives: heuristic filters are arbitrary and lack a principled basis, while classical filters are static and would require a tedious, layer-specific tuning process to be effective.

#### 4.4 MODEL ANALYSIS

Efficiency and Robustness. We provide comprehensive results on real-world datasets (Appendix.E, Table.21), including the Speed, Memory and MACs. DMANet demonstrates a balance between performance and efficiency. Compared to other models, it requires fewer MACs and less memory while achieving better accuracy. Meanwhile, on synthetic data (Appendix.E, Table.20), we performed a fine-grained analysis of computational costs. An ablation study revealed that our anti-aliasing filter is not a performance bottleneck, introducing negligible overhead (a worst-case latency increase of only 2.4%). The robustness of our model was validated through a series of noise injection experiments, detailed in Appendix.G. We introduced five types of synthetic noise (e.g., high-frequency) at various intensities  $\epsilon$ . Our anti-aliasing architecture demonstrated great resilience, as its performance degraded gracefully with increasing noise, thereby validating its ability to effectively mitigate signal disturbances.

**Dependency Modeling.** Figure.4 presents feature dependency heatmaps from the ECL dataset, which reveal the effect of anti-aliasing filter. The filter smooths fine-grained dependencies that are susceptible to aliasing during downsampling. By suppressing these potentially noisy or misleading correlations, the filtering process accentuates the underlying structural patterns in both the temporal and channel. It allows the subsequent layers to more easily extract stable and meaningful features from a cleaner, more coherent representation. Detailed dependency passing analysis in Appendix.H.

# 5 CONCLUSION

This paper presents DMANet, a novel architecture that tackles the critical aliasing problem in multiscale time series forecasting through a Decomposition-Prevention-Fusion framework, employing pre-sampling anti-aliasing based on Equivalent Sampling Rate and post-sampling interpolation for high-fidelity features. Extensive experiments on diverse benchmarks demonstrate DMANet's state-of-the-art performance and robustness, validating the significance of the anti-aliasing design. DMANet offers a promising direction by explicitly integrating signal processing principles to enhance time series analysis robustness.

# 6 ETHICS STATEMENT

Our research is primarily foundational, focusing on a technical challenge within time series analysis, i.e., the problem of aliasing in multiscale deep learning models. We have considered the ethical implications of our work and believe that it follows the scientific standards.

#### 6.1 SOCIETAL IMPACT

The primary goal of DMANet is to improve the fidelity and reliability of time series forecasting models by mitigating the spectral distortion caused by aliasing. This has a positive social impact by enhancing the trustworthiness of predictive systems in high-sensitivity domains. For example, in industrial fault diagnosis, preventing the introduction of spurious frequency patterns can improve the accuracy of diagnostic tools and support expert decision-making. Similarly, more reliable models are beneficial in fields such as economics, transportation planning, and weather forecasting. We acknowledge that, like any advanced forecasting technology, our methods could potentially be misused. However, our work is a general-purpose technical improvement, not an application-specific tool, and we advocate for its responsible use in future research and applications.

#### 6.2 Data Usage

All experiments were carried out on publicly available and well-established benchmark datasets, e.g., ETT, Weather, ECL, PEMS, ILI, COVID-19. A complete list and description of these datasets are provided in the Appendix.C. In this study, no sensitive or private user data was used, thus avoiding concerns related to privacy and data protection.

#### 6.3 BIAS AND FAIRNESS

Although our work does not directly address dataset bias, it contributes to model fairness by tackling a source of technical error. By avoiding the aliasing problem, our model is less likely to learn from misleading artifacts in the data. This enhances the model's robustness and ensures its predictions are based on a more faithful representation of the underlying signal, which is a prerequisite for fair and reliable decision-making.

#### 7 REPRODUCIBILITY STATEMENT

#### 7.1 CODE

The complete source codes for DMANet, including model implementations and scripts to reproduce experimental results, are available in our anonymous repository at https://anonymous.4open.science/r/DMANet-ED7A. The repository includes instructions for setting up the environment, preparing the data, and running the training and evaluation scripts. Upon acceptance, the repository will be made public and accessible.

# 7.2 Datasets

Our study utilizes multiple publicly available real-world datasets for long-term and short-term fore-casting, including ETT, Weather, ECL, Solar-Energy, PEMS, ILI, COVID-19, and others. Detailed descriptions, statistics, sources, and data-splitting protocols (Train/Validation/Test ratios) for each dataset are provided in Appendix.C. The data processing follows the established protocols of previous benchmark studies to ensure fair comparison.

#### 7.3 EXPERIMENTAL SETUP

The Section.4.1 outlines the overall setup, while Appendix.C provides more implementation details, including the hyperparameter search spaces for all tasks, the specific configurations for each dataset, detailed descriptions of the baseline models, and the fair comparison settings. All experiments were conducted on a single NVIDIA GeForce RTX 3090 24GB GPU using the PyTorch framework

with the same version. In addition, the full results shown in Appendix.D complement the summary tables in the main text. Furthermore, we rigorously validate our results through statistical tests on experiments conducted with five random seeds, confirming that DMANet's superior performance is statistically significant with 99% confidence.

# REFERENCES

- Imran Ahmed, Gwanggil Jeon, and Francesco Piccialli. From artificial intelligence to explainable artificial intelligence in industry 4.0: A survey on what, how, and where. *IEEE Transactions on Industrial Informatics*, 18(8):5031–5042, 2022. doi: 10.1109/TII.2022.3146552.
- Linwei Chen, Lin Gu, and Ying Fu. When semantic segmentation meets frequency aliasing. In *The Twelfth International Conference on Learning Representations*, 2024a. URL https://openreview.net/forum?id=SYBdkHcXXK.
- Peng Chen, Yingying Zhang, Yunyao Cheng, Yang Shu, Yihang Wang, Qingsong Wen, Bin Yang, and Chenjuan Guo. Pathformer: Multi-scale transformers with adaptive pathways for time series forecasting. *arXiv preprint arXiv:2402.05956*, 2024b.
- Tao Dai, Beiliang Wu, Peiyuan Liu, Naiqi Li, Jigang Bao, Yong Jiang, and Shu-Tao Xia. Periodicity decoupling framework for long-term series forecasting. *International Conference on Learning Representations*, 2024.
- Abhimanyu Das, Weihao Kong, Andrew Leach, Shaan Mathur, Rajat Sen, and Rose Yu. Long-term forecasting with tide: Time-series dense encoder. *arXiv preprint arXiv:2304.08424*, 2023.
- Luo Donghao and Wang Xue. ModernTCN: A modern pure convolution structure for general time series analysis. *International Conference on Learning Representations*, 2024.
- Emadeldeen Eldele, Mohamed Ragab, Zhenghua Chen, Min Wu, and Xiaoli Li. Tslanet: Rethinking transformers for time series representation learning. In *International Conference on Machine Learning*, 2024.
- Julia Grabinski, Steffen Jung, Janis Keuper, and Margret Keuper. Frequencylowcut pooling plug and play against catastrophic overfitting. pp. 36–57, Berlin, Heidelberg, 2022a. Springer-Verlag. ISBN 978-3-031-19780-2. doi: 10.1007/978-3-031-19781-9\_3. URL https://doi.org/10.1007/978-3-031-19781-9\_3.
- Julia Grabinski, Steffen Jung, Janis Keuper, and Margret Keuper. Frequencylowcut pooling—plug & play against catastrophic overfitting. In *European Conference on Computer Vision*, 2022b. URL https://arxiv.org/abs/2204.00491.
- Lu Han, Xu-Yang Chen, Han-Jia Ye, and De-Chuan Zhan. SOFTS: Efficient multivariate time series forecasting with series-core fusion. In *Advances in Neural Information Processing Systems*, 2024.
- Sungwon Han, Seungeon Lee, Meeyoung Cha, Sercan O Arik, and Jinsung Yoon. Retrieval augmented time series forecasting. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=GUDnecJdJU.
- Yifan Hu, Peiyuan Liu, Peng Zhu, Dawei Cheng, and Tao Dai. Adaptive multi-scale decomposition framework for time series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 17359–17367, 2025.
- Songtao Huang, Zhen Zhao, Can Li, and LEI BAI. TimeKAN: KAN-based frequency decomposition learning architecture for long-term time series forecasting. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=wTLc79YNbh.
- Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y. Zhang, Xiaoming Shi, Pin-Yu Chen, Yuxuan Liang, Yuan-Fang Li, Shirui Pan, and Qingsong Wen. Time-LLM: Time series forecasting by reprogramming large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=Unb5cVPtae.

- Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2022.
  - Feifei Kou, Jiahao Wang, Lei Shi, Yuhan Yao, Yawen Li, Suguo Zhu, Zhongbao Zhang, and Junping Du. CFPT: Empowering time series forecasting through cross-frequency interaction and periodicaware timestamp modeling. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=hHYjiOJFum.
  - Chenghan Li, Mingchen Li, and Ruisheng Diao. TVNet: A novel time series analysis method based on dynamic convolution and 3d-variation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=MZDdTzN6Cy.
  - Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting long-term time series forecasting: An investigation on linear mapping, 2023a. URL https://arxiv.org/abs/2305.10721.
  - Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting long-term time series forecasting: An investigation on linear mapping. *arXiv preprint arXiv:2305.10721*, 2023b.
  - Zhe Li, Zhongwen Rao, Lujia Pan, and Zenglin Xu. Mts-mixers: Multivariate time series forecasting via factorized temporal and channel mixing. *arXiv* preprint arXiv:2302.04501, 2023c.
  - Qinglong Liu, Cong Xu, Wenhao Jiang, Kaixuan Wang, Lin Ma, and Haifeng Li. Timestacker: A novel framework with multilevel observation for capturing nonstationary patterns in time series forecasting. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=5RYSqSKz9b.
  - Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *ICLR*, 2022a. URL https://openreview.net/forum?id=0EXmFzUn5I.
  - Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. *Advances in neural information processing systems*, 35: 9881–9893, 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. *International Conference on Learning Representations*, 2024.
  - Xiaowen Ma, Zhen-Liang Ni, Shuai Xiao, and Xinghao Chen. Timepro: Efficient multivariate long-term time series forecasting with variable- and time-aware hyper-state. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=s69Ei2VrIW.
  - Mohammad Amin Morid, Olivia R. Liu Sheng, and Joseph Dunbar. Time series prediction using deep learning methods in healthcare. *ACM Trans. Manage. Inf. Syst.*, 14, 2023.
  - Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=JbdcOvTOcol.
  - Harry Nyquist. Certain topics in telegraph transmission theory. *Transactions of the American Institute of Electrical Engineers*, 47(2):617–644, 1928.
  - Xihao Piao, Zheng Chen, Taichi Murayama, Yasuko Matsubara, and Yasushi Sakurai. Fredformer: Frequency debiased transformer for time series forecasting. In *KDD*, pp. 2400–2410, 2024a. URL https://doi.org/10.1145/3637528.3671928.
  - Xihao Piao, Zheng Chen, Taichi Murayama, Yasuko Matsubara, and Yasushi Sakurai. Fredformer: Frequency debiased transformer for time series forecasting. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2400–2410, 2024b.

- Xiangfei Qiu, Xingjian Wu, Yan Lin, Chenjuan Guo, Jilin Hu, and Bin Yang. Duet: Dual clustering enhanced multivariate time series forecasting. *arXiv preprint arXiv:2412.10859*, 2024.
  - Omer Berat Sezer, Mehmet Ugur Gudelek, and Ahmet Murat Ozbayoglu. Financial time series forecasting with deep learning: A systematic literature review: 2005–2019. *Applied soft computing*, 90:106181, 2020.
  - Zongjiang Shang, Ling Chen, Binqing Wu, and Dongliang Cui. Ada-MSHyper: Adaptive multi-scale hypergraph transformer for time series forecasting. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=RNbrIQ0se8.
  - Claude Elwood Shannon. Communication in the presence of noise. *Proceedings of the IRE*, 37(1): 10–21, 1949.
  - Jingzhe Shi, Qinwei Ma, Huan Ma, and Lei Li. Scaling law for time series forecasting. *Advances in Neural Information Processing Systems*, 2024.
  - Wanneng Shu, Ken Cai, and Neal Naixue Xiong. A short-term traffic flow prediction model based on an improved gate recurrent unit neural network. *IEEE Transactions on Intelligent Transportation Systems*, 23(9):16654–16665, 2021.
  - Kristofers Volkovs, Evalds Urtans, and Vairis Caune. Primed unet-lstm for weather forecasting. In *Proceedings of the International Conference on Advances in Artificial Intelligence*, ICAAI '23, pp. 13–17, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400708985. doi: 10.1145/3633598.3633601. URL https://doi.org/10.1145/3633598.3633601.
  - Hao Wang, Lichen Pan, Yuan Shen, Zhichao Chen, Degui Yang, Yifei Yang, Sen Zhang, Xinggao Liu, Haoxuan Li, and Dacheng Tao. FreDF: Learning to forecast in the frequency domain. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=4A9IdSa1ul.
  - Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. MICN: Multiscale local and global context modeling for long-term series forecasting. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=zt53IDUR1U.
  - Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y. Zhang, and JUN ZHOU. Timemixer: Decomposable multiscale mixing for time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024a. URL https://openreview.net/forum?id=7oLshfEIC2.
  - Yuxuan Wang, Haixu Wu, Jiaxiang Dong, Guo Qin, Haoran Zhang, Yong Liu, Yunzhong Qiu, Jianmin Wang, and Mingsheng Long. Timexer: Empowering transformers for time series forecasting with exogenous variables. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b. URL https://openreview.net/forum?id=INAeUQ041T.
  - Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *CoRR*, abs/2106.13008, 2021. URL https://arxiv.org/abs/2106.13008.
  - Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *ICLR*, 2023.
  - Qiang Wu, Gechang Yao, Zhixi Feng, and Shuyuan Yang. Peri-midformer: Periodic pyramid transformer for time series analysis. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=5iUxMVJVEV.
  - Zhijian Xu, Ailing Zeng, and Qiang Xu. FITS: modeling time series with 10k parameters. In *ICLR*, 2024.

- Kun Yi, Jingru Fei, Qi Zhang, Hui He, Shufeng Hao, Defu Lian, and Wei Fan. Filternet: Harnessing frequency filters for time series forecasting. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a. URL https://openreview.net/forum?id=uqL2D9idAD.
- 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. *Advances in Neural Information Processing Systems*, 36, 2024b.
- Tianyi Yin, Jingwei Wang, Yunlong Ma, Han Wang, Chenze Wang, Yukai Zhao, Min Liu, Weiming Shen, and Yufeng Chen. Apollo-forecast: Overcoming aliasing and inference speed challenges in language models for time series forecasting, 2024. URL https://arxiv.org/abs/2412.12226.
- Wenzhen Yue, Yong Liu, Xianghua Ying, Bowei Xing, Ruohao Guo, and Ji Shi. Freeformer: Frequency enhanced transformer for multivariate time series forecasting. *arXiv* preprint arXiv:2501.13989, 2025.
- Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 11121–11128, 2023.
- Xinyu Zhang, Shanshan Feng, Jianghong Ma, Huiwei Lin, Xutao Li, Yunming Ye, Fan Li, and Yew Soon Ong. Frnet: Frequency-based rotation network for long-term time series forecasting. In *KDD*, pp. 3586–3597, 2024. URL https://doi.org/10.1145/3637528.3671713.
- Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *ICLR*, 2022.
- Tian Zhou, Ziqing MA, Xue Wang, Qingsong Wen, Liang Sun, Tao Yao, Wotao Yin, and Rong Jin. FiLM: Frequency improved legendre memory model for long-term time series forecasting. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 12677–12690. Curran Associates, Inc., 2022a. URL https://proceedings.neurips.cc/paper\_files/paper/2022/file/524ef58c2bd075775861234266e5e020-Paper-Conference.pdf.
- Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International Conference on Machine Learning*, pp. 27268–27286. PMLR, 2022b.
- Tian Zhou, Peisong Niu, Liang Sun, Rong Jin, et al. One fits all: Power general time series analysis by pretrained lm. *Advances in neural information processing systems*, 36:43322–43355, 2023.
- Yifan Zhou, Zeqi Xiao, Shuai Yang, and Xingang Pan. Alias-free latent diffusion models: Improving fractional shift equivariance of diffusion latent space. In *CVPR*, 2025.

# A PRELIMINARIES

#### A.1 PROBLEM STATEMENT

**Time Series.** Time series  $\mathbf{X} \in \mathbb{R}^{C \times N}$  refers to a sequence of data points ordered by time, where N denotes the total number of timestamps and C represents the number of channels at each timestamp. Time series forecasting involves predicting future data points based on historical time series observations. The historical observations can be represented as  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L] \in \mathbb{R}^{C \times L}$ , and L is the length of the historical look-back window. The future data for the next  $L_{\text{next}}$  time steps, denoted as  $\hat{X}_o = [\mathbf{x}_{L+1}, \mathbf{x}_{L+2}, \dots, \mathbf{x}_{L+L_{\text{next}}}] \in \mathbb{R}^{C \times L_{\text{next}}}$ , correspond to the forecast horizon. Given these, time series forecasting models are required to learn mapping functions  $\mathbf{F}: X \in \mathbb{R}^{C \times L} \to \hat{X}_o \in \mathbb{R}^{C \times L_{\text{next}}}$ .

Aliasing. This issue arises when different high-frequency components in a continuous signal are indistinguishably mapped to the same low-frequency components after sampling or improper downsampling. Formally, let the sampling interval be  $\Delta t$ , with the Nyquist frequency defined as  $f_{\mathrm{Nyquist}} = \frac{1}{2\Delta t}$ . Any frequency component  $f > f_{\mathrm{Nyquist}}$  in the signal will alias to a spurious frequency  $\tilde{f} = |f - k \cdot f_s|$  in the sampled sequence X, where  $f_s = \frac{1}{\Delta t}$  is the sampling rate, and  $k \in \mathbb{Z}^+$  ensures  $\tilde{f} \leq f_{\mathrm{Nyquist}}$ . This may occur when the sampling rate or downsampling operations fail to meet the Nyquist criterion, that is, the sampling frequency must be at least twice the highest frequency in the original signal. If not resolved, high-frequency components would fold back into lower frequencies during downsampling, creating spurious artifacts.

#### A.2 PROBLEM DESCRIPTION: ALIASING IN MULTI-SCALE TIME SERIES DOWNSAMPLING

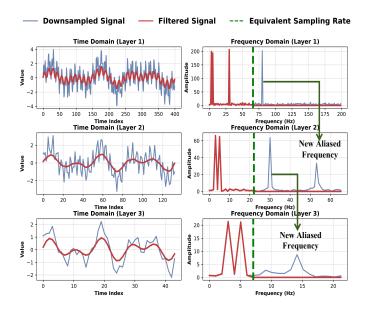


Figure 5: Left: filtering effect of the anti-aliasing filters; Right: emergence of new aliased frequencies.

In Figure.5, we present a case study exploring the critical role of anti-aliasing filters in signal preservation during multi-scale downsampling. By downsampling a synthetic signal containing both high-frequency and low-frequency components, we demonstrate the occurrence of aliasing during the reduction of the sampling rate.

The synthetic time series signal used in the study consists of several frequency components: low-frequency components (3 Hz and 5 Hz), high-frequency components (30 Hz and 80 Hz), and Gaus-

sian noise. The signal is initially sampled at a rate of 400 Hz. Subsequently, we perform multi-scale downsampling at different levels (each with a window size of 3), resulting in sampling rates of 400 Hz, 133 Hz, and 44 Hz.

According to the Nyquist sampling theorem, the Nyquist frequency is half of the sampling rate. Therefore, a 400 Hz sampling rate is sufficient to accurately sample the frequency components of the original signal. The calculation of aliasing frequencies is derived from the spectral periodicity characteristics of the Nyquist sampling theorem Nyquist (1928), based on the formula:

$$f_{\text{alias}} = |f_o - k \cdot f_s|, \tag{12}$$

where  $f_{\text{alias}}$  is the aliased frequency,  $f_o$  is the original high-frequency component, k is an integer representing the multiple mapping to the sampling frequency, and  $f_s$  is the sampling rate.

In the first layer, the Nyquist frequency is 200 Hz, corresponding to a sampling rate of 400 Hz. Given that the highest frequency component of the signal is 80 Hz, which is well below the Nyquist frequency, no aliasing occurs; all frequency components can be accurately sampled and reconstructed. In practical applications, an anti-aliasing filter limits frequency components above 66 Hz, thereby preventing aliasing and removing high-frequency noise.

In the second layer, the Nyquist frequency is reduced to 66.67 Hz (corresponding to a sampling rate of approximately 133.33 Hz), which results in aliasing of the original 80 Hz high-frequency component. According to the aliasing formula (12), for  $f_o = 80$  Hz and with k = 1:

$$f_{\text{alias}} = |80 - 1 \times 133.33| \approx 53.33 \,\text{Hz}.$$
 (13)

This calculation indicates that the 80 Hz component folds into the lower frequency region, specifically within the 50–60 Hz range, thereby introducing non-original frequency components and causing spectral distortion. The anti-aliasing filter in this layer effectively removes frequencies above 22 Hz to mitigate this issue.

In the third layer, the Nyquist frequency further decreases to 22.22 Hz (with a corresponding sampling rate of approximately 44.44 Hz), leading to the aliasing of the original 30 Hz component. Using the aliasing formula with k=1:

$$f_{\text{alias}} = |30 - 1 \times 44.44| \approx 14.44 \,\text{Hz},$$
 (14)

indicating that the 30 Hz component folds around 14 Hz. Additionally, due to the interaction between the sampling rate and the sampling process, frequency components in the 10–15 Hz range cannot be accurately represented, even though they lie below the Nyquist frequency. This aliasing phenomenon becomes particularly significant as the frequencies approach the Nyquist limit. Nevertheless, the anti-aliasing filter is still able to extract the true frequency information with reasonable accuracy, thereby alleviating the impact of aliasing.

#### A.3 AN INTUITIVE EXPLANATION OF ALIASING-RELATED CONCEPTS

This appendix provides a detailed explanation of the core signal processing concepts illustrated in Figure.6, which motivate the design of DMANet. We structure this explanation to clarify the relationship between sampling, spectral overlap, aliasing, and our proposed solution.

#### A.3.1 THE CORE PROBLEM: SPECTRAL OVERLAP AND ALIASING

The central challenge DMANet addresses is aliasing, a form of signal distortion that occurs during downsampling. To fully understand this phenomenon, it is crucial to distinguish between its **physical cause spectral overlap** and its **perceptual consequence aliasing**. These two terms describe different links in a cause-and-effect chain, where aliasing is the direct result of spectral overlap due to improper sampling Zhou et al. (2025).

• Sampling and Spectral Replicas: When a signal is sampled, its original spectrum is periodically replicated along the frequency axis, creating what are known as spectral replicas. Grabinski et al. (2022b) The act of sampling also limits the frequency range we can observe without distortion to a new, narrower baseband (from 0 Hz to the new Nyquist frequency).

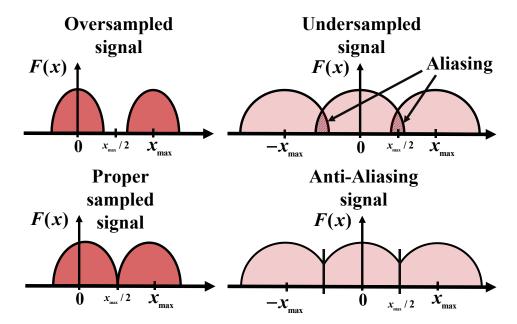


Figure 6: A conceptual illustration of the sampling process in the frequency domain. **Top-left:** Oversampling provides a wide guard band, preventing aliasing. **Top-right:** Undersampling causes spectral replicas to overlap, leading to aliasing where high frequencies (hatched areas) are misrepresented as low frequencies. **Bottom-left:** Proper (or critical) sampling meets the Nyquist criterion exactly, with replicas touching but not overlapping. **Bottom-right:** An anti-aliasing filter removes high-frequency content before sampling, ensuring that even with a lower sampling rate, no overlap occurs.

- The Cause: Spectral Overlap. If the sampling rate is too low to satisfy the Nyquist criterion, the separation between these spectral replicas becomes insufficient. This prevents the formation of a safety margin (Guard Band), causing them to physically overlap. This physical overlap, illustrated in the undersampled signal panel of Figure 6, is the root cause of the problem.
- The Consequence: Aliasing. This overlap is the direct cause of aliasing. Any high-frequency component from the original signal that exceeds the new Nyquist frequency is folded back into the new, observable low-frequency baseband. This process causes the original high-frequency information to appear as a spurious low frequency. This spurious frequency then mixes with the true low-frequency components within the baseband, becoming indistinguishable from them. Ultimately, this misrepresentation of high-frequency information as low-frequency information corrupts the signal's fidelity and is the core problem we address in our paper.

#### A.3.2 A DEEPER DIVE INTO THE SCENARIOS OF FIGURE.6

# Figure.6 visualizes four key scenarios:

- The Undesirable Case (Top-Right): The undersampled scenario is precisely the adverse outcome our work aims to prevent. The resulting aliasing (hatched areas) erroneously introduces spurious low-frequency patterns, a distortion that can severely hinder analysis and forecasting.
- The Ideal and Safe Case (Top-Left): The oversampled scenario is ideal because it successfully avoids aliasing. The empty space between the spectral replicas is a Guard Band, which does not imply information loss but rather a safe margin ensuring that the original spectrum can be unambiguously recovered.
- The Theoretical Boundary (Bottom-Left): Proper sampled (or Critical Sampling) occurs when the sampling rate is exactly twice the signal's maximum frequency  $(f_s = 2 \cdot f_{max})$ .

The spectral replicas touch edge-to-edge without overlapping. While theoretically sound, operating at this exact boundary is risky in practice.

• The DMANet Approach (Bottom-Right): The Anti-Aliasing signal panel illustrates the core principle of our work. Before an operation that would otherwise cause undersampling, an anti-aliasing filter is applied. This low-pass filter removes the high-frequency components (the part of the spectrum above  $x_{max}/2$ ) that would cause overlap. After this pre-filtering, even a lower sampling rate can be safely applied without generating aliasing artifacts.

# B PROOFS OF THE EQUIVALENT SAMPLING RATE (ESR)

#### B.1 NOTATIONS AND SIGNAL MODELING

#### B.1.1 ORIGINAL DISCRETE-TIME SIGNAL

First, we define  $X[n] \in \mathbb{R}^{C_i}$  as a discrete-time signal with  $C_i$  channels. The *i*-th channel of X[n], that is,  $X_i[n]$ , is obtained by sampling a continuous-time signal  $x_i(t)$ :

$$X_i[n] = x_i(nT_s), T_s = \frac{1}{f_s}, i = 1, \dots, C_i,$$
 (15)

where  $f_s$  is the original sampling rate and  $T_s$  is the sampling period. It is assumed that each continuous-time signal  $x_i(t)$  is band-limited to  $|\omega| \le \omega_B$ .

#### B.1.2 MODULE STRUCTURE

The input signal X[n] passes sequentially through a depth-wise convolution and a point-wise convolution defined in the architecture of DMANet:

$$X[n] \xrightarrow{\text{DepthwiseConv1d}(K,S)} U[m] \xrightarrow{\text{PointwiseConv1d}} V[m].$$
 (16)

First, the DepthwiseConv1d operation, characterized by a kernel  $h_i[k]$  of length K for each i-th input channel, where  $1 \leq k \leq K$ , and a stride S, transforms X[n] into an intermediate signal  $U[m] \in \mathbb{R}^{C_i}$ . Subsequently, this intermediate signal U[m] is processed by a PointwiseConv1d operation. This second stage uses a convolution matrix  $W \in \mathbb{R}^{C_o \times C_i}$ , with elements  $w_{j,i}$ , to map the  $C_i$  channels of U[m] to the  $C_o$  output channels, producing the final output signal  $V[m] \in \mathbb{R}^{C_o}$ . Consequently, the output sampling rate of the entire module is  $f_s' = f_s/S$ .

#### B.2 DOWNSAMPLING AND ALIASING CONDITIONS

As a baseline reference, consider directly downsampling the original signal X[n] by a factor S without any convolution filtering to obtain the signal Y[m]:

$$Y[m] = X[Sm]. (17)$$

The new sampling rate is  $f_s' = f_s/S$ . To avoid aliasing caused directly by downsampling, the bandwidth  $B = \omega_B/(2\pi)$  of the original continuous-time signal must satisfy the Nyquist-Shannon sampling theorem requirement, which is defined as:

$$B \le \frac{f_s'}{2} = \frac{f_s}{2S}.\tag{18}$$

Expressed in terms of normalized angular frequency  $\Omega_B = \omega_B T_s = 2\pi B T_s$ , we solve this equation and represent the condition to avoid aliasing as:

$$\Omega_B \le \frac{\pi}{S}.\tag{19}$$

This condition applies to ideal direct downsampling, assuming that perfect anti-aliasing filtering has been performed before downsampling to remove frequency components above  $\pi/S$ , corresponding to  $\frac{f_s}{2S}$ . Our proposed DMANet contains depth-wise and point-wise modules which perform filtering in its workflow, and its behaviors are more complex.

#### B.3 LINEAR MAPPING WITH DEPTH-WISE CONVOLUTION AND POINT-WISE CONVOLUTION

The operations of these two modules can be expressed as a series of linear mappings.

#### B.3.1 Depth-wise Convolution

The output of the depth-wise convolution for the *i*-th channel,  $U_i[m]$ , is computed as follows:

$$U_i[m] = \sum_{k=0}^{K-1} h_i[k] X_i[mS+k].$$
 (20)

Here, m is the index for the output sequence. Due to the stride S, the output  $U_i[m]$  depends on the input  $X_i[n]$  in the range from n = mS to n = mS + K - 1.

#### **B.3.2** Point-wise Convolution

The j-th output channel  $V_i[m]$  is obtained by a linear combination of  $U_i[m]$ :

$$V_j[m] = \sum_{i=1}^{C_i} w_{j,i} U_i[m].$$
(21)

Then, we substitute equation 20 into the above equation to get the expression of  $V_i[m]$ :

$$V_{j}[m] = \sum_{i=1}^{C_{i}} w_{j,i} \left( \sum_{k=1}^{K} h_{i}[k] X_{i}[mS+k] \right) = \sum_{i=1}^{C_{i}} \sum_{k=1}^{K} w_{j,i} h_{i}[k] X_{i}[mS+k].$$
 (22)

#### B.3.3 UNIFIED LINEAR MAPPING

To form a unified linear transformation at each output time m, we construct an input vector  $\mathbf{x}_m$  that contains all original input samples involved in computing V[m]:

$$\mathbf{x}_{m} = \begin{bmatrix} X_{1}[mS] \\ \vdots \\ X_{1}[mS+K-1] \\ \vdots \\ X_{C_{i}}[mS] \\ \vdots \\ X_{C_{i}}[mS+K-1] \end{bmatrix} \in \mathbb{R}^{C_{i}K}.$$

This is a column vector formed by stacking K consecutive samples, starting from mS, from each of the  $C_i$  channels. Concurrently, a weight matrix  $G \in \mathbb{R}^{C_o \times (C_i K)}$  is constructed. For the j-th row of G, its elements correspond to  $w_{j,i}h_i[k]$  and are arranged according to the order of the elements in  $\mathbf{x}_m$ . Specifically, if the p-th element of  $\mathbf{x}_m$  is  $X_i[mS+k]$ , then the weight in G corresponds to the output  $V_j[m]$  and this input element is  $G_{j,p} = w_{j,i}h_i[k]$ . Thus, the output V[m] can be defined as:

$$V[m] = G\mathbf{x}_m, \ V[m] \in \mathbb{R}^{C_o}. \tag{23}$$

This equation shows that at each output time m, the output vector V[m] is a linear mapping of a local window  $\mathbf{x}_m$  of the input signal.

# B.4 RANK CONSTRAINT AND DEGREES OF FREEDOM COUNTING

A fundamental property of linear algebra states that the rank of the matrix G, denoted as rank(G), is limited by its dimensions:

$$\operatorname{rank}(G) \le \min\{\text{number of rows, number of columns}\} = \min\{C_o, C_i K\}.$$
 (24)

We assume that the values of the weights  $h_i[k]$  and  $w_{j,i}$  are generic. They can be learned and are not overly sparse or linearly dependent, so that matrix G can achieve its theoretically maximum possible rank. Then, the dimension of independent information, also named the degrees of freedom, that the module extracts from the  $C_iK$ -dimensional input window  $\mathbf{x}_m$  and transmits to the  $C_o$ -dimensional output V[m] at each output time m is:

$$D = \operatorname{rank}(G) = \min\{C_i K, C_o\}. \tag{25}$$

This D represents the maximum dimension of linearly independent information from the input segment  $\mathbf{x}_m$  that the system can distinguish or represent, without considering noise or specific signal statistics.

To relate this total degree of freedom D to each channel of the input signal, we can average it over the  $C_i$  input channels. Thus, the equivalent temporal degrees of freedom  $\alpha$  contributed by each input channel to produce one output sample V[m] is:

$$\alpha = \frac{D}{C_i} = \frac{\min\{C_i K, C_o\}}{C_i} = \min\left\{\frac{C_i K}{C_i}, \frac{C_o}{C_i}\right\} = \min\left\{K, \frac{C_o}{C_i}\right\}. \tag{26}$$

Here,  $\alpha$  can be understood as: for each input channel, its information, under the combined effect of temporal processing through kernel length K and inter-channel mapping via  $C_o/C_i$ , is refined or compressed to be equivalent to  $\alpha$  independent information units. These units contribute to the final output sample V[m]. The bottleneck here is determined by the smaller of K (temporal context length per channel) and  $C_o/C_i$  (channel transformation ratio).

#### B.5 DEFINITION OF EQUIVALENT SAMPLING RATE

The actual output sampling rate of the module for each output channel is  $f_s' = f_s/S$ . At each output sampling instant, we have determined that each input channel contributes  $\alpha = \min\{K, C_o/C_i\}$  equivalent temporal degrees of freedom.

The Equivalent Sampling Rate  $f_{\rm ESR}$  is defined as a rate such that if each of the original  $C_i$  input channels were sampled at  $f_{\rm ESR}$ , and each sample carried one independent degree of freedom, then its total degrees of freedom throughput would match that of the current depth-wise and point-wise modules.

The total rate of generating degrees of freedom is:  $D = \min\{C_i K, C_o\} \times \frac{f_s}{S}$ . If  $C_i$  channels each operate at an equivalent sampling rate of  $f_{\rm ESR}$ , their total degrees of freedom rate is  $C_i \times f_{\rm ESR}$ :

$$C_i \times f_{\text{ESR}} = \min\{C_i K, C_o\} \times \frac{f_s}{S}.$$
 (27)

Then, we can get the solve for  $f_{\rm ESR}$ :

$$f_{\text{ESR}} = \frac{\min\{C_i K, C_o\}}{C_i} \times \frac{f_s}{S} = \min\left\{K, \frac{C_o}{C_i}\right\} \times \frac{f_s}{S}.$$
 (28)

If we normalize the original sampling rate  $f_s$  to 1, we obtain the normalized ESR:

$$ESR_{norm} = \frac{1}{S} \min \left\{ K, \frac{C_o}{C_i} \right\}. \tag{29}$$

Based on this equivalent sampling rate  $f_{\rm ESR}$ , we can define an equivalent Nyquist frequency  $f_{\rm Nyq.ESR}$ . This frequency represents the maximum bandwidth that the input signal can accommodate without information loss due to module structural limitations:

$$f_{\text{Nyq-ESR}} = \frac{f_{\text{ESR}}}{2} = \frac{f_s}{2S} \min \left\{ K, \frac{C_o}{C_i} \right\}.$$
 (30)

We can use  $f_{\rm ESR}$  to quantify the information processing capability or information retention degree of the downsampling module consisting of depth-wise convolution and point-wise convolution relative to each input channel. It provides a useful metric to compare the effective information throughput of modules with different parameter configurations with K, S,  $C_i$ , and  $C_o$ . It is important to note that the anti-aliasing significance of  $f_{\rm Nyq,ESR}$  also depends on whether the depth-wise convolution kernel  $h_i[k]$  can effectively act as a low-pass filter to attenuate frequency components above  $f_{\rm Nyq,ESR}$ . If  $h_i[k]$  is not an ideal low-pass filter, the frequency components of the original signal above  $f_{\rm Nyq,ESR}$ , even if not completely filtered out by  $h_i[k]$ , may not be accurately represented by the output V[m] due to subsequent dimensionality reduction.

## C IMPLEMENTATION DETAILS

We summarized details of datasets, evaluation metrics, experiments in this section.

#### C.1 Datasets details

We evaluated the performance of different models on several well-established datasets for long-term forecasting, including Weather, Electricity, Solar-Energy, PeMS(PEMS03, PEMS04, PEMS07, PEMS08), and the ETT series (ETTh1, ETTh2, ETTm1, ETTm2). Furthermore, to demonstrate DMANet's capability in handling highly non-stationary data, we conducted an extensive series of supplementary experiments on short-term forecasting across datasets from various domains. These include Health & Medical (ILI, COVID-19), Web Events (Wiki, Website), Finance (NASDAQ, SP500, DowJones), Market (CarSales), Energy (Power), and Society (Unemp). We detail the descriptions of the dataset in Table.6.

#### C.2 BASELINE DETAILS

Acknowledging that the performance of different methods varies across scenarios, we conducted a comprehensive comparison of various approaches under three distinct settings: long-term forecasting with a lookback window of 96, long-term forecasting with a lookback window of 720, and short-term forecasting. The evaluated methods are categorized as follows:

- Frequency-domain methods: TimeStacker Liu et al. (2025), FilterNet Yi et al. (2024a), FITS Xu et al. (2024), Fredformer Piao et al. (2024a), FEDformer Zhou et al. (2022b).
- CNN-based methods: ModernTCN Donghao & Xue (2024), TVNet Li et al. (2025), TSLANet Eldele et al. (2024), TimesNet Wu et al. (2023), PDF Dai et al. (2024), MICN Wang et al. (2023).
- MLP-based methods: SOFTS Han et al. (2024), TimeMixer Wang et al. (2024a), DLinear Zeng et al. (2023), TiDE Das et al. (2023), RLinear Li et al. (2023b), MTS-Mixer Li et al. (2023c)
- Transformer-based methods: TimeXer Wang et al. (2024b),iTransformer Liu et al. (2024), Crossformer Zhang & Yan (2022), Pathformer Chen et al. (2024b), Stationary Liu et al. (2022b), Pyraformer Liu et al. (2022a), Autoformer Wu et al. (2021)
- LLM-based methods: GPT4TS Zhou et al. (2023), Time-LLM Jin et al. (2024)
- KAN-based methods: TimeKAN Huang et al. (2025)
- Mamba-based methods: TimePro Ma et al. (2025)
- **Retrieval-Augmented methods:** RAFT Han et al. (2025)

#### C.3 IMPLEMENTATION DETAILS.

Regarding evaluation metrics, we used mean square error (MSE) and mean absolute error (MAE) for both long-term and short-term forecasting. All experiments were conducted using PyTorch on a single NVIDIA GeForce RTX 3090 24GB GPU. We applied an early stopping strategy to all baselines when the validation loss did not decrease for three consecutive epochs. Notably, inspired by FreDF Wang et al. (2025), we argue that formulating the loss function in the frequency domain is advantageous for learning an anti-aliasing architecture. Consequently, we directly adopted the frequency-domain MAE as the loss function for both long-term and short-term forecasting. More detailed settings can be found in Appendix.C.5.

#### C.4 FAIR COMPARISON SETTINGS.

To ensure a fair comparison and address challenges related to scaling laws, we maintained a consistent lookback window of 96 for all experiments in Table.9 and Table 10, and 720 for all experiments in Table.11 and Table.12. Our baseline comparisons mimic the experimental protocols established in TimesNet Wu et al. (2023), including same data processing and splitting procedures. For most

Table 6: Detailed dataset descriptions and statistics. **Dim** denotes the number of variates for each dataset. **Frequency** refers to the time interval between consecutive steps. **Split** indicates the data partitioning ratio (Train/Validation/Test). **Prediction len.** represents the prediction lengths. Our long-term forecasting employs a fixed input length of 96 or 720. For the majority of datasets, we evaluate across prediction horizons of 96, 192, 336, 720. A distinct setting is applied to the PeMS datasets, which are evaluated on shorter horizons of 12, 24, 48. For short-term forecasting, we adopt two settings: one with an input of 12 steps to predict 3, 6, 9, 12 steps, and another with an input of 36 steps to predict 24, 36, 48, 60 steps.

Dataset	Dim	Frequency	Total len.	Split	Prediction len.	Information
ETTh1, ETTh2	7	Hourly	17420	6:2:2	{96,192,336,720}	Electricity
ETTm1, ETTm2	7	15 mins	69680	6:2:2	{96,192,336,720}	Electricity
Weather	21	10 mins	52696	7:1:2	{96,192,336,720}	Weather
ECL	321	Hourly	26304	7:1:2	{96,192,336,720}	Electricity
Solar-Energy	137	10 mins	52560	7:1:2	{96,192,336,720}	Energy
PEMS03	358	5 mins	26209	6:2:2	{12,24,48}	Transportation
PEMS04	307	5 mins	16992	6:2:2	{12,24,48}	Transportation
PEMS07	883	5 mins	28224	6:2:2	{12,24,48}	Transportation
PEMS08	170	5 mins	17856	6:2:2	{12,24,48}	Transportation
ILI	7	Weekly	966	7:1:2	{24,36,48,60}	Health
COVID-19	55	Daily	335	7:1:2	{3,6,9,12}	Health
NASDAQ	12	Daily	3914	7:1:2	{24,36,48,60}	Finance
SP500	5	Daily	8077	7:1:2	{24,36,48,60}	Finance
DowJones	27	Daily	6577	7:1:2	{24,36,48,60}	Finance
CarSales	10	Daily	6728	7:1:2	{24,36,48,60}	Market
Power	2	Daily	1186	7:1:2	{24,36,48,60}	Energy
Website	4	Daily	2167	7:1:2	{3,6,9,12}	Web
Wiki	99	Daily	730	7:1:2	{3,6,9,12}	Web
Unemp	53	Monthly	531	6:2:2	{3,6,9,12}	Society

methods, we adopted the results reported in their original papers. For some methods that did not report results on the Solar-Energy dataset, we reproduced their performance using their official code repositories. The results for FITS Xu et al. (2024) and FreTSWang et al. (2025) were replicated from the FilterNet report Yi et al. (2024a); for other methods, we used the long-term prediction results provided in the iTransformer repository Liu et al. (2024). These results are based on the experimental configurations provided in the original paper or official code for each model. We verified that all hyperparameters for these baselines were selected from their respective official repositories, ensuring consistency with our fair comparison setup, where the only variations were the input and output sequence lengths.

For the experiments with the lookback window extended to 720, we referred to established baseline results: results in Table.11 were replicated from DUET Qiu et al. (2024), the results for GPT4TS Zhou et al. (2023) and TimeLLM Jin et al. (2024) in Table.12 were replicated from TSLANet Eldele et al. (2024), and the remaining results in Table.12 were replicated from TVNet Li et al. (2025). For short-term forecasting, we followed the results from the FreEformer repository Yue et al. (2025).

#### C.5 HYPERPARAMETER SETTINGS.

**Primary Long-term Forecasting Task** For our model hyperparameter selection, in 96 lookback window long-term forecasting, we fixed  $d_{\rm model}=512$ , downsampling layer l to 2, depth-wise convolution kernal size K to 3, stride s to 2, and set the proportion of channel changes c to 0.5. And we only performed a limited search on the encoder layers E, learning rate LR, and batch size. Detailed configurations for each dataset can be found in Table.7.

Other Long-term Forecasting Tasks For long-term forecasting with an extended 720 lookback window, as well as for the 96 lookback forecasting on PEMS datasets and 336 lookback univariate forecasting tasks, we implemented a more extensive hyperparameter search. This search was conducted for each forecast horizon within a given dataset to find the optimal configuration. The search space was defined as follows:  $d_{model} \in \{256, 512\}$ , Learning Rate  $LR \in \{1 \times 10^{-3}, 2 \times 10^{-3}, 5 \times 10^{-3}, 1 \times 10^{-2}, 2 \times 10^{-2}\}$ , Encoder Layers  $E \in \{1, 2, 3\}$ , Downsampling Layers  $l \in \{2, 3, 4\}$ , Batch Size  $e \in \{8, 16, 32, 64\}$ . Other hyperparameters, such as the convolutional kernel size and stride, remained fixed across all experiments, consistent with the settings used in the primary 96 lookback forecasting task. In contrast to all baseline lookback windows searched from  $\{192,336,512,672,720\}$  etc., We provide long-term forecasting for the fixed 720 lookback window.

**Short-term Forecasting** We implemented a more extensive hyperparameter search like Other Long-term Forecasting Tasks. This search was conducted for each forecast horizon within a given dataset to find the optimal configuration. The search space was defined as follows: Downsampling Layers is fixed 2,  $d_{model} \in \{256, 512\}$ , Learning Rate  $LR \in \{1 \times 10^{-3}, 2 \times 10^{-3}, 5 \times 10^{-3}, 1 \times 10^{-2}, 2 \times 10^{-2}\}$ , Encoder Layers  $E \in \{1, 2\}$ , Batch Size  $\in \{2, 4, 8, 16\}$ . Other hyperparameters, such as the convolutional kernel size and stride, remained fixed across all experiments, consistent with the settings used in the primary 96-lookback forecasting task.

**Ablation Study on Pre-Sampling Filtering** To validate our ESR-based filtering approach, we conducted an ablation study comparing it against alternatives that do not adhere to the Nyquist sampling theorem. Each experimental group differs from our full DMANet only in the cutoff frequency determination method within the Pre-Sampling Filtering module; all other structures and parameters remain identical. We categorize the compared methods into two groups: heuristic and classical filters.

HEURISTIC FILTERS These methods serve as simple, non-theoretical baselines. They are designed to mimic intuitive or simplistic approaches to filtering that one might adopt without a rigorous signal processing foundation.

- Max: For each time series in the batch, this filter identifies the frequency bin with the maximum amplitude and sets the cutoff frequency to twice its index. All components below this dynamic cutoff are preserved, while those above are zeroed out.
- **Random:** This filter applies a stochastic mask to the frequency spectrum, where each frequency component is independently dropped with a probability of p = 0.5.

CLASSICAL FILTERS These methods serve as benchmarks against well-established, theoretically-grounded filtering techniques. To ensure a fair comparison, a normalized cutoff frequency of 0.4 was used across all classic filter variants, preserving the lowest 80% of the frequency band.

- **Ideal:** A sharp cutoff filter where all frequency components above the cutoff frequency are set to zero.
- **Butterworth:** Known for its maximally flat passband, providing high-fidelity signal preservation. We used a 4th-order filter.
- Gaussian: A smooth filter often used to avoid ringing artifacts, with a sigma of 0.15.
- Chebyshev (Type I): Achieves a steeper rolloff than Butterworth at the cost of introducing ripples in the passband. We used a 4th-order filter with 0.5 dB of passband ripple.

Table 7: Experiment configuration of DMANet in 96 lookback window. All the experiments use the ADAM optimizer with the default hyperparameter configuration for  $(\beta_1, \beta_2)$  as (0.9, 0.999).

Dataset / Configurations	Mo	del l	Hyper-parameter		Trainii	ng Process	
	$ \overline{E} $	l	$d_{ m model}$	LR*	Loss	Batch Size	Epochs
ETTh1	1	2	512	$2*10^{-2}$	MAE	8	15
ETTh2	1	2	512	$1*10^{-2}$	MAE	8	15
ETTm1	1	2	512	$2*10^{-3}$	MAE	16	15
ETTm2	2	2	512	$5*10^{-3}$	MAE	32	15
Weather	1	2	512	$5*10^{-3}$	MAE	16	15
Electricity	2	2	512	$1*10^{-3}$	MAE	8	15
Solar-Energy	2	2	512	$5*10^{-3}$	MAE	16	15

<sup>\*</sup> LR means the initial learning rate.

# D FULL RESULTS

#### D.1 ERROR BARS

To evaluate the performance stability and robustness of DMANet, we conducted multiple independent runs with five different random seeds and compared its performance against the second-best model, TimeMixer. The results, averaged over four prediction horizons (96, 192, 336, and 720), are presented in Table 8. We report the mean and standard deviation of the MSE and MAE metricsacross the five experiments, as well as the confidence level of DMANet's superiority over TimeMixer. This performance improvement is statistically significant, with a 99% confidence level in all evaluated scenarios.

Table 8: Standard deviation and statistical tests for our DMANet method and second-best method (TimeMixer) on five datasets.

Metric		MSE			MAE	
Dataset	DMANet	TimeMixer	Confidence	DMANet	TimeMixer	Confidence
ETTm1	0.376±0.005	0.386±0.003	99%	0.388±0.003	0.399±0.001	99%
ETTm2	0.269±0.007	$0.278 \pm 0.001$	99%	$0.311 \pm 0.005$	$0.325 \pm 0.001$	99%
Weather	0.238±0.005	$0.245 \pm 0.001$	99%	$0.263 \pm 0.005$	$0.276 \pm 0.001$	99%
Electricity	0.171±0.002	$0.182 \pm 0.002$	99%	$0.264 \pm 0.002$	$0.272 \pm 0.002$	99%
Solar-Energy	0.228±0.003	$0.235 \pm 0.001$	99%	$0.249 \pm 0.002$	$0.292 \pm 0.001$	99%

#### D.2 Long-term Forecasting

Here, Table.9, Table.10, Table.11 and Table.12 present comprehensive evaluation results for long-term forecasting, including both configurations with fixed lookback windows L=96 and extended window settings L=720 designed to adhere to the scaling law inherent to TSF. In the L=96 fixed-window experiments, **consistent hyperparameters** were maintained across all forecast horizons within each dataset. By contrast, the L=720 experiments employed horizon-specific hyperparameter adjustments to enhance model adaptability while preserving scaling law compliance. Under both experimental paradigms, DMANet consistently demonstrates superior performance with statistically significant margins, thereby empirically validating its effectiveness and robustness. Notably, even when handling extended sequence lengths through augmented lookback windows, DMANet retains its inherent capability to adaptively model critical dependencies within extended temporal sequences.

The results for PESM dataset forecasting, presented in Table.13 for a lookback window of L=96 and the forecasting horizon  $T\in\{12,24,48\}$ , demonstrate the exceptional capability of DMANet. Across all four PEMS datasets, DMANet consistently outperforms all baselines. This superiority is quantified by average reductions of 14.4% in MSE and 5.7% in MAE compared to a strong baseline, iTransformer. We attribute this robust performance to our convolutional architecture's inherent proficiency in preserving localized features and mitigating the interference of high-frequency noise, which are critical for high-dimensional short-term prediction.

#### D.3 SHORT-TERM FORECASTING

The short-term forecasting results, presented in Table.15, validate the superiority of DMANet in handling highly non-stationary time series. Across a diverse set of challenging datasets including ILI (health), COVID-19 (pandemic), DowJones (finance), and Unemp (society), DMANet consistently achieves state-of-the-art performance, securing the top rank in 17 out of 20 metrics. It significantly outperforms other methods, including strong frequency-domain baselines like Fredformer and FilterNet. We attribute this exceptional capability in short-term and non-stationary forecasting to DMANet's synergistic design: its convolutional architecture excels at preserving local features, while the anti-aliasing structure effectively mitigates disruptive high-frequency noise. This robust performance on volatile, real-world data underscores the effectiveness of our approach in capturing the transient and complex patterns inherent to non-stationary signals.

# D.4 Univariate Forecasting

Here we provide the univariate forecasting results on ETT datasets. There is a target feature oil temperature within those datasets, which is the univariate time series that we are trying to forecast. As shown in Table.14, the anti-aliasing depth-wise convolution has better temporal modelling capabilities, allowing DMANet to achieve better performance than the state-of-the-art CNN-based ModerTCN in univariate forecasting tasks.

Table 9: Full results of long-term forecasting with a 96-step lookback window (Part I). The input sequence length L is set to 96 for all baselines. All results are averaged across four different forecasting horizon:  $T \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively. Among them, - means that the code has not yet been open sourced. We will put the summary table in the appendix of the next version.

Me	odels	DMA Ou			Stacker 125		eXer 24b		former 024		Mixer 24a		rNet 24a		ormer 24a		TS 24		eTS 24b
M	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	336	0.354 0.384	$\frac{0.372}{0.394}$	0.364 0.389	0.367 0.391	$0.362 \\ 0.395$	$0.383 \\ 0.407$	$0.377 \\ 0.426$	0.368 0.391 0.420	$\frac{0.361}{0.390}$	$0.381 \\ 0.404$	$0.364 \\ 0.396$	$0.383 \\ 0.406$	$0.363 \\ 0.395$	$0.380 \\ 0.403$	$0.392 \\ 0.424$	$0.393 \\ 0.414$	$0.382 \\ 0.421$	0.397 0.426
Щ									0.459										
	96								0.264										
ETTm2	192 336	0.231 0.289	0.288 0.325	0.235 0.293	0.292 0.329	0.237 0.296	$0.299 \\ 0.338$	$0.250 \\ 0.311$	$0.309 \\ 0.348$	$0.237 \\ 0.298$	$0.299 \\ 0.340$	$0.240 \\ 0.297$	$0.300 \\ 0.339$	$0.241 \\ 0.302$	$0.300 \\ 0.340$	$0.247 \\ 0.307$	$0.305 \\ 0.342$	$0.260 \\ 0.373$	0.329 0.405
Щ									0.407										
ETTh1	336	0.417 0.457	$\frac{0.420}{0.440}$	$\frac{0.429}{0.459}$	0.416 0.436	$0.429 \\ 0.468$	$0.435 \\ 0.448$	$0.441 \\ 0.487$	0.405 0.436 0.458 0.491	0.429 0.484	$0.421 \\ 0.458$	0.436 0.476	$0.422 \\ 0.443$	$0.440 \\ 0.472$	0.425 $0.440$	$0.436 \\ 0.478$	$0.423 \\ 0.444$	$0.453 \\ 0.503$	0.443 0.475
	Avg.	0.428	0.429	0.433	0.423	0.437	0.437	0.454	0.447	0.447	0.440	0.440	0.432	0.445	0.432	0.447	0.448	0.488	0.474
ETTh2	336 720	0.349 0.393 0.418	<b>0.374 0.410</b> 0.437	$\begin{array}{c} \underline{0.373} \\ 0.407 \\ \underline{0.412} \end{array}$	0.385 0.416 <b>0.431</b>	0.363 0.414 <b>0.408</b>	0.389 0.423 <u>0.432</u>	0.380 0.428 0.427	0.349 0.400 0.432 0.445	$0.372 \\ \underline{0.386} \\ \underline{0.412}$	0.392 0.414 0.434	0.369 0.420 0.430	0.395 0.432 0.446	0.370 <b>0.385</b> 0.419	$0.390 \\ \underline{0.413} \\ 0.439$	0.381 0.426 0.431	0.396 0.438 0.446	0.472 0.564 0.815	0.475 0.528 0.654
-	96								0.407										
Weather	192 336	<b>0.199</b> <u>0.256</u>	<b>0.238</b> <u>0.282</u>	$0.207 \\ 0.261$	0.241 0.281	0.204 0.261	$0.247 \\ 0.290$	$0.221 \\ 0.278$	0.254 0.296 0.349	0.208 <b>0.251</b>	$0.250 \\ 0.287$	$0.214 \\ 0.268$	$0.252 \\ 0.293$	$0.211 \\ 0.267$	$0.251 \\ 0.292$	$0.213 \\ 0.269$	$0.254 \\ 0.294$	$0.223 \\ 0.272$	0.275 0.316
	Avg.	0.236	0.262	0.243	0.264	0.241	0.271	0.258	0.279	0.240	0.271	0.248	0.278	0.246	0.272	0.249	0.276	0.255	0.363
Electricity	336 720	0.157 0.175 0.210	0.250 0.269 0.301	0.176 0.195 0.235	0.262 0.278 0.310	0.157 0.176 0.211	0.256 0.275 0.306	$\begin{array}{c} \underline{0.162} \\ \underline{0.178} \\ 0.225 \end{array}$	0.240 0.253 <b>0.269</b> 0.317	0.166 0.185 0.225	0.256 0.277 0.310	0.185 0.202 0.242	0.270 0.286 0.319	$0.165 \\ \underline{0.177} \\ 0.213$	$0.258 \\ \underline{0.273} \\ \underline{0.304}$	0.200 0.214 0.255	0.280 0.295 0.327	0.187 0.202 0.237	0.276 0.292 0.325
				0.194	0.275				0.270										
Solar-Energy	336 720	0.184 0.220 0.247 0.257	0.242 0.266 0.270	- - -	- - -	0.236 0.252 <u>0.244</u>	0.301 0.307 0.305	0.233 0.248 0.249	$\begin{array}{c} \underline{0.237} \\ 0.261 \\ \underline{0.273} \\ \underline{0.275} \end{array}$	0.222 0.231 0.223	0.283 0.292 0.285	0.259 0.284 0.284	0.284 0.298 0.298	0.226 0.254 0.249	0.259 0.277 0.284	0.397 0.433 0.429	0.387 0.410 0.396	0.283 0.299 0.298	0.338 0.344 0.351
S	Avg.	0.227	0.249	-	-	0.237	0.302	0.233	0.262	0.216	0.280	0.263	0.286	0.232	0.274	0.397	0.398	0.283	0.338

Table 10: Full results of long-term forecasting with a 96-step lookback window (Part II). The input sequence length L is set to 96 for all baselines. All results are averaged across four different forecasting horizon:  $T \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

M	odels	DMANet Ours		ePro	Time 20	KAN 25		FTS 24		DF 25		nTST 123		esNet	DLi 20	near 23	MI(	
M	etric	MSE MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	336	0.308 0.343 0.354 0.372 0.384 0.394 0.447 0.431	0.367 0.402	0.383 0.409	0.357 0.382	0.383 0.401	$0.375 \\ 0.405$	0.389 0.412	$0.373 \\ 0.402$	$0.385 \\ 0.404$	0.367 0.399	$\begin{array}{c} 0.385 \\ 0.410 \end{array}$	$0.374 \\ 0.410$	$0.387 \\ 0.411$	$0.382 \\ 0.415$	$0.391 \\ 0.415$	$0.403 \\ 0.436$	0.408 0.431
	Avg.	0.373 0.385	0.391	0.400	0.376	0.395	0.393	0.403	0.392	0.399	0.415	0.400	0.400	0.406	0.404	0.408	0.423	0.422
ETTm2	336	0.165 0.244 0.231 0.288 0.289 0.325 0.385 0.383	0.242 0.303	0.303 0.342	$\frac{0.239}{0.301}$	$0.299 \\ 0.340$	0.246 0.319	$0.306 \\ 0.352$	0.241 0.298	0.298 0.334	$0.241 \\ 0.305$	$0.302 \\ 0.343$	$0.249 \\ 0.321$	$0.309 \\ 0.351$	$0.284 \\ 0.382$	$0.361 \\ 0.429$	$0.284 \\ 0.381$	0.361 0.429
	Avg.	0.268 0.310	0.281	0.326	0.277	0.322	0.287	0.330	0.278	0.319	0.281	0.326	0.291	0.333	0.354	0.402	0.305	0.349
ETTh1	336	0.370 <b>0.391</b> 0.417 <b>0.420</b> 0.457 0.440 0.468 0.465	0.427 0.472	0.429 0.450	0.414 0.445	0.420 0.434	$0.435 \\ 0.480$	$0.431 \\ 0.452$	$0.430 \\ 0.474$	$0.427 \\ 0.451$	0.460 0.501	$0.445 \\ 0.466$	$0.436 \\ 0.491$	$0.429 \\ 0.469$	$0.446 \\ 0.489$	$0.441 \\ 0.467$	$0.454 \\ 0.493$	0.464 0.487
	Avg.	0.428 0.429	0.438	0.438	0.417	0.427	0.449	0.442	0.437	0.435	0.469	0.454	0.458	0.450	0.461	0.457	0.475	0.480
ETTh2	336 720	0.280 0.329 0.349 0.374 0.393 0.410 0.418 <u>0.437</u>	0.367 0.419 0.427	0.394 0.431 0.445	0.375 0.423 0.443	0.392 0.435 0.449	0.373 0.410 <b>0.411</b>	0.394 0.426 <b>0.433</b>	0.363 0.419 0.415	0.385 0.426 0.437	0.388 0.426 0.431	0.400 0.433 0.446	0.402 0.452 0.462	0.414 0.452 0.468	0.482 0.591 0.839	0.479 0.541 0.661	0.492 0.607 0.824	0.492 0.555 0.655
		0.361 0.388																
Weather	336 720	0.148 0.191 0.199 0.238 0.256 0.282 0.339 0.336	0.216 0.273 0.351	0.254 0.296 0.346	0.207 0.263 <b>0.338</b>	$\begin{array}{c} \underline{0.249} \\ \underline{0.290} \\ \underline{0.340} \end{array}$	0.217 0.282 0.356	0.253 0.300 0.351	0.220 0.275 0.356	0.253 0.294 0.347	0.225 0.278 0.354	0.259 0.297 0.348	0.219 0.280 0.365	0.261 0.306 0.359	0.237 0.282 0.345	0.295 0.331 0.382	0.239 0.285 0.351	0.299 0.336 0.388
_	Avg. 96	0.236 0.262																
Electricity	192 336 720	0.139 0.234 0.157 0.250 0.175 0.269 0.210 0.301	0.156 0.172 0.209	0.249 0.267 <b>0.299</b>	0.182 0.197 0.236	0.273 0.286 0.320	0.158 0.178 0.218	0.248 0.269 0.305	0.159 <b>0.172</b> <b>0.204</b>	<b>0.247 0.263</b> <u>0.294</u>	0.199 0.215 0.256	0.289 0.305 0.337	0.184 0.198 0.220	0.289 0.300 0.320	0.210 0.223 0.258	0.305 0.319 0.350	0.189 0.198 0.217	0.302 0.312 0.330
_		0.170 0.264																
Solar-Energy	336 720	0.184 0.217 0.220 0.242 0.247 0.266 0.257 0.270	0.231 0.250 0.253	0.263 0.281 0.285	0.285 0.315 0.313	0.326 0.338 0.340	0.229 0.243 0.245	0.253 0.269 0.272	0.276 0.301 0.308	0.288 0.306 0.316	0.267 0.290 0.289	0.310 0.315 0.317	0.296 0.319 0.338	0.318 0.330 0.337	0.320 0.353 0.357	0.398 0.415 0.413	0.278 0.298 0.299	0.354 0.375 0.379
	Avg.	0.227 0.249	0.232	0.266	0.292	0.331	0.229	<u>U.256</u>	0.279	0.292	0.270	0.307	0.301	0.319	0.330	0.401	0.283	0.338

Table 11: Full results of long-term forecasting with a 720-step lookback window (Part I) The input length L is fixed 720 for optimal horizon in the scaling law of TSF Shi et al. (2024). All results are averaged across four different forecasting horizon:  $T \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

Mode	els D	MANet Ours		DF 024		former 124		ormer 24b		TS 024	Time			hTST 123		former 122		esNet 023		near 123		onary 22b
Metri	c MS	E MAI	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
= 19 19 33 72	0.3 0.3	22 0.36 52 0.38	0.286 4 0.321 1 0.354 0 0.408	0.364 0.383	0.341 0.374	0.380 0.396	0.337 0.374	0.363 0.384	0.337 0.368	0.365 0.384	0.335 0.368	$0.372 \\ 0.386$	$0.329 \\ 0.362$	0.368 0.390	0.374 0.413	$0.410 \\ 0.432$	0.392 0.423	0.404 0.426	0.336 0.367	0.366 0.386	0.494 0.577	0.451 0.490
Av	g.   <b>0.3</b>	41 0.37	<b>4</b> <u>0.342</u>	0.376	0.361	0.390	0.357	<u>0.375</u>	0.357	0.377	0.356	0.380	0.349	0.381	0.464	0.456	0.408	0.415	0.356	0.378	0.531	0.472
2 19 33 72 72	0.2	14 0.28 54 0.32	0.163 7 0.219 0 0.269 3 0.349	0.290 0.330	0.242 0.282	$0.312 \\ 0.337$	0.219 0.267	$0.288 \\ 0.319$	0.219 0.272	$0.291 \\ 0.326$	0.225 0.277	$0.298 \\ 0.332$	0.221 0.276	$0.293 \\ 0.327$	0.369 0.588	$0.416 \\ 0.600$	0.254 0.313	$0.310 \\ 0.345$	0.224 0.277	$0.304 \\ 0.337$	0.338 0.432	0.373 0.416
Av	g.   <b>0.2</b>	45 0.30	0.250	0.313	0.269	0.327	0.253	0.308	0.254	0.313	0.257	0.318	0.256	0.314	0.501	0.505	0.292	0.331	0.259	0.324	0.383	0.390
Meather 19 33 72 72	0.1	39 0.23 39 0.27	8   0.147 7   0.193 5   0.245 7   0.323	0.240 0.280	0.200 0.252	$0.248 \\ 0.287$	0.191 0.243	$0.235 \\ 0.274$	0.215 0.261	$0.261 \\ 0.295$	0.192 0.247	$0.243 \\ 0.284$	$0.191 \\ 0.242$	$0.239 \\ 0.279$	0.198 0.258	$0.260 \\ 0.314$	0.219 0.278	$0.262 \\ 0.302$	0.216 0.258	$0.273 \\ 0.307$	0.241 0.341	0.290 0.341
Av	g.   0.2	18 0.25	0.227	0.263	0.232	0.270	0.225	0.257	0.244	0.280	0.226	0.264	0.224	0.261	0.234	0.292	0.255	0.282	0.242	0.293	0.293	0.315
90   19   33   72	0.1	45 0.24 50 0.25	7   0.128 2   0.147 8   0.165 0   0.199	0.242 0.260	0.154 0.169	0.250 0.265	0.157 0.170	$0.253 \\ 0.267$	0.154 0.170	0.250 0.268	0.168 0.189	0.269 0.291	$0.158 \\ 0.168$	0.260 0.267	0.146 0.165	$0.243 \\ 0.264$	0.180 0.204	0.280 0.304	0.154 0.169	$0.251 \\ 0.268$	0.180 0.204	0.283 0.305
Av	g.   0.1	54 0.25	<u>0.160</u>	0.253	0.163	0.258	0.168	0.261	0.169	0.265	0.184	0.284	0.171	0.270	0.171	0.263	0.190	0.290	0.167	0.264	0.194	0.295
is   90   19   33   72	0.1 0.1 0.1 0.1	33 0.23 95 0.24 93 0.24	5   0.181 0.200 3   0.208 4   0.212	0.259 0.269 0.275	0.193 0.203 0.223	0.257 0.266 0.281	0.196 0.195 0.208	0.220 0.220 0.237	0.229 0.241 0.248	0.267 0.273 0.277	0.201 0.190 0.203	0.259 0.256 0.261	0.204 0.212 0.215	0.302 0.293 0.307	0.208 0.212 0.215	0.226 0.239 0.256	0.206 0.208 0.232	0.276 0.284 0.294	0.220 0.234 0.243	0.282 0.295 0.301	0.395 0.410 0.377	0.386 0.394 0.376
Av			0.200		1						'				<u> </u>							<u> </u>

Table 12: Full results of long-term forecasting with a 720-step lookback window (Part II). The input length L is fixed 720 for optimal horizon in the scaling law of TSF Shi et al. (2024). All results are averaged across four different forecasting horizon:  $T \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

M	odels		ANet irs)		Net (25)	RLi (20	near 23a)		Mixer 23c)		CN (23)		rnTCN 024)	FEDf	ormer 22b)		AFT 025)		ANet		24TS		-LLM (24)
N	letric	MSE	MAE	MSE	MAE	MSE	MAE																
ETTm1	96 192 336 720	0.322 0.352	$0.364 \\ 0.381$	0.326 0.365	$0.367 \\ 0.391$	0.355 0.370	$0.363 \\ 0.383$	0.354 0.384	$0.386 \\ 0.405$	0.359 0.398	$0.387 \\ 0.413$	0.332 0.365	$0.368 \\ 0.391$	0.365 0.392	$0.415 \\ 0.425$	0.329 0.355	$0.367 \\ 0.383$	0.328 0.355	0.349 0.370 0.389 0.425	$0.332 \\ 0.366$	$0.372 \\ 0.394$	$0.310 \\ 0.352$	$0.358 \\ 0.384$
	Avg.	0.341	0.374	0.348	0.379	0.358	0.376	0.370	0.395	0.383	0.406	0.351	0.381	0.382	0.422	0.348	0.378	0.348	0.383	0.348	0.383	0.329	0.372
ETTm2	96 192 336 720	0.214 0.264	$0.287 \\ 0.320$	0.220 0.272	0.293 0.316	0.219 0.273	$0.290 \\ 0.326$	0.241 0.297	$0.303 \\ 0.338$	0.245 0.295	$0.316 \\ 0.350$	0.222 0.272	$0.293 \\ 0.324$	0.252 0.324	$0.318 \\ 0.364$	0.219 0.275	0.296 0.336	0.224 0.275	0.259 0.297 0.329 0.380	$0.229 \\ 0.286$	$0.301 \\ 0.341$	0.219 0.271	0.293 0.329
	Avg.	0.245	0.307	0.251	0.311	0.256	0.314	0.277	0.325	0.277	0.336	0.253	0.314	0.292	0.343	0.254	0.320	0.256	0.316	0.226	0.326	0.251	0.313
Weather	96 192 336 720	0.189 0.239	$0.237 \\ 0.275$	0.194 0.235	0.238 0.277	0.218 0.265	$0.260 \\ 0.294$	0.199 0.249	0.248 0.291	0.220 0.275	$0.283 \\ 0.328$	0.196 0.238	0.245 0.277	0.275 0.339	$0.329 \\ 0.377$	0.211 0.260	0.264 0.302	0.193 0.245	0.197 0.241 0.282 0.337	$0.204 \\ 0.254$	$0.248 \\ 0.286$	0.189 0.262	0.235 0.279
-	Avg.	0.218	0.257	0.221	0.261	0.247	0.279	0.235	0.272	0.242	0.298	0.224	0.264	0.310	0.357	0.241	0.286	0.325	0.337	0.237	0.270	0.225	0.257
Electricity	96 192 336 720	0.145 0.160 0.182	0.242 0.258 0.280	0.165 0.164 0.190	0.241 0.269 0.284	0.154 0.171 0.209	0.248 0.264 0.297	0.163 0.176 0.212	0.261 0.277 0.308	0.168 0.196 0.203	0.279 0.308 0.312	0.143 0.161 0.191	0.239 0.259 0.286	0.197 0.213 0.233	0.311 0.328 0.344	0.149 0.161 0.197	0.247 0.259 0.297	0.152 0.168 0.205	0.229 0.244 0.262 0.293	0.153 0.169 0.206	0.251 0.266 0.297	0.160 0.160 0.192	0.248 0.248 0.298
$\perp$	Avg.	0.154	0.252	0.165	0.254	0.169	0.261	0.173	0.272	0.182	0.292	<u>0.156</u>	0.253	0.207	0.321	0.160	0.259	0.165	0.257	0.167	0.263	0.158	0.252

Table 13: Full results of long-term forecasting with a 96-step lookback window (Part III). The input sequence length L is set to 96 for all baselines. All results are averaged across four different forecasting horizon:  $T \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

Mo	odels	DM	ANet	iTrans	former	Fredf	ormer	Time	Mixer	Patch	nTST	Crossi	former	Time	esNet	TI	DE	DLi	near	Fre	TS	FEDf	ormer
M	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
EMS03		0.086	0.193	0.093	$\begin{array}{c} \underline{0.174} \\ \underline{0.201} \\ \textbf{0.236} \end{array}$	0.094	0.205	0.113	0.226	0.142	0.259	0.121	0.240	0.118	0.223	0.257	0.371	0.201	0.317	0.127	0.198	0.241	0.275
PE	Avg.	0.094	0.200	0.096	<u>0.204</u>	0.105	0.214	0.127	0.235	0.151	0.265	0.138	0.253	0.119	0.271	0.271	0.380	0.219	0.295	0.137	0.234	0.167	0.291
PEMS04		0.082	0.185	0.095	$\begin{array}{r} \underline{0.183} \\ \underline{0.205} \\ \underline{0.233} \end{array}$	0.117	0.224	0.128	0.243	0.153	0.257	0.131	0.256	0.103	0.215	0.292	0.398	0.224	0.340	0.144	0.258	0.177	0.293
H	Avg.	0.086	0.190	0.098	0.207	0.125	0.215	0.144	0.254	0.162	0.273	0.145	0.267	0.109	0.220	0.307	0.405	0.236	0.350	0.148	0.265	0.195	0.308
MS07		0.074	0.174	0.088	0.165 0.190 0.215	0.089	0.192	0.111	0.219	0.150	0.262	0.139	0.247	0.101	0.204	0.271	0.383	0.210	0.329	0.127	0.239	0.125	0.244
PE	Avg.	0.080	0.179	0.088	0.190	0.096	0.197	0.140	0.244	0.166	0.270	0.181	0.272	0.106	0.208	0.297	0.394	0.241	0.343	0.142	0.247	0.133	0.282
EMS08	48	0.085 0.121	0.192 0.235	0.115 0.186	0.182 0.219 0.235	<u>0.112</u> <u>0.174</u>	0.214 0.267	0.137 0.265	0.246 0.343	0.224 0.321	0.281 0.354	0.215 0.315	0.260 0.335	0.141 0.198	0.238 0.283	0.318 0.497	0.409 0.510	0.248 0.440	0.353 0.470	0.152 0.247	0.256 0.331	0.210 0.320	0.310 0.394
Ь	Avg.	0.090	0.198	0.127	0.212	0.122	0.222	0.164	0.263	0.238	0.289	0.232	0.270	0.150	0.244	0.347	0.421	0.281	0.366	0.165	0.264	0.234	0.326

Table 14: Univariate long-term forecasting results on ETT datasets. Following PatchTST and ModerTCN, input length is fixed as 336 and prediction lengths are  $T \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

_																			
Μ	odels	DM	ANet	Mode	rnTCN	iTrans	former	Time	Mixer	Patcl	nTST	DLi	near	Pyraf	ormer	FEDf	ormer	Autof	ormer
	letric		MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	96	0.026	0.122	0.026	0.121	0.029	0.127	0.029	0.128	0.029	0.126	0.028	0.123	0.127	0.281	0.033	0.140	0.056	0.183
LI	192	0.039	0.150	0.040	0.152	0.045	0.162	0.044	0.160	0.043	0.158	0.045	0.156	0.205	0.343	0.058	0.186	0.081	0.216
ETTm1	336						0.189												
Ξ	720	0.072	0.203	0.073	0.206	0.080	0.218	0.081	0.218	0.080	0.217	0.080	0.210	0.387	0.485	0.102	0.250	0.110	0.267
	Avg.	0.047	0.162	0.048	0.163	0.053	0.174	0.053	0.173	0.052	0.171	0.054	0.168	0.255	0.392	0.069	0.202	0.081	0.221
	96	0.063	0.182	0.065	0.183	0.071	0.193	0.068	0.187	0.071	0.192	0.063	0.183	0.074	0.208	0.067	0.198	0.065	0.189
2	192	0.093	0.228	0.095	0.232	0.109	0.248	0.101	0.236	0.102	0.237	0.092	0.227	0.116	0.252	0.102	0.245	0.118	0.256
ETTm2	336	0.117	0.260	0.119	0.261	0.141	0.289	0.133	0.278	0.130	0.274	0.119	0.261	0.143	0.295	0.130	0.279	0.154	0.305
딤	720	0.167	0.317	0.173	0.323	0.190	0.343	0.183	0.332	0.179	0.328	0.175	0.320	0.197	0.338	0.178	0.325	0.182	0.335
	Avg.	0.110	0.247	0.113	0.250	0.128	0.268	0.121	0.258	0.121	0.258	0.112	0.248	0.133	0.273	0.119	0.262	0.130	0.271
	96	0.054	0.176	0.055	0.179	0.059	0.185	0.057	0.181	0.056	0.181	0.056	0.180	0.099	0.277	0.079	0.215	0.071	0.206
-	192						0.208												
ETTh1	336	0.073	0.215	0.074	0.214	0.084	0.223	0.085	0.227	0.094	0.242	0.098	0.244	0.198	0.370	0.119	0.270	0.107	0.258
E	720	0.082	0.227	0.086	0.232	0.089	0.236	0.083	0.227	0.101	0.250	0.189	0.359	0.209	0.348	0.142	0.299	0.126	0.283
	Avg.	0.069	0.205	0.071	0.206	0.076	0.213	0.074	0.210	0.082	0.221	0.104	0.247	0.170	0.335	0.111	0.257	0.105	0.252
	96	0.121	0.269	0.124	0.274	0.136	0.287	0.133	0.283	0.130	0.276	0.131	0.279	0.152	0.303	0.128	0.271	0.153	0.306
2	192	0.154	0.310	0.164	0.321	0.187	0.342	0.190	0.341	0.181	0.331	0.176	0.329	0.197	0.370	0.185	0.330	0.204	0.351
ETTh2	336	0.174	0.336	0.171	0.336	0.219	0.374	0.226	0.379	0.226	0.379	0.209	0.367	0.238	0.385	0.231	0.378	0.246	0.389
됴	720	0.211	0.371	0.228	0.384	0.253	0.403	0.241	0.396	0.253	0.406	0.276	0.426	0.274	0.435	0.278	0.420	0.268	0.409
	Āvg.	0.165	0.322	0.172	0.329	0.199	0.352	0.198	0.350	0.198	0.348	0.198	0.350	0.215	0.373	0.206	0.350	0.218	0.364

Table 15: Full results of short-term forecasting on supplementary datasets from domains including Health & Medical (ILI, COVID-19), Web Events (Wiki, Website), Finance (NASDAQ, SP500, DowJones), Market (CarSales), Energy (Power), and Society (Unemp). The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

M	lodel	DMA	Net	Timel	Mixer	Filter	rNet	FI	ΓS	DLi	near	Fredfo	ormer	Patch	TST
M	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ITI	24 36 48 60	1.746 1.718 1.744 1.842	0.813 0.817 0.826 0.839	2.110 2.084 <u>1.961</u> 1.926	0.879 0.890 <u>0.866</u> <u>0.878</u>	2.190 1.902 2.051 2.151	0.882 0.925	3.718 3.994 4.543	1.422 1.554	3.172	1.194 1.232	<u>1.925</u>	0.922 0.913	2.123 2.001	0.849 0.912 0.883 0.895
_	Avg	1.763	0.824	2.020	0.878	2.073	0.885		1.465	3.083	1.217	1.947	0.899	2.128	0.885
Covid19	3 6 9 12	1.098 1.735 2.167 2.640	0.489 0.625 0.722 0.843	1.237 2.003 2.594 3.103	0.547 0.739 0.860 <u>0.981</u>	1.195 1.839 2.537 <u>2.782</u>	0.897	2.683 3.147	0.919 1.050	3.803	0.909 1.053 1.160 1.288	1.165 1.465 2.145 2.833	0.548 0.685 0.845 0.984		0.573 0.762 0.916 1.030
	Avg	<u>1.910</u>	0.670	2.234	0.782	2.088	0.780	2.875	0.979	3.483		1.902	0.765	2.221	0.820
NASDAQ	24 36 48 60	0.200 0.233	0.296 0.323	0.122 0.183 <b>0.200</b> 0.238	$\begin{array}{c} 0.221 \\ 0.279 \\ \underline{0.298} \\ \underline{0.328} \end{array}$	0.224 0.259	0.273 0.314 0.340	0.140 0.184 0.234 0.282	0.284 0.324 0.357	0.196 0.244 0.318	0.344 0.401	$\begin{array}{c} 0.170 \\ \hline 0.218 \\ 0.262 \end{array}$	0.306 0.339	0.265	0.314 0.339
	Avg	0.177	0.273	0.186	0.281	0.197	0.289		0.302	0.228	0.331	0.194	0.285	0.198	0.286
Wiki	3 6 9 12	6.116 6.419 6.665 6.824	0.372 0.388 0.402 0.411		0.418 0.426	6.697 6.899	0.401 0.416 0.426	7.470 8.326 8.869 9.394	0.544 0.564 0.608	6.254 6.579 6.776 6.927	0.467 0.508 0.513	6.190 6.696 6.768 7.168	0.387 0.404 <u>0.411</u> 0.424	6.112 6.425 6.743 6.814	0.380 0.395 0.426 0.414
	Avg	6.506	0.393	6.572	0.409	6.572	0.411	8.515		6.634		6.705	0.406	6.523	0.404
SP500	24 36 48 60	0.153 0.205 0.250 0.293		$\begin{array}{c} \underline{0.159} \\ \underline{0.218} \\ \underline{0.264} \\ 0.322 \end{array}$	0.288 0.343 0.367 0.416	0.280	$\frac{0.341}{0.384}$	0.193 0.259 0.324 0.391	0.389 0.439	0.291	0.363	0.283	0.315 0.365 0.394 0.438	0.221	0.298 <u>0.341</u> 0.397 <u>0.409</u>
-	Avg	0.225	0.329	0.241	0.353			0.291			0.391	0.261	0.378	0.246	0.361
DowJones	24 36 48 60	10.422 13.975	0.800 0.917	8.327 11.192 15.278 20.997	0.813 0.945	12.011 14.814	$0.823 \\ 0.933$	11.907 15.821	$0.837 \\ 0.969$	10.986 14.157	0.922	14.696	0.808 0.921	11.210 14.866	0.935
_	Avg	11.957	0.850	13.948											
CarSales	24 36 48 60	0.318 0.332 0.346 0.357	0.314 0.327 0.340 0.352		0.331 0.343			0.373 0.385	0.360	0.368 0.382	0.365 0.379	0.333 0.349	0.335	0.332	0.330
	Avg	0.338	0.333	0.338	0.336	0.336	0.335	0.379	0.365	0.373	0.368	0.340	0.338	0.338	0.335
Power	24 36 48 60	1.408		1.341 1.420 1.567 1.609	0.963	1.590 1.680	0.968 1.009		0.994 1.052	1.518 1.610	0.957 0.995	1.652	0.953 1.008	1.710	1.020
	Avg	1.373		1.484		1.614									
Website	3 6 9 12 Avg	0.112 0.154 0.199	0.238 0.265 0.294	0.086 0.124 0.159 0.204 0.143	0.248 0.275 0.306	0.116 0.151 <u>0.194</u>	0.242 0.269 0.297	0.235 0.276 0.409	0.356 0.372 0.484	0.182 0.220 0.255	0.302 0.330 0.355	0.116 0.150	0.241 0.268 0.301	0.121 0.157	0.246 0.273 0.302
Unemp	3 6 9 12	0.036 0.081 0.127	0.117 0.180 0.234	0.015 0.057 0.109 0.195	0.154 0.213 0.293	0.043 0.107 0.155	0.130 0.216 0.255	0.229 0.369 0.475	0.345 0.443 0.500	0.115 0.191 0.240	0.255 0.329 0.386	0.046 0.095 <u>0.148</u>	0.139 0.198 <u>0.250</u>	0.043 0.093 0.164	
	Avg	0.064	0.146	0.094	0.183	0.079	0.166	0.308	0.394	0.154	0.292	0.075	0.163	0.078	0.160

#### D.5 RESULTS FOR HYPERPARAMETER ANALYSIS

In this section, we also explore the effect of the hyperparameters used in our experiments, including the depth-wise convolution kernal size K, the depth-wise convolution kernal stride size s, the channel change c and  $\lambda$  on the loss function.

The channel change c signifies the alteration in the number of channels during the downsampling process, where values below 1 denote a reduction in channel quantity, whereas values exceeding 1 indicate channel expansion.

For  $\lambda$ , according to FreDF Wang et al. (2025), the loss function is a weighted sum of the time-domain MSE and the frequency-domain MAE.  $\lambda$  represents the proportion of the frequency MAE in the loss function, and  $(1 - \lambda)$  represents the proportion of the time-domain MSE.

We show the experimental results from Table.16 to Table.19.

Table 16: Impact of kernal size. A lower MSE or MAE indicates a better performance.

Models	Metrics	Wea	ther	ET	Th2	ET	Гт2	
		96	336	96	336	96	336	
K = 1	MSE MAE	0.151 0.194	0.274 0.325	0.289 0.334	0.393 0.411	0.169 0.248	0.289 0.326	
K = 3	MSE MAE	0.148 0.191	0.256 0.282	0.280 0.329	0.393 0.410	0.165 0.244	0.289 0.325	
K=5	MSE MAE	0.149 0.192	0.259 0.285	0.286	0.393 0.410	0.168 0.247	0.296 0.331	
K = 7	MSE MAE	0.150 0.193	0.258 0.283	0.286	0.394 0.411	0.167 0.245	0.292 0.327	

Table 17: Impact of stride size. A lower MSE or MAE indicates a better performance.

Models	Metrics	Wea	ther	ET	Th2	ET	Гт2
		96	336	96	336	96	336
s=1	MSE MAE	0.151 0.194	0.259 0.284	0.281 0.331	0.423 0.421	0.170 0.247	0.297 0.331
s=2	MSE MAE	0.148 0.191	0.256 0.282	0.280 0.329	0.393 0.410	0.165 0.244	0.289 0.325
s=3	MSE MAE	0.149 0.192	0.259 0.284	0.274 0.325	0.393 0.410	0.167 0.246	0.290 0.326
s=4	MSE MAE	0.148 0.190	0.257 0.282	0.279 0.326	0.395 0.411	0.167 0.246	0.291 0.327

Table 18: Impact of channel change. A lower MSE or MAE indicates a better performance.

Models	Metrics	Wea	ather	ET	Th2	ET	Γm2
		96	336	96	336	96	336
c = 0.25	MSE MAE	0.149 0.193	0.259 0.284	0.278 0.326	0.396 0.412	0.167 0.245	0.291 0.327
c = 0.5	MSE MAE	0.148 0.191	0.256 0.282	0.280 0.329	0.393 0.410	0.165 0.244	0.289 0.325
c = 1	MSE MAE	0.149 0.191	0.261 0.287	0.284 0.333	0.408 0.415	0.171 0.250	0.293 0.328
c=2	MSE MAE	0.149 0.192	0.257 0.283	0.281 0.332	0.401 0.415	0.169 0.247	0.294 0.327
c = 4	MSE MAE	0.152 0.196	0.261 0.287	0.292 0.337	0.398 0.415	0.173 0.250	0.293 0.328

Table 19: Impact of  $\lambda$  in loss. A lower MSE or MAE indicates a better performance.

Models	Metrics	Wea	ther	ET	Th2	ET	Гт2
		96	336	96	336	96	336
$\lambda = 0.1$	MSE MAE	0.149 0.191	0.257 0.283	0.289 0.333	0.392 0.412	0.168 0.246	0.297 0.332
$\lambda = 0.3$	MSE MAE	0.149 0.192	0.259 0.284	0.289 0.332	0.394 0.411	0.167 0.246	0.297 0.332
$\lambda = 0.5$	MSE MAE	0.149 0.191	0.259 0.284	0.286 0.331	0.394 0.412	0.169 0.246	0.298 0.332
$\lambda = 0.7$	MSE MAE	0.149 0.191	0.259 0.283	0.290 0.332	0.396 0.414	0.169 0.247	0.297 0.332
$\lambda = 1$	MSE MAE	0.148 0.191	0.256 0.282	0.280 0.329	0.393 0.410	0.165 0.244	0.289 0.325

1838

1840

1841

1842 1843

1844 1845

1846

1847

1848

1849

1850 1851

1852

1853

1855

1857

1870 1871

1872

1873

1874

1875 1876

1877 1878

1879

1880

1881

1882

1883

1884

1885

1888

1889

C = 96

C = 192

C = 336

# E MORE DETAILS OF COMPUTATIONAL COSTS

To comprehensively evaluate the efficiency and scalability of DMANet, we conducted controlled experiments on both synthetic and real-world datasets. Our analysis focuses on two key aspects: the computational overhead of our proposed components and the overall model's performance compared to state-of-the-art methods.

#### E.1 EFFICIENCY AND SCALABILITY ANALYSIS ON SYNTHETIC DATA

We first use synthetic data to perform a fine-grained analysis under controlled conditions, isolating the impact of sequence length and channel dimensions. With fixed hyperparameters (look-back window=96, batch size=64, etc.), we measure inference speed (ms) and peak GPU memory (MB) under two scenarios: (1) fixing the number of channels C while varying the sequence length T, and (2) fixing T while varying C. Each experiment was repeated 500 times for stability. The results are presented in Table.20.

Table 20: Inference Speed (ms) and Memory Usage (MB) Comparison Across Different Models and Configurations. Values for speed are reported as mean  $\pm$  std over 500 runs.

	DMAN	Net	w/o-E	SR	Chebysl	nev
Configuration	Speed	Memory	Speed	Memory	Speed	Memory
T = 256	$1.376\pm0.122$	15.71	1.309±0.359	15.71	1.916±0.150	15.71
T = 512	$1.372 \pm 0.107$	31.29	$1.325 \pm 0.115$	31.29	$1.942 \pm 0.173$	31.29
T = 1024	$1.677 \pm 0.137$	65.38	$1.600 \pm 0.346$	65.38	$2.249 \pm 0.182$	65.38
T = 2048	$2.915 \pm 0.151$	146.49	$2.846 {\pm} 0.256$	146.49	$3.576 \pm 0.019$	146.49
C = 48	$1.518\pm0.103$	61.36	1.444±0.222	61.36	1.523±0.206	61.36
C = 96	$1.982 \pm 0.141$	122.40	$1.882 \pm 0.154$	122.40	$2.640 \pm 0.187$	122.40
C = 192	$3.731\pm0.047$	251.77	$3.580 \pm 0.072$	251.77	$4.250 \pm 0.182$	251.77
C = 336	$7.029 \pm 0.025$	471.27	$6.760 \pm 0.043$	471.27	$7.399 \pm 0.167$	471.27
	Linea	ır	TransC	onv	Attenti	on
Configuration	Speed	Memory	Speed	Memory	Speed	Memory
T = 256	$1.228\pm0.352$	15.89	$1.405\pm0.128$	15.92	$3.065 \pm 0.234$	75.04
T = 512	$1.201 \pm 0.158$	32.02	$1.605 \pm 0.124$	31.68	$3.412 \pm 0.239$	280.09
T = 1024	$1.838 \pm 0.196$	68.39	$2.191 \pm 0.168$	66.70	$8.710 \pm 0.117$	1084.94
T = 2048	$5.079 \pm 0.117$	158.52	$3.818 \pm 0.089$	147.73	$29.417 \pm 0.224$	4270.15
C = 48	1.414±0.225	61.43	1.915±0.161	61.33	$3.892 \pm 0.186$	298.42

From these results, we draw two key conclusions:

 $1.794\pm0.078$ 

 $3.374\pm0.054$ 

 $5.642\pm0.150$ 

122.31

252.63

471.27

1. The computational overhead of our dynamic anti-aliasing (ESR Filter) is negligible. A direct comparison between DMANet and its ablated version (w/o ESR) reveals that the peak memory usage is nearly identical across all configurations. The time overhead introduced by the ESR filter is minimal, with a worst-case relative increase of only 2.4% (at T=2048). Furthermore, DMANet is consistently faster than the variant using a classical Chebyshev filter. This empirically proves that our dynamic anti-aliasing mechanism is a computationally lightweight strategy that does not introduce a performance bottleneck.

 $2.432\pm0.303$ 

 $4.418\pm0.281$ 

 $8.116\pm0.019$ 

119.46

236.45

414.28

 $5.278\pm1.450$ 

 $8.956\pm0.035$ 

 $15.971\pm0.102$ 

332.12

404.04

518.38

**2.** The efficiency and scalability of frequency-domain interpolation for upsampling. We further validate our choice of upsampling mechanism by comparing it with common alternatives. The Attention-based method is not viable for long sequences due to the explosive, quadratic growth in its memory and time costs. While a simple Linear layer is fast, it scales poorly when processing very long sequences (e.g., at T=2048, its speed degrades significantly). Although Transposed Con-

volution is lightweight, our method is faster in most scenarios. In conclusion, our chosen frequency-domain interpolation achieves an excellent balance of cost-effectiveness and scalability across different data shapes.

# E.2 EFFICIENCY COMPARISON WITH STATE-OF-THE-ART MODELS ON REAL-WORLD DATASETS

Next, we benchmark the overall efficiency of DMANet against leading SOTA models on real-world datasets. We fix the input and prediction lengths (T=96, F=96) to ensure a fair comparison and report on memory usage, multiply-accumulate operations (MACs), and inference speed, alongside predictive accuracy.

Table 21: A comparison of Speed, memory consumption (Memory) and multiply-accumulate operations (MACs) for DMANet and five other models. To ensure a fair comparison, we fix the prediction length F=96 and the input length T=96.

Dataset		]	ETTm2					Weather		
Metric	Memory	MACs	Speed (s / iter)	MSE	MAE	Memory	MACs	Speed (s / iter)	MSE	MAE
iTransformer	4.45MB	14.16M	0.0117	0.182	0.265	62.63MB	682.50M	0.0178	0.175	0.215
TimeMixer	71.91MB	12.15M	0.0350	0.182	0.267	169.20MB	41.40M	0.0623	0.163	0.209
TimesNet	215.78MB	72.60G	0.1483	0.181	0.261	103.49MB	18.07G	0.0730	0.169	0.220
PatchTST	146.60MB	4.33G	0.0305	0.180	0.264	153.64MB	1.24G	0.0683	0.189	0.230
DLinear	0.42MB	0.30M	0.0033	0.193	0.292	0.97MB	0.30M	0.0035	0.196	0.256
DMANet (Ours)	0.98MB	0.53M	0.0262	0.175	0.253	2.63MB	0.89M	0.0226	0.156	0.201

As shown in Table.21, DMANet demonstrates a state-of-the-art balance between efficiency and performance. Compared to Transformer-based models (iTransformer, PatchTST) and CNN-based models (TimesNet, TimeMixer), DMANet requires significantly less memory and fewer MACs while achieving superior forecasting accuracy. While DLinear is exceptionally fast due to its simple architecture, DMANet provides a substantial accuracy improvement with only a marginal increase in memory and MACs. These results confirm that the lightweight and scalable design choices validated in our synthetic experiments translate directly to a highly competitive and efficient model for real-world applications.

#### F More details of Pre-Sampling Filtering

To comprehensively evaluate the robustness of our model and its generalization ability to different types of signal disturbance, we synthesized noise and superimposed it onto the original clean signals  $x_{\text{clean}}$  to generate noisy signals  $x_{\text{noisy}}$  for model testing. The synthetic noise was generated using a unified framework that supports multiple noise types, with precise control over the intensity of the noise through parameters. Specifically, we implemented the following noise types:

- **Frequency-Domain Noise:** It includes High-frequency noise, Low-frequency noise, and Broadband noise. This type is generated by taking the Fast Fourier Transform (FFT) of the original signal, generating a band-limited or broadband random Gaussian noise spectrum in the frequency domain, and then converting it back to the time domain via Inverse Fast Fourier Transform (IFFT). The frequency band division for high- and low-frequency noise is controlled by the  $r_{\rm cut}$  parameter, defined as the cutoff proportion in the frequency space.
- **Trend Noise:** Simulates slow-varying, non-periodic disturbances. This noise is generated by creating a low-order (e.g., quadratic) polynomial with random coefficients to simulate the trend component in the time series and adding it to the original signal.
- **Seasonal Noise:** Simulates periodic disturbances. This noise is generated by superimposing one or more sine waves with predefined base frequencies specified by the parameter  $f_{\text{seasonal}}$ , each having a random initial phase.

The noise intensity is precisely controlled by  $\epsilon$ , which defines the desired ratio of noise energy  $E_{\text{noise}}$  to clean signal energy  $E_{\text{clean}}$ , i.e.,  $E_{\text{noise}}/E_{\text{clean}}$ . After generating the noise, which can be denoted as noise, the noise energy is calculated and scaled accordingly to ensure that the noise added to the clean signal has a relative energy level consistent with  $\epsilon$ . The final noisy signal  $x_{\text{noisy}}$  is obtained by adding the scaled noise  $\mathbf{noise}_{\text{scaled}}$  to the original clean signal:  $x_{\text{noisy}} = x_{\text{clean}} + \mathbf{noise}_{\text{scaled}}$ .

Then, we systematically analyze the performance of the model when faced with various signal distortions. In our experiments, concretely, we fixed the  $r_{\rm cut}$  at 0.3, set  $f_{\rm seasonal}$  to  $\{1/24, 1/12\}$ , and used  $\epsilon$  values of  $\{0.1, 0.2, 0.5\}$  in different experimental groups. The results are shown in Table.22.

Comparative Study of Anti-Aliasing Strategies. To further investigate our proposed Equivalent Sampling Rate mechanism and explore efficient anti-aliasing strategies, we conducted a comparative study on the Weather dataset using a 96-step lookback to predict a 720-step horizon. We benchmarked three distinct anti-aliasing configurations:

- **DMANet** (**ESR-based**): Our proposed model, which uses the architecture-aware ESR to dynamically determine the cutoff frequency for a sharp filter.
- **DMANet\_but** (**Butterworth**): A variant where the ESR-based filter is replaced by a traditional 4th-order Butterworth low-pass filter, a well-established mathematical filter known for its maximally flat passband.
- **DMANet\_mix** (**Fusion-based**): A hybrid model that first uses ESR to partition the spectrum and then processes the high- and low-frequency bands through separate convolutional layers before fusing them, designed to explore the utility of preserved high-frequency information.
- **DMANet\_wo** (**No Filter**): A baseline variant that removes the anti-aliasing filter entirely, processing the raw input directly through the network to assess the necessity and impact of frequency-domain filtering.

The results under various noise conditions are summarized in Table.22. Overall, most of configurations demonstrate notable robustness, with only graceful performance degradation as noise intensity increases. This highlights the general effectiveness of incorporating a pre-sampling filtering stage to enhance noise resistance.

Our ablation study reveals a insight into the effectiveness of different anti-aliasing strategies. Theoretically, one might expect the Butterworth filter (DMANet\_but), with its maximally flat passband, to excel at handling low-frequency and trend noise by preserving the signal fidelity in that band Yin et al. (2024). Conversely, our ESR-based hard-cutoff filter (DMANet) should be superior against high-frequency and seasonal noise due to its removal of aliasing-prone components.

Table 22: Robustness analysis of DMANet variants under different types and intensities of synthetic noise on the Weather dataset. All experiments use a 96-step lookback to predict a 720-step horizon.

Model Variant	Noise Type	$\epsilon =$	1%	$\epsilon =$	5%	$\epsilon =$	10%
Wilder Variable	rioise Type	MSE	MAE	MSE	MAE	MSE	MAE
	Seasonal	0.343	0.339	0.342	0.341	0.341	0.343
	Trend	0.344	0.340	0.345	0.345	0.351	0.358
DMANet	All (Broadband)	0.345	0.341	0.345	0.342	0.346	0.344
	Low-Frequency	0.345	0.340	0.347	0.342	0.349	0.347
	High-Frequency	0.344	0.340	0.343	0.340	0.343	0.341
	Seasonal	0.346	0.340	0.342	0.340	0.340	0.343
	Trend	0.347	0.342	0.348	0.347	0.354	0.359
DMANet_but	All (Broadband)	0.349	0.342	0.347	0.343	0.346	0.344
	Low-Frequency	0.352	0.344	0.349	0.345	0.350	0.347
	High-Frequency	0.345	0.341	0.343	0.340	0.343	0.343
	Seasonal	0.353	0.343	0.350	0.345	0.350	0.348
	Trend	0.348	0.342	0.353	0.350	0.356	0.359
DMANet_mix	All (Broadband)	0.352	0.344	0.353	0.344	0.357	0.348
	Low-Frequency	0.354	0.345	0.351	0.345	0.353	0.348
	High-Frequency	0.358	0.346	0.353	0.345	0.352	0.347
	Seasonal	0.348	0.343	0.350	0.346	0.346	0.349
	Trend	0.348	0.343	0.351	0.351	0.357	0.361
$DMANet\_wo$	All (Broadband)	0.350	0.344	0.352	0.346	0.347	0.345
	Low-Frequency	0.355	0.346	0.355	0.349	0.353	0.349
	High-Frequency	0.350	0.343	0.350	0.345	0.351	0.344

Interestingly, our empirical results in Table.22 show that while performance is competitive on seasonal and high-frequency noise, DMANet consistently and significantly outperforms DMANet\_but on trend and low-frequency noise. This seemingly counter-intuitive result highlights a critical limitation of applying classical filters naively within a deep learning pipeline. While the Butterworth filter is static and optimally preserves its predefined passband, it is architecture-agnostic. It may still pass frequencies that, while low, are too high for the subsequent strided convolution to process without aliasing. In contrast, our ESR-based approach is architecture-aware. It does not aim to be a perfect mathematical filter in isolation; its sole purpose is to perfectly prepare the signal for the next layer. By dynamically calculating a precise cutoff based on the network's own parameters, it ensures that no aliasing occurs at any stage, even if this means a slightly more aggressive filtering. This architectural synergy proves to be more practically effective.

Furthermore, the fusion-based DMANet\_mix consistently underperforms the other two variants. This result empirically supports our design rationale for employing a strict cutoff strategy: for a lightweight model, it is more effective to concentrate its limited capacity on core, learnable patterns rather than attempting to fit the complex and often noisy dynamics of high-frequency information. As observed in prior work like FITS Xu et al. (2024), removing a significant portion of high-frequency components largely preserves a time series' dominant trends. The poor performance of DMANet\_mix indicates that simply preserving and processing this high-frequency content is less effective than principled filtering, likely because this band is dominated by noise that the model cannot distinguish from a true signal.

Collectively, these results validate that our ESR-based approach provides the most robust and adaptive solution. By dynamically and precisely removing only the frequencies that would cause aliasing, it not only focuses the model on the most decisive, learnable patterns but also achieves this with superior adaptability compared to static classical filters, all without the need for manual filter design.

# G More details of Our Method

#### G.1 THE RATIONALE FOR THE EMBEDDING FIRST ARCHITECTURE

A critical challenge in multi-scale time series analysis is the fusion of features from different scales without introducing signal distortion. A common approach, which we term multi-scale first, involves downsampling the raw signal and then embedding each scale. However, this seemingly intuitive process hides a significant pitfall: the upsampling step required for feature fusion inevitably causes spectral distortion due to its reliance on a limited reconstruction basis. To circumvent this fundamental issue, our DMANet adopts a principled **embedding first** architecture, ensuring all operations are conducted with high fidelity within a unified feature space.

# G.1.1 THE PITFALL OF PREMATURE MULTI-SCALE DECOMPOSITION

The multi-scale first approach, seen in models like TimeMixer Wang et al. (2024a), begins by decomposing the raw signal X into a set of time series  $\{X_m \in \mathbb{R}^{C \times s_m}, s_m < L\}$ . While feasible, the core flaw lies in the subsequent step of unifying these scales for feature fusion. To restore the original length L, each short sequence  $s_m$  must be upsampled using a linear layer,  $g_m : \mathbb{R}^{s_m} \to \mathbb{R}^L$ . This process is inherently problematic due to its limited representational capacity:

- Limited Basis Vectors: The weight matrix  $W \in \mathbb{R}^{L \times s_m}$  of the upsampling layer provides only  $s_m$  column vectors. These vectors form the *entire basis* available to reconstruct the output signal. Consequently, all reconstructed signals are confined to a very small,  $s_m$ -dimensional subspace of the target space  $\mathbb{R}^L$ .
- **Deformed Basis Vectors:** To approximate the diverse signals in the training data from this constrained basis, the model is forced to learn complex, **non-smooth**, **and oscillatory** basis vectors as a poor compromise.
- **Inevitable Spectral Distortion:** When a signal is reconstructed as a linear combination of these deformed basis vectors, it unavoidably inherits their unnatural properties. This leads to severe **spectral distortion**, corrupting the signal's fidelity and polluting the final prediction.

# G.1.2 EMBED FIRST: A PRINCIPLED APPROACH IN A UNIFIED FEATURE SPACE

Our DMANet architecture is designed to completely avoid the aforementioned reconstruction problem by first establishing a unified workspace for all operations.

- 1. **Defining a Unified Workspace:** We begin by projecting the information-complete raw signal  $X \in \mathbb{R}^{C \times L}$  into a new feature basis space using a linear layer. This generates a feature sequence  $X' \in \mathbb{R}^{C \times T}$ , creating a unified and consistent workspace for all subsequent synergistic operations.
- High-Fidelity Operations: Within this consistent feature space, all core operations are performed in a principled manner:
  - Downsampling: Our anti-aliasing downsampling module pre-filters features in the frequency domain before reducing resolution, preventing information aliasing and ensuring reliable feature transfer across scales.
  - Upsampling: To restore resolution for feature fusion, we employ zero-padding in the frequency domain. This method is equivalent to ideal interpolation and relies on the Fourier basis (sines and cosines)—a fixed, universal, and complete orthogonal basis. Adhering to the Nyquist-Shannon sampling theorem, this ensures the smoothest possible reconstruction, free from the uncontrolled high-frequency artifacts generated by the alternative approach.

By ensuring all features are derived and processed with high-fidelity operations within the same basis space, we maintain inherent consistency and make feature fusion fundamentally more reliable.

#### G.1.3 EMPIRICAL VALIDATION

To validate our theoretical analysis, we conducted a comprehensive comparison between our DMANet (Embedding First) and the alternative architecture (Multi-Scale First). We also performed an ablation study by removing the initial embedding module (w/o embed) to verify the effectiveness of operating within a latent space.

The results in Table.23 provide strong empirical support for our design.

Table 23: Comparative analysis and ablation study for the Embedding First architecture. Our full DMANet model is compared against the Multi-Scale First approach and a variant without the initial embedding module.

Model	Metric	ETTh1	ETTm1	Weather	Elect	Wiki	ILI	Unemp	Dowjone
DMANet (Embedding First)	MSE	0.428	0.373	0.236	0.172	6.506	1.763	0.064	11.957
	MAE	0.429	0.385	0.262	0.265	0.393	0.824	0.146	0.850
Multi-Scale First	MSE	0.441	0.385	0.242	0.181	6.555	2.097	0.074	12.420
	MAE	0.435	0.391	0.268	0.273	0.407	0.858	0.167	0.860
w/o embed	MSE	0.436	0.389	0.249	0.188	6.551	2.084	0.073	12.330
	MAE	0.427	0.392	0.274	0.277	0.406	0.879	0.163	0.857

The Multi-Scale First approach consistently underperforms our model. This performance gap is a direct, practical consequence of the spectral distortion introduced by its unprincipled, basis-limited reconstruction step.

Furthermore, we acknowledge that the initial linear mapping carries a potential risk of losing some temporal dependencies. This is a deliberate design choice, and its justification is twofold. First, we incorporate a learnable positional encoding to preserve crucial temporal context. Second, as our ablation study will demonstrate, the benefits of analyzing the series in a latent space—where patterns are more suitable for anti-aliasing and feature extraction—outweigh the alternative of operating directly on the raw signal. The placeholder for the w/o embed results in Table.23 will provide strong evidence for this superiority.

#### G.2 Principled Anti-Aliasing via Dynamic Frequency Cutoff

For any given downsampling layer l in DMANet, its anti-aliasing operation is the application of a low-pass filter with a mathematically-derived, strict cutoff frequency. Given the layer's convolutional parameters—kernel size k, stride s, and channel ratio c—we first calculate its Effective Sampling Ratio (ESR $^l$ ) using Equation.3 to determine its true signal processing capability. This allows us to establish a new Nyquist frequency,  $f_{\rm Nyquist}^l$ . As shown in Equation.4, all frequency components above this threshold are strictly zeroed out via a frequency-domain mask. Our method does not partially retain or vaguely attenuate high-frequency components; it employs a principled cutoff scheme where the threshold is dynamically determined for each layer.

#### G.2.1 THE RATIONALE: FOCUSING ON LEARNABLE CORE PATTERNS

The core motivation for this strict cutoff strategy is to concentrate the model's capacity on learnable, core patterns. High-frequency information in time series often contains significant noise or stochastic fluctuations that are difficult to model, and typically exceed the learning capacity of a lightweight model. Attempting to fit these complex dynamics can hinder the model from capturing the more decisive, underlying trends.

As observed in prior work like FITS Xu et al. (2024), removing a significant portion of high-frequency components largely preserves the overall shape and dominant trends of a time series. Our strategy builds on this insight: by proactively simplifying the learning task, we focus the model's limited capacity on the low-frequency periodic and trend patterns that are most critical for the forecasting task, thereby achieving both efficient and accurate predictions.

#### G.2.2 DYNAMIC ADAPTABILITY AND PARAMETER-FREE DESIGN

A key advantage of our method is its dynamic nature. In a complex multi-scale architecture, different layers may employ varying downsampling parameters (k, s, c). Our framework automatically

derives a matching, optimal cutoff frequency for each specific layer. This ensures that the antialiasing protection remains effective and theoretically grounded across any architectural variation, eliminating the need for tedious, manual parameter tuning required by classical filters or the randomness of heuristic approaches.

Furthermore, this framework possesses theoretical flexibility. By adjusting the convolution parameters, the ESR can be controlled to retain more, or even all, frequency components. For instance, if parameters are set such that ESR = 1 (e.g.,  $s = \min(K, C_{\text{out}}/C_{\text{in}})$ ), the cutoff frequency matches the original signal's Nyquist frequency, meaning no valid frequency components are attenuated.

#### G.2.3 Addressing the High-Frequency Information Trade-off

We acknowledge that this design is built upon a core trade-off: we filter high-frequency components to prevent aliasing at the cost of potentially discarding useful information. This is a deliberate choice motivated by the efficiency and robustness goals for a lightweight model.

We also recognize that high-frequency information can be critical in certain scenarios, such as forecasting sharp spikes or in contexts where high-frequency harmonics are themselves key features. It is precisely for this reason that we deliberately conducted extensive supplementary experiments across a diverse range of domains (including Electricity, Weather, Transportation, Health, Web, Market, Energy, Society, Finance, etc.). The goal was to proactively probe the application boundaries of our method and provide a clear reference for its practical use.

To further investigate this trade-off, we will conduct a controlled experiment comparing our strict cutoff method with an alternative that handles frequencies differently. As shown in Table.24, we will compare our standard DMANet against a variant where, after identifying the cutoff frequency, both the low-frequency and the zeroed-out high-frequency components are independently passed through linear layers and then fused. This will help quantify the practical impact of the information contained in the high-frequency bands.

Table 24: Ablation study on the handling of high-frequency components. We compare our strict cutoff method with a variant that uses linear fusion for high and low frequencies.

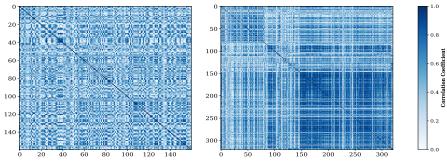
Model	Metric	ETTh1	ETTm1	Weather	Elect	Wiki	ILI	Unemp	Dowjone
DMANet (Strict Cutoff)	MSE	0.428	0.373	0.236	0.172	6.506	1.763	0.064	11.957
	MAE	0.429	0.385	0.262	0.265	0.393	0.824	0.146	0.850
DMANet_mix (Fusion-based)	MSE	0.434	0.374	0.237	0.171	6.528	1.986	0.076	12.288
	MAE	0.433	0.385	0.263	0.265	0.398	0.867	0.161	0.858

In summary, DMANet's anti-aliasing employs a precise, dynamic cutoff strategy tailored to each layer's actual sampling capability. This design combines theoretical robustness with practical, parameter-free convenience, and its effectiveness is validated through extensive empirical analysis.

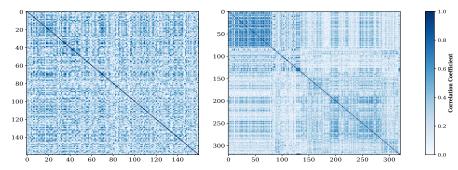
# H MORE DETAILS OF DEPENDENCY MODELING

We visualize the temporal dependencies and channel-wise relationships within a batch of the Electricity in Figure.7 and Figure.8 and for Weather in Figure.9 and Figure.10, comparing their states before and after processing by DMANet's components. To further illustrate the differences between scenarios with and without the anti-aliasing filter, we selected the Electricity dataset to visualize the temporal dependency differences of upsampling before and after applying the anti-aliasing filter in Figure.11, as well as the channel dependency correlation differences of downsampling with and without the anti-aliasing filter in Figure.12.

DMANet tends to leverage more effective dependencies to capture future trends. Comparing the cases with and without the anti-aliasing filter, the figures reveal that the pre-processing anti-aliasing operation, acting as a low-pass filter, smooths or attenuates fine-grained dependencies that are susceptible to aliasing during sampling. This process helps to highlight the main temporal dependency patterns and channel relationships. Furthermore, convolution, leveraging its local receptive field, focuses on local patterns at neighboring time points. Thus, the combination of filtering and convolutional downsampling effectively extracts stable temporal features.



(a) w / Pre-Sampling filtering. Left: Downsample corr: 0.432, Right: Upsample corr: 0.534.



(b) w/o Pre-Sampling filtering. Left: Downsample corr: 0.311, Right: Upsample corr: 0.254.

Figure 7: Visualization for channel dependency modeling on Electricity in the first layer of the second multiscale encoder block.

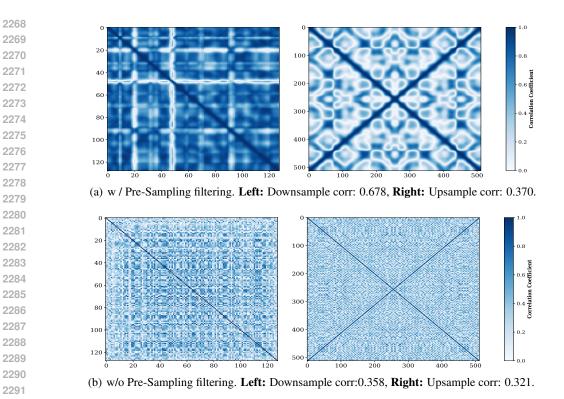


Figure 8: Visualization for temporal dependency modeling on Electricity in the first layer of the second multiscale encoder block.

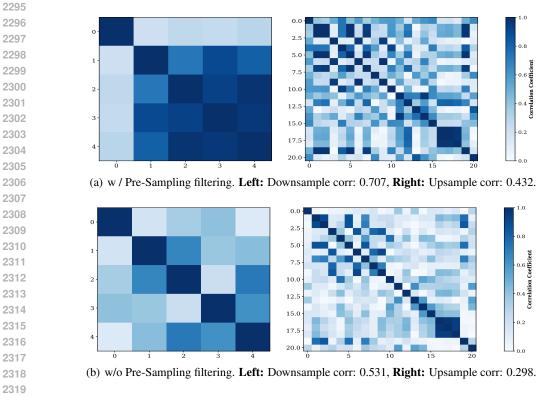


Figure 9: Visualization for channel dependency modeling on Weather in the first layer of the first multiscale encoder block.

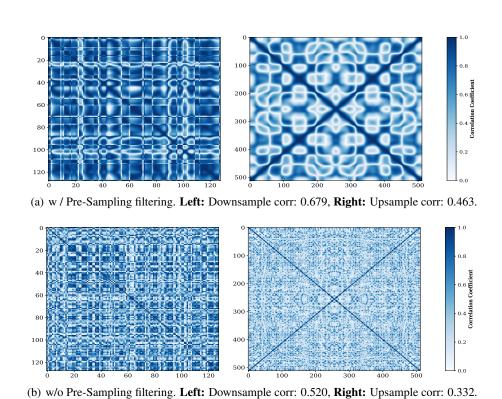


Figure 10: Visualization for temporal dependency modeling on Weather in the first layer of the first multiscale encoder block.

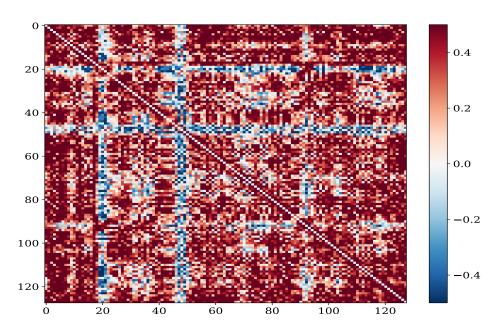


Figure 11: Temporal dependency differences in up-sampling with or without the application of an anti-alias filter on Electricity. Red indicates increased dependency after use anti-alias filter.

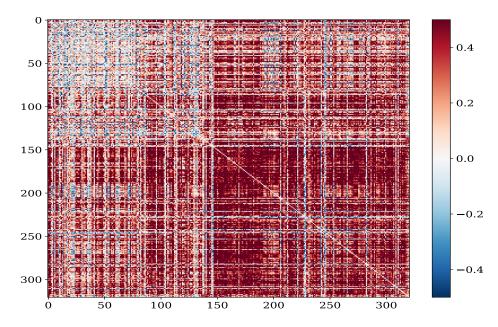


Figure 12: Channel dependency differences in down-sampling with or without the application of an anti-alias filter on Electricity. Red indicates increased dependency after use anti-alias filter.

# I THE USE OF LARGE LANGUAGE MODELS

Large Language Models were employed as general-purpose assistive tools throughout the research process. Specifically, LLMs were used to polish the language and improve the readability of this manuscript, including refining grammar, improving clarity, and restructuring sentences for better readability. The authors take full responsibility for the content of this paper.