# Introducing Spectral Attention for Long-Range Dependency in Time Series Forecasting

**Bong Gyun Kang**[1][*]     **Dongjun Lee**[1][*]     **HyunGi Kim**[2]     **DoHyun Chung**[3]

**Sungroh Yoon**[1,2][†]

[1] Interdisciplinary Program in Artificial Intelligence, Seoul National University
[2] Department of Electrical and Computer Engineering, Seoul National University
[3] Department of Future Automotive Mobility, Seoul National University

## Abstract

Sequence modeling faces challenges in capturing long-range dependencies across diverse tasks. Recent linear and transformer-based forecasters have shown superior performance in time series forecasting. However, they are constrained by their inherent inability to effectively address long-range dependencies in time series data, primarily due to using fixed-size inputs for prediction. Furthermore, they typically sacrifice essential temporal correlation among consecutive training samples by shuffling them into mini-batches. To overcome these limitations, we introduce a fast and effective Spectral Attention mechanism, which preserves temporal correlations among samples and facilitates the handling of long-range information while maintaining the base model structure. Spectral Attention preserves long-period trends through a low-pass filter and facilitates gradient to flow between samples. Spectral Attention can be seamlessly integrated into most sequence models, allowing models with fixed-sized look-back windows to capture long-range dependencies over thousands of steps. Through extensive experiments on 11 real-world time series datasets using 7 recent forecasting models, we consistently demonstrate the efficacy of our Spectral Attention mechanism, achieving state-of-the-art results.

## 1   Introduction

Time series forecasting (TSF) stands as a core task in machine learning, ubiquitous in our lives through applications such as weather forecasting, traffic flow estimation, and financial investment. Over time, TSF techniques have evolved from statistical [5, 10, 18] and machine learning approaches [2, 9, 20] to deep learning models like Recurrent Neural Networks [15, 25, 41] and Convolution based Networks [12, 45]. Following the success of Transformers [44] in various domains, Transformer-based models have also become mainstream in the time series domain [23, 27, 36, 49, 54, 55, 56]. Recently, methodologies based on Multi-layer Perceptron have received renewed attention [7, 8, 28, 52]. However, despite the advancements, long-term dependency modeling in TSF remains challenging [53].

Unlike image models, where data are randomly sampled from the image distribution and are thus independent of each other [16], TSF models sample data from the continuous signal, dependent

---

[*]Denotes equal contribution. {luckypanda, elite1717}@snu.ac.kr
[†]Corresponding author. sryoon@snu.ac.kr

Figure 1: (a) Training data are sampled for each time step from continuous sequences, exhibiting high temporal correlations. (b) Conventional approaches simply ignore this temporal information with a shuffled batch. (c) We address the temporal correlation between the samples for the first time, enabling the model to consider long-range dependencies that surpass the look-back window.

on the time variable $t$ as shown in Figure 1a. This leads to a high level of correlation between each training sample, which consists of a fixed-sized look-back window before $t$ (as input) and the subsequent prediction sequence after $t$ (as label). Therefore, the conventional approach of shuffling the consecutive samples into mini-batches deprives the model of utilizing the crucial inherent temporal correlation between the samples. This restricts the model's consideration to only a fixed-size look-back window for sequence modeling, limiting the ability to address long-range dependencies (Figure 1b).

Recent studies pointed out that simply increasing the look-back window leads to substantially detrimental effects such as increased model size and longer training and inference times [53]. This is particularly challenging for transformer forecasters, which exhibit quadratic time/memory complexity [36, 44, 54], and even for efficient models using techniques like Sparse Attention, which have O(nlogn) complexity [23, 27, 55]. Furthermore, if the commonly used look-back window of 96 is extended fivefold, the model can only utilize time steps of less than 500, making it difficult to consider long-range dependencies spanning thousands or the entire dataset. Also, increasing the look-back window may not be beneficial, often leading to decreased performance [53], highlighting the fact that current models are not sufficient in capturing long-range dependencies.

To address this limitation, we propose Spectral Attention, which can be applied to most TSF models and enables the model to utilize long-range temporal correlations in sequentially obtained training data. With the stream of consecutive training samples (Figure 1c), Spectral Attention stores an exponential moving average (EMA) of the activations with various smoothing factors at each time step. This serves as a low-pass filter, inherently embedding long-range information over a thousand steps. Attention is then applied to the stored EMA activations of various smoothing factors (low-pass filters with different frequencies), enabling the model to learn which periodic trends to consider when predicting the future, thereby enhancing its performance. Spectral Attention is even applicable to models such as iTransformer [31], which do not preserve the temporal order of time series data internally.

We further extend Spectral Attention, where computations depend on the activation of the previous timesteps, to Batched Spectral Attention, enabling parallel training across multiple timesteps. This extension makes Spectral Attention faster and more practical and allows for the direct utilization of temporal relationships among consecutive data samples within a mini-batch in the training base model. In Batched Spectral Attention, the EMA is unfolded over time to perform Spectral Attention simultaneously across multiple time steps. This unfolding allows gradients at time $t$ to propagate through the Spectral Attention module to the previous time step within the mini-batch, achieving effects similar to Backpropagation Through Time (BPTT) [47] and extends the model's effective input window.

**Our approach** preserves the base TSF model architecture and learning objective while enabling the model to leverage long-term trends spanning thousands of steps. By effectively utilizing the temporal correlation of training samples, our method allows gradients to flow back in time beyond the look-back window, extending to the entire mini-batch. Also, our method requires little additional training time and memory. We conducted experiments on 7 recent TSF models and 11 real-world datasets and demonstrated consistent performance enhancement in all architecture, achieving state-of-the-art results. We summarize the main contributions as follows:

- We propose Spectral Attention, which addresses long-range dependencies spanning thousands of steps through frequency filtering and attention mechanisms, leveraging the temporal correlation among the consecutive samples.

- We propose Batched Spectral Attention, which enables parallel training across multiple timesteps and expends the effective input window, allowing the gradient to flow through time within the mini-batch.

- Batched Spectral Attention is applicable to most existing TSF models and practical in real-world scenarios with minimal additional memory and comparable training time. Also, it allows finetuning with a trained TSF model.

- Batched Spectral Attention demonstrates consistent model-agnostic performance improvements, particularly showcasing superior performance on datasets with significant long-term trend variations.

## 2 Related Works

**Classic TSF models.** Statistical TSF methods, such as ARIMA [35], Holt-Winters [18], and Gaussian Process [10], assume that temporal variations adhere to predefined patterns. However, their practical applicability is largely limited by the complex nature of real-world data. Machine learning approaches, such as Support Vector Machines [4] and Random Forests [14] have proven to be effective even compared to early artificial neural networks [13, 17, 40]. Convolutional network-based methods leverage convolution kernels to capture temporal variations sliding along the temporal dimension [12, 45]. Recurrent neural network (RNN) grasp changes over time via state transitions across different time steps. [25, 41]. However, RNN-based models exhibit limitations in modeling long-range dependencies due to challenges such as vanishing gradients and error accumulation [30, 43]. Recently, transformer and linear-based models have emerged as alternatives, demonstrating superior performance compared to recurrent models [46, 53].

**Transformer and Linear based models.** Transformer-based models [44] address temporal relationships between time points using the attention. LogSparseTransformer [27], Reformer [23], and Informer [55] have been proposed to make the Transformer architecture efficient, addressing the quadratic time complexity. The Autoformer [49] incorporates series decomposition as an inner block of Transformer and aggregates similar sub-series by utilizing the Auto-Correlation mechanism. PatchTST [36] introduces patching, a channel-independent approach that processes each variable separately and focuses on cross-time attention. Crossformer [54] utilizes a channel-dependent approach to learn cross-variate dependencies. This is achieved through the use of cross-time and cross-dimension attention. iTransformer [31] applies attention and FFN in an inverted way, where attention handles correlations between channels and FFN handles the temporal information. Recently, to address Transformers' potential difficulties in capturing long-range dependencies [53], methodologies based on the linear model and Multi-Layer Perceptron (MLP) structures have emerged. DLinear [53] utilizes the decomposition method introduced in Autoformer and predicts by adding the output of two linear layers for each seasonal and trend element. TiDE [7] proposes an architecture based on MLP residual blocks that combines information from dynamic and static covariates with a look-back window for encoding, followed by decoding. TSMixer [8] performs forecasting by repeatedly mixing time and features using an MLP. RLinear [28] comprises of a single linear layer with RevIN [21] for normalization and de-normalization.

**Frequency-utilizing models.** Using the frequency domain for TSF is a well-established approach [3, 42, 19]. Conventional approaches leverage frequency information during the preprocessing stage [37] or decompose time series based on frequency filtering [39]. In deep TSF models, research has also explored architectural advancements that are aware of the frequency information. SAAM [34], which is applicable to RNNs, performs FFT and autocorrelation on the input signal. WaveFroM [51] uses discrete wavelet transform (DWT) to project time series into wavelet domains of various scales and performs forecasting through graph convolution and dilated convolution. FEDformer [56] adopts a mixture-of-experts strategy to refine the decomposition of seasonal and trend components and introduce sparse attention mechanisms in the frequency domain. TimesNet [48] transforms 1D time series into 2D tensors utilizing multi-periodicity by identifying dominant frequencies through Fourier Transform, modeling temporal variations effectively. FreTS [52] leverages frequency-domain MLPs to achieve global signal analysis and compact energy representation, addressing the limitations of

Figure 2: (a) Plug-in Spectral Attention (SA) module takes a subset of intermediate feature $F$ and returns $F'$ with long-range information beyond the look-back window. The model is trained end-to-end, and gradients flow through the SA module. (b) To capture the long-range dependency, SA stores momentums of feature $F$ generated from the sequential inputs. Multiple momentum parameters $\alpha_i$ capture dependencies across various ranges. (c) SA module computes $F'$ by attending multiple low-frequency ($M^{\alpha_i}$) and high-frequency ($F - M^{\alpha_i}$) components and feature ($F$) using learnable Spectral Attention Matrix (SA-Matrix)

point-wise mappings and information bottlenecks in conventional MLP-based methods. FITS [50] employs interpolation within the complex frequency domain to construct a concise and robust model. While these models leverage frequency information, they are limited in modeling long-range dependencies, as the frequency conversion is confined to the look-back window. On the other hand, our Spectral Attention is the first to achieve long-range dependency modeling beyond the look-back window by incorporating consecutive data streams during model training.

## 3 Methods

**Problem Statement.** In multivariate time series forecasting, time series data is given $\mathbb{D}_T$ : $\{x_1, ..., x_T\} \in \mathbb{R}^{T \times N}$ at time $T$ with $N$ variates. Our goal is, at arbitrary future time $t > T$, to predict future $S$ time steps $Y_t = \{x_{t+1}, ..., x_{t+S}\} \in \mathbb{R}^{S \times N}$. To achieve this goal, TSF model $f$ utilizes length $L$ look-back window as input $X_t = \{x_{t-L+1}, ..., x_t\} \in \mathbb{R}^{L \times N}$ making prediction $P_t = f(X_t), P \in \mathbb{R}^{S \times N}$.

Model is trained with the training dataset $D_T = \{(X_t, Y_t) | L \le t \le T - S\}$. While conventional methods typically randomly sample each $X, Y$ from $D_T$ to constitute the mini-batch, we utilize sequential sampling to incorporate temporal correlations between samples into the learning process.

### 3.1 Spectral Attention

Spectral Attention ($SA$) can be applied to every TSF model that satisfies the aforementioned problem statement. This base TSF model is represented by $P = f(X)$, and $SA$ can be applied to arbitrary activation $F$ within the model. The base model can be reformulated as $P = f_2(F, E)$ and $F, E = f_1(X)$. $F, E$ are intermediate state and SA module takes an arbitrary subset $F$ as input and transforms it into $F'$ of the same size; $F' = SA(F), P' = f_2(F', E)$. The resulting SA plugged model $f_{SA}$ is depicted in Figure 2a.

With $X_t$ as the base model input, $SA$ takes $D$-dimensional feature vector $F_t \in \mathbb{R}^D$ as input. $SA$ updates the exponential moving average (EMA) $M_t \in \mathbb{R}^{K \times D}$ of $F_t$ in its internal memory with the $K$ smoothing factors $\{\alpha_1, ..., \alpha_K\} \in \mathbb{R}^K$ ($\alpha_1 < \cdots < \alpha_K$) as shown in Figure 2b.

$$M_{t+1}^{k,i} = \alpha_k \times M_t^{k,i} + (1 - \alpha_k) \times F_t^i \tag{1}$$

EMA retains the trend of features over long-range time periods based on the smoothing factor. It operates as a low-pass filter, with the -3db (half) cut-off frequency of $freq_{cut} = \frac{1}{2\pi} cos^{-1} \left[ 1 - \frac{(1-\alpha)^2}{2\alpha} \right]$, effectively preserving the trend over 6,000 period with $\alpha = 0.999$ .

To represent high-frequency patterns contrasting with the low-pass filtered long-range pattern $M_t$, we generated $H_t \in \mathbb{R}^{K \times D}$ by subtracting $M_t$ from $F_t$.

$$H_t^{k,i} = F_t^i - M_t^{K-k-1,i} \tag{2}$$

4

Figure 3: BSA module takes a sequentially-sampled mini batch $\{X_t, ...X_{t+B-1}\}$ and computes the corresponding EMA momentums $\{M_t, ...M_{t+B-1}\}$ over time. This is done via single matrix multiplication enabling parallelization. We made the momentum parameter $\alpha_i$ learnable, allowing the model to directly learn the periodicity of the information essential for the future prediction.

$SA$ contains learnable parameters: *sa-matrix* $\in \mathbb{R}^{(2K+1)\times D}$, which learns what frequency the model should attend to for each feature. $2 \times H_t, F_t, 2 \times M_t$ are concatenated on dimension 0, resulting in $\mathbb{R}^{(2K+1)\times D}$, which is then weighted summed with *sa-matrix* on dimension 0, generating output $F'_t$ (Figure 2c).

$$F'_t = sum(softmax(\text{sa-matrix}, dim\,0) \cdot concat((2 \times H_t, F_t, 2 \times M_t), dim\,0), dim\,0) \quad (3)$$

The *sa-matrix* is initialized so that $softmax(\text{sa-matrix})$ resembles a Gaussian distribution on axis 0. This results in symmetric value on axis 0 (*sa-matrix*$^{K+1-i}$ = *sa-matrix*$^{K+1+i}$) and makes $SA$ an identity function on initialization ($\because H^k + M^{K-k-1} = F$).

$$F = SA_{init}(F) \quad (4)$$

$SA$ allows the model to attend to multiple frequencies of its feature signal, enabling it to focus on either long-range dependencies or high-frequency patterns as needed and shift the feature $F$ distribution on the frequency domain. By initializing $SA$ as an identity function, the model can be fine-tuned with the already trained base model, allowing efficient implementation.

### 3.2 Batched Spectral Attention

Batched Spectral Attention ($BSA$) enables batch training over multiple time steps. The main concept involves unfolding EMA, which facilitates gradients to flow across consecutive samples in a mini-batch, akin to BPTT. This enables efficient parallel training and promotes the model to extract long-range information beneficial for future prediction, extending the effective look-back window. The overall flow of $BSA$ is depicted in Figure 3.

With mini-batch of size $B$, consecutive samples $X_{[t,t+B-1]} = \{X_t, ...X_{t+B-1}\} \in \mathbb{R}^{B \times S \times N}$ are given as input. Following aforementioned $SA$ setting, $BSA$ takes $F_{[t,t+B-1]} = \{F_t, ...F_{t+B-1}\} \in \mathbb{R}^{B \times D}$ as input. In the next step, $BSA$ utilizes $F_{[t,t+B-1]}$ and the stored $M_t \in \mathbb{R}^{K \times D}$ to generate $M_{t+b}(0 \le b \le B)$ by unfolding the Equation 1.

$$M_{t+b}^{k,i} = \alpha_k^b \times M_t^{k,i} + (1-\alpha_k)\alpha_k^{b-1} \times F_t^i + \cdots + (1-\alpha_k) \times F_{t+b-1}^i \quad (5)$$

This equation can be transformed to calculate $M_{[t,t+B]} \in \mathbb{R}^{(B+1)\times K \times D}$ in parallel as follows.

$$M_{[t,t+B]}^{:,k,i} = lower\text{-}triangle(A^k) \times concat((M_t^{k,i}, F_{[t,t+B-1]}^{:,i}), dim\,0) \quad (6)$$

$$A \in \mathbb{R}^{K \times (B+1) \times (B+1)}, \; A^{k,p,q} = (1-\alpha_k)^{I\{q>0\}}\alpha_k^{p-q} \quad (7)$$

$A$ refers to unfolding matrix and $I$ refers to indicator function. $M_{t+B}$ is stored in $BSA$ for the next mini-batch input. $F'_{[t,t+B-1]}$ is computed in parallel, similar to Equation 2 and 3, using $F_{[t,t+B-1]}$, $B_{[t,t+B-1]}$, and *sa-matrix*$\in \mathbb{R}^{(2K+1)\times D}$. The $lower\text{-}triangle$ function prevents gradients from the past timestep from flowing into future models, aligning with the characteristics of time-series data.

5

At the beginning of each training epoch, $M_0$ is initialized to $F_0$ for all $\alpha$, enhancing stability for subsequent EMA accumulation. Since $SA$ is proposed to address long-range dependencies in training, it lacks sufficient information in the early stages when not enough steps have been seen. Therefore, for the stability of training, we linearly warm up the learning rate for the first $1/(1 - max(smoothing\ factors))$ timesteps at the beginning of each training epoch. The overall training of both the base model and $BSA$ is conducted according to Algorithm 1.

In SA, the *smoothing factors* $\{\alpha_1, ..., \alpha_K\}$ $\in \mathbb{R}^K$ were given as scalar values, whereas in $BSA$, they are expressed by learnable parameters. This is because $BSA$ can utilize additional past information in training beyond the look-back window by incorporating a batch-sized time window, allowing it to determine the extent of long-range dependency required for training. To keep the smoothing factors between 0 and 1, we initialized learnable parameters by applying an inverse sigmoid to the initial smoothing factors and then applied a sigmoid function in training.

So far, we assume the feature $F$ from the base model as a vector. However, the output of the intermediate layers of the model is often represented as a tensor with two or more dimensions. In real practice, we use additional channel dimensions in $BSA$ to process the activation tensor, which acts as applying multiple $BSA$ modules simultaneously.

---

**Algorithm 1** Batched Spectral Attention (1 epoch)

---

**Input:** Trained up to $(E-1)$th epoch
TSF model $f_\theta$, $BSA$ with *sa-matrix*, *smoothing factors*
Train data: $\mathfrak{D}_{tr} = \{(X_0, Y_0), ..., (X_{tr}, Y_{tr})\}$
Valid data: $\mathfrak{D}_{val} = \{(X_{tr+1}, Y_{tr+1}), ..., (X_{val}, Y_{val})\}$
mini-batch size: B
**Train:**
$st = 0, ed = B - 1$
initialize $M_0$ in $BSA$ with $X_0$ and $f_\theta$ $(:= f_{2\theta} \cdot f_{1\theta})$
**for** $X_{[st,ed]}, Y_{[st,ed]}$ in $\mathfrak{D}_{tr}$ **do** {: *train phase*}
    $P_{[st,ed]} \leftarrow f_{2\theta} \cdot BSA \cdot f_{1\theta}(X_{[st,ed]})$
    $\mathcal{L} = \mathcal{L}_{\text{mse}}(P_{[st,ed]}, Y_{[st,ed]})$
    Compute $\nabla\mathcal{L}$ and adjust learning rate
    update $f_\theta$, *sa-matrix*, *smoothing factors*
    $st+ = B, ed = min(ed + B, tr)$
**end for**
$st = tr + 1, ed = tr + B$
**for** $X_{[st,ed]}, Y_{[st,ed]}$ in $\mathfrak{D}_{val}$ **do** {: *Validation phase*}
    $P_{[st,ed]} \leftarrow f_{2\theta} \cdot BSA \cdot f_{1\theta}(X_{[st,ed]})$
    $\mathcal{L} = \mathcal{L}_{\text{mse}}(P_{[st,ed]}, Y_{[st,ed]})$
    accumulate $\mathcal{L}$ and calculate validation loss $\mathcal{L}_{val}$
    $st+ = B, ed = min(ed + B, val)$
**end for**
**Output:** Trained up to $E$ th epoch
$f_\theta$, *sa-matrix*, $\{\alpha_1, ..., \alpha_K\}$, $\mathcal{L}_{val}$ for $E$ th epoch

---

### 3.3 Consecutive dataset split

In the TSF scenario, the entire time series data $\{x_1, ..., x_{T_{end}}\}$ is divided into train, validation, and test sets in chronological order. Let $T_{tr}$ and $T_{val}$ denote the last time in the train and validation data respectively. Model training occurs using data for t in $[1, T_{tr}]$, while model selection for the best model can utilize data for t in $[T_{tr+1}, T_{val}]$. The test set comprises predicting time steps $[T_{val+1}, T_{end}]$, which are not accessible during training or validation. However, since each training sample consists of the look-back window of size L and a prediction window of size S, the training input samples are restricted to $X_{[L,T_{tr}-S]}$. Validation samples and test input samples range from $X_{[T_{tr},T_{val}-S]}$, and from $X_{[T_{val},T_{end}-S]}$, respectively. While this approach is plausible for independent data like images, it is unnatural for sequential data, as it leaves unreachable gaps $(X_{[T_{tr}-S+1,T_{tr}-1]}, X_{[T_{val}-S+1,T_{val}-1]})$, undermining the consecutive characteristics. We filled the missing gaps, making train, validation, and test sets consecutive so that our $BSA$ model could update momentum continuously. For fair evaluation, added samples were not used for either model training, validation, or performance assessment. Detailed explanations of model validation and evaluation are provided in Appendix A.3. The full code is available at https://github.com/DJLee1208/BSA_2024.

## 4 Experiments

We first evaluate BSA using state-of-the-art TSF models and various real-world time series forecasting scenarios in section 4.1. Next, to demonstrate that BSA effectively addresses long-range dependencies and is robust to distribution shift, we perform experiments on synthetic signals of various frequencies in section 4.2. Finally, we analyze the BSA's performance variations depending on the insertion sites within the base model, examine computation and memory costs, and conduct an ablation study in section 4.3.

**Datasets.** We use eleven real-world public datasets: Weather, Traffic, ECL, ETT (4 sub-datasets; h1, h2, m1, m2), Exchange, PEMS03, EnergyData, and Illness [6, 26, 29, 49]. In the Illness dataset,

Table 1: Forecasting results averaged across prediction lengths $S \in \{96, 192, 336, 720\}$ and three seeds. Illness and Exchange datasets use different prediction lengths due to short data lengths. Higher performance between the base model and BSA is bolded. The Avg. column shows the mean across all datasets. Red indicates p-value < 0.05 in paired t-test. Abbreviations are as follows, We.; Weather, Tr.: Traffic, Eh1: ETTh1, Eh2: ETTh2, Em1: ETTm1, Em2: ETTm2, Ex.:Exchange, PE.: PEMS03, En.: EnergyData, Il.: Illness. Full results are provided in Appendix B.1.

| Model | | Metric | We. | Tr. | ECL | Eh1 | Eh2 | Em1 | Em2 | Ex. | PE. | En. | Il. | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DLinear | base | MSE | 0.244 | 0.737 | 0.227 | 0.529 | 0.349 | 0.462 | 0.248 | 0.088 | 0.436 | 0.874 | 2.849 | 0.640 |
| | | MAE | 0.308 | 0.451 | 0.325 | 0.507 | 0.406 | 0.447 | 0.332 | 0.211 | 0.502 | 0.683 | **1.111** | 0.480 |
| | BSA | MSE | **0.220** | **0.691** | **0.217** | **0.527** | **0.341** | **0.441** | **0.221** | **0.081** | **0.384** | **0.838** | **2.797** | **0.614** |
| | | MAE | **0.281** | **0.435** | **0.316** | **0.507** | **0.398** | 0.447 | **0.310** | **0.207** | **0.464** | **0.667** | 1.114 | **0.468** |
| RLinear | base | MSE | 0.248 | 0.723 | 0.231 | 0.551 | 0.309 | 0.488 | 0.220 | 0.085 | 0.992 | 0.816 | 2.637 | 0.664 |
| | | MAE | 0.274 | 0.442 | 0.319 | **0.516** | 0.370 | 0.457 | 0.307 | 0.204 | 0.740 | 0.630 | 1.013 | 0.479 |
| | BSA | MSE | **0.237** | **0.665** | **0.212** | **0.545** | **0.304** | **0.475** | **0.210** | 0.085 | **0.674** | **0.797** | **2.605** | **0.619** |
| | | MAE | **0.267** | **0.404** | **0.301** | 0.517 | 0.370 | **0.453** | **0.303** | **0.203** | **0.607** | **0.622** | **1.005** | **0.459** |
| FreTS | base | MSE | 0.240 | 0.527 | 0.195 | 0.513 | 0.310 | 0.460 | 0.234 | 0.095 | 0.259 | 0.976 | 5.073 | 0.807 |
| | | MAE | **0.286** | 0.322 | **0.285** | 0.506 | 0.393 | 0.450 | **0.321** | 0.232 | 0.354 | 0.735 | 1.603 | 0.499 |
| | BSA | MSE | **0.236** | **0.502** | **0.192** | **0.511** | **0.308** | **0.431** | 0.241 | **0.089** | **0.229** | **0.940** | **4.930** | **0.783** |
| | | MAE | 0.292 | **0.317** | 0.287 | **0.505** | **0.387** | **0.442** | 0.326 | **0.219** | **0.324** | **0.710** | **1.591** | **0.491** |
| TimesNet | base | MSE | 0.260 | **0.620** | 0.203 | 0.576 | 0.386 | 0.508 | 0.247 | 0.094 | **0.232** | 0.880 | **2.719** | 0.611 |
| | | MAE | 0.284 | **0.333** | 0.301 | 0.545 | 0.427 | 0.482 | 0.326 | 0.219 | **0.311** | 0.662 | 0.942 | 0.439 |
| | BSA | MSE | **0.252** | 0.624 | **0.199** | **0.569** | **0.378** | **0.499** | **0.229** | 0.094 | 0.234 | **0.861** | 2.720 | **0.605** |
| | | MAE | **0.279** | 0.335 | **0.298** | **0.540** | **0.418** | **0.476** | **0.314** | 0.219 | 0.314 | **0.655** | **0.938** | **0.435** |
| iTransformer | base | MSE | 0.256 | **0.428** | 0.181 | 0.542 | 0.321 | 0.466 | 0.224 | 0.091 | 0.262 | 0.833 | **2.454** | 0.551 |
| | | MAE | 0.277 | **0.284** | 0.272 | 0.518 | 0.380 | 0.454 | 0.314 | 0.211 | 0.345 | 0.639 | **0.947** | 0.422 |
| | BSA | MSE | **0.236** | 0.428 | **0.176** | **0.537** | **0.317** | 0.466 | **0.219** | **0.090** | **0.198** | **0.786** | 2.540 | **0.545** |
| | | MAE | **0.268** | 0.288 | **0.270** | **0.517** | **0.378** | **0.452** | **0.310** | 0.211 | **0.298** | **0.619** | 0.957 | **0.415** |
| Crossformer | base | MSE | 0.240 | 0.565 | **0.182** | 0.517 | 0.323 | 0.485 | 0.243 | 0.220 | 0.212 | 1.178 | 4.946 | 0.828 |
| | | MAE | 0.292 | 0.292 | **0.276** | **0.520** | 0.402 | 0.470 | 0.342 | 0.345 | 0.299 | **0.817** | 1.518 | 0.507 |
| | BSA | MSE | **0.227** | **0.554** | 0.186 | **0.520** | **0.317** | **0.468** | 0.245 | **0.219** | **0.201** | 1.204 | **4.809** | 0.814 |
| | | MAE | **0.284** | **0.288** | 0.280 | 0.522 | **0.396** | **0.469** | **0.340** | 0.340 | **0.290** | 0.823 | **1.489** | **0.502** |
| PatchTST | base | MSE | 0.255 | **0.467** | 0.198 | **0.535** | 0.317 | 0.473 | 0.220 | 0.087 | 0.361 | 0.839 | 2.016 | 0.524 |
| | | MAE | 0.278 | **0.296** | 0.286 | **0.516** | 0.375 | 0.456 | 0.310 | 0.204 | 0.413 | 0.642 | 0.880 | 0.423 |
| | BSA | MSE | **0.236** | 0.467 | **0.189** | 0.539 | **0.314** | **0.458** | **0.214** | **0.086** | **0.244** | **0.788** | **1.974** | **0.501** |
| | | MAE | **0.268** | 0.297 | **0.285** | 0.520 | **0.374** | **0.451** | **0.306** | **0.202** | **0.341** | **0.621** | **0.866** | **0.412** |

the look-back window is set to 36, and the forecasting lengths are 24, 36, 48, and 60. For the other datasets, the look-back window is set to 96, and the forecasting lengths are 96, 192, 336, and 720. Train, validation, and test split ratio are 0.6, 0.2, 0.2 for the ETT dataset and 0.7, 0.1, 0.2 for the Weather, Traffic, ECL, Exchange, PEMS03, EnergyData, and Illness datasets. The Weather, Traffic, and ECL datasets settings follow the TimesNet paper [48]. Details on datasets are provided in Appendix A.1.

**Baseline models.** As a benchmark, we apply BSA to 7 recent or well-known forecasting models. (1) Linear based models: DLinear [53], RLinear [28], FreTS [52] (2) Convolution based methods: TimesNet [48] (3) Transformer based methods: iTransformer [31], Crossformer [54], PatchTST [36]. For each dataset, the model structure configuration is based on the Time-Series-Library [48]. Details on the base model configurations are provided in Appendix A.2.

**Training details.** We first train the base model for more than 30 epochs (20 epochs for the Traffic dataset) using Adam [22] to ensure that the validation MSE saturates, while also conducting an extensive hyperparameter search. Then, we fine-tune with the BSA module attached to the pre-trained base model. All experiments are based on the average values of three random seeds. We provide further details in Appendix A.3.

## 4.1 Real world datasets

We demonstrate that BSA improves forecasting performance by providing the ability to address long-range dependencies, regardless of the base model architecture. Table 1 presents the averaged values over prediction lengths for both the base and the BSA applied model across 11 time series datasets and 7 TSF models. We apply the BSA module at the beginning of the base model's activation, with the number of BSA channels matching the number of channels in the data (position 1 in Figure

| (a) SA-Matrix heat map | (b) Weather-Temp.(°C) | (c) Weather-SWDR(W/m) | (d) ETTh1-HULL |

Figure 4: This figure illustrates the analysis of the SA-matrix of the DLinear model trained on the 720-step prediction task for the Weather and ETTh1 datasets. Panel (a) shows the heatmap of the SA-matrix, and (b)-(d) show the attention and FFT graphs.

7). BSA effectively improves MSE and MAE across all architectures. The average performance gain, in terms of MSE, ranged from as low as 0.96% to as high as 7.2%, with linear-based models demonstrating relatively high performance. Paired t-tests demonstrate statistically significant (p-value $< 0.05$) improvements in most MSE and all MAE across various models. This result emphasizes the model-agnostic versatility of BSA. Furthermore, BSA exhibited overall performance gain across all datasets, with the enhancement being statistically significant in 82% of cases. BSA's higher gain on ETTm compared to ETTh, both derived from the same signal but with different sampling rates, further indicates BSA's effectiveness in handling long-range dependencies.

To understand how BSA addresses long-range dependencies, we conduct an internal inspection. In Figure 4a, we present the heat map of the trained SA-Matrix of the DLinear (Temperature(°C) channel, Weather data). The positive x-axis corresponds to the $log_{10}$ values of the periods preserved by the low-pass filter, derived with the smoothing factor. Negative values correspond to high-frequency components. The blue graph in Figure 4b represents the SA-Matrix averaged over the feature dimension, illustrating the frequencies to which BSA attends overall. The red graph represents the result of applying the Fast Fourier Transform (FFT) and denoising to the raw signal. The blue graph skewed towards the low-frequency side indicates that the BSA effectively captures the long-range trend of the data. Figure 4c depicts the graph for the SWDR (Short Wave Downward Radiation per unit area, W/m) channel in the same SA-Matrix. While not identical to Figure 4b, it also exhibits strong attention towards the low-frequency pattern. In contrast, Figure 4d, the FFT graph for the HULL (High UseLess Load channel, ETTh1 data), shows that the signal itself lacks long-range trends, resulting in the symmetric SA-Matrix. This result demonstrates that BSA operates as intended, learning the low-frequency components of the signal for future prediction. We provided other graphs and detailed information on the graph plotting method in Appendix B.2.

## 4.2 Synthetic datasets

To further demonstrate that BSA learns long-range dependencies beyond the look-back window, we add sine waves with periods of 100, 300, and 1000 to the natural data while maintaining the mean and standard deviation (Refer to Appendix C for details on synthetic data generation). Figure 5 illustrates the performance improvement of BSA over the iTransformer model on the ETTh1 and



Figure 5: Results of the iTransformer model on synthetic (a) ETTh1 and (b) ETTh2 datasets. The x-axis is the prediction length (96, 192, 336, 720), and the y-axis is the performance improvement (%) compared to the base model. Each color represents the different periods of the sine wave added to the natural data. 0 indicates original data and serves as the baseline.

8

(a) LUFL original     (b) LUFL period 100     (c) LUFL period 300     (d) LUFL period 1000

Figure 6: Attention and FFT graphs on LUFL channel of the original and synthetic ETTh1 data (iTransformer, 720-step prediction). (a) is from the original data and (b)-(d) are from the synthetic data created by adding sine waves with periods of 100, 300, and 1000, respectively. The red arrows on the FFT graphs show the added synthetic signals. Full visualization is provided in Appendix C.1

ETTh2 datasets. The x-axis is the prediction length, and each line represents the period of the added sine wave. Performance improvement is calculated as $100\times$(base MSE - BSA MSE) / base MSE. While the base model with a 96-length look-back window is expected to learn the 100-period trend, BSA outperformed it, especially for 96 and 192-step predictions (green line). The yellow line (period 300) shows nearly a 30% performance improvement across all prediction lengths. While the base model fails to learn the long-range interactions within a period of 300, BSA captures and utilizes the underlying trend. BSA also learns the 1000-period signal (red line) and demonstrates substantial improvements, especially in long prediction-length (336, 720) tasks. These results show that BSA effectively learns long-range patterns beyond the look-back window, essential for future prediction.

Figure 6 is generated from LUFL (Low UseFul Load) channel of ETTh1 data and with added sine waves of periods 100, 300, and 1000. The red arrow shows the added synthetic trend in the FFT graph (red line). As long-range trends are introduced, the attention graph (blue line) shifts from resembling symmetric Gaussian to a low-frequency bias. Longer sine wave periods cause greater shifts, prioritizing long-range information for predictions.

## 4.3  Analysis and ablation studies

| Location | iTransformer | | DLinear | |
|---|---|---|---|---|
| | MSE | MAE | MSE | MAE |
| 1 | 0.2357 | 0.2682 | **0.2196** | **0.2815** |
| 2 | 0.2352 | 0.2681 | - | - |
| 3 | 0.2329 | **0.2654** | - | - |
| 4 | 0.2508 | 0.2752 | 0.2327 | 0.2994 |
| Query-5 | **0.2326** | 0.2671 | - | - |
| Key-6 | 0.2538 | 0.2762 | - | - |
| Value-7 | 0.2415 | 0.2702 | - | - |
| baseline | 0.2556 | 0.2766 | 0.2444 | 0.3084 |

Table 2: Performance analysis on BSA insertion site. Each number corresponds with the insertion site in Figure 7.



Figure 7: Schematic diagram of the BSA insertion site on Transformer.

The BSA module offers high flexibility as it can be applied to arbitrary activations of the base model. In Table 2, we analyzed the performance changes by applying BSA at various locations within the model. Each location corresponds to a number in the Transformer architecture depicted in Figure 7. While we uniformly applied BSA to position 1 for the main result Table 1, the results in Table 2 suggest the potential for further performance enhancement by applying BSA at appropriate locations. Additionally, while there is variability depending on the placement, the performance consistently remains higher compared to the baseline, demonstrating the stability of our method. We provided a full result in Appendix D.1.

BSA shows consistent performance improvement across varying look-back window (input) lengths. Table 3 demonstrates BSA's superiority for look-back window lengths of 48, 96, and 192. Notably, while the baseline model's performance drops significantly with shorter inputs, BSA maintains high performance.

9

Table 3: BSA's performance improvement (%, denoted with *) compared to base model across three input lengths $\{48, 96, 192\}$ (iTransformer, average over 4 prediction lengths and 3 seeds). The full result, including other models, is in Appendix D.2.

| Input-length | Weather | | PEMS03 | |
|---|---|---|---|---|
| | MSE*(%) | MAE*(%) | MSE*(%) | MAE*(%) |
| 48 | 12.08 | 5.21 | 29.45 | 17.35 |
| 96 | 7.78 | 3.03 | 24.55 | 13.66 |
| 192 | 6.87 | 3.34 | 9.63 | 5.18 |

Table 4: Computational cost increase with the BSA (%), averaged over 3 seeds and 4 prediction lengths on the PEMS03 dataset. TN: TimesNet, iTF: iTransformer, CF: Crossformer, PTST: PatchTST. The full result, including other datasets, is in Appendix D.3.

| Increase (%) | TN | iTF | CF | PTST |
|---|---|---|---|---|
| Time | 0.21 | 2.24 | -2.04 | -0.29 |
| Memory | 0.34 | 2.34 | 0.17 | 0.38 |
| Parameter | 0.16 | 4.88 | 1.96 | 2.25 |

We also demonstrate the impact of using BSA on the training time, peak memory, and the number of model parameters. The extra computation required for BSA is constant with data length and linear with a look-back window length (c.f. quadratic for the base model with transformer architecture). Table 4 demonstrates the computational cost of the BSA is quite small, showing less than a 5% increase even with the large PEMS03 dataset.

We conduct an ablation study on the three key components that constitute BSA (Table 5). "BPTT" refers to whether gradients can flow between samples within the mini-batch. Without BPTT, BSA learns similarly to SA. "SFs" denotes whether to use multiple smoothing factors. Lastly, "Learn SF" indicates whether the smoothing factor is treated as learnable. The results indicate that each component significantly contributes to the performance improvement of BSA.

Table 5: Ablation study with 3 components in BSA (Weather data, 720-step prediction, iTransformer).

| BPTT | SFs | Learn SF | MSE | MAE |
|---|---|---|---|---|
| baseline | | | 0.3551 | 0.3473 |
| | | | 0.3480 | 0.3452 |
| ✓ | | | 0.3401 | 0.3452 |
| ✓ | ✓ | | 0.3320 | 0.3395 |
| ✓ | ✓ | ✓ | **0.3263** | **0.3358** |

## 5 Conclusion

Our study addresses the challenges in handling long-range dependencies inherent in time series data by introducing a fast and effective Spectral Attention mechanism. By preserving temporal correlations and enabling the flow of gradients between samples, this mechanism facilitates the model in capturing crucial long-range interactions essential for accurate future predictions. Therefore, our research paves the way for fixed-sized input models to effectively handle long-range dependencies extending far beyond the input window. Through extensive experimentation, we demonstrated that our Spectral Attention mechanism enhances performance across various base architectures, with its ability to grasp long-term dependencies being the key factor behind this improvement. BSA effectively tackles long-term fluctuations, complementing the base model's capacity to manage intricate yet short-term patterns. This integrated model holds promise for improving real-world application performance. For instance, it could boost weather forecast accuracy by simultaneously capturing minute-by-minute weather changes and seasonal variations. Moreover, when predicting deterioration from a patient's real-time data, it can consider medications with lengthy onset times. Our study has limitations: we did not analyze the impact of BSA placement within the base model in detail. Also, BSA's performance gains may be limited when applied to datasets with only high-frequency information within the look-back window. These issues should be addressed in future research.

## 6 Acknowledgement

# References

[1] Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018.

[2] Nesreen K Ahmed, Amir F Atiya, Neamat El Gayar, and Hisham El-Shishiny. An empirical comparison of machine learning models for time series forecasting. *Econometric reviews*, 29(5-6):594–621, 2010.

[3] Peter Bloomfield. *Fourier analysis of time series: an introduction*. John Wiley & Sons, 2004.

[4] Bernhard E Boser, Isabelle M Guyon, and Vladimir N Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 144–152, 1992.

[5] George EP Box and Gwilym M Jenkins. Some recent advances in forecasting and control. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 17(2):91–109, 1968.

[6] Luis M Candanedo, Véronique Feldheim, and Dominique Deramaix. Data driven prediction models of energy use of appliances in a low-energy house. *Energy and buildings*, 140:81–97, 2017.

[7] 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.

[8] Vijay Ekambaram, Arindam Jati, Nam Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. Tsmixer: Lightweight mlp-mixer model for multivariate time series forecasting. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 459–469, 2023.

[9] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.

[10] Agathe Girard, Carl Rasmussen, Joaquin Q Candela, and Roderick Murray-Smith. Gaussian process priors with uncertain inputs application to multiple-step ahead time series forecasting. *Advances in neural information processing systems*, 15, 2002.

[11] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

[12] Pradeep Hewage, Ardhendu Behera, Marcello Trovati, Ella Pereira, Morteza Ghahremani, Francesco Palmieri, and Yonghuai Liu. Temporal convolutional neural (tcn) network for an effective weather forecasting using time-series data from the local weather station. *Soft Computing*, 24:16453–16482, 2020.

[13] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.

[14] Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on document analysis and recognition*, volume 1, pages 278–282. IEEE, 1995.

[15] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[16] Tzu-Ming Harry Hsu, Hang Qi, and Matthew Brown. Measuring the effects of non-identical data distribution for federated visual classification. *arXiv preprint arXiv:1909.06335*, 2019.

[17] MJC Hu and Halbert E Root. An adaptive data processing system for weather forecasting. *Journal of Applied Meteorology and Climatology*, 3(5):513–523, 1964.

[18] Rob J Hyndman and George Athanasopoulos. *Forecasting: principles and practice*. OTexts, 2018.

[19] A. Kahraman, G. Yang, and P. Hou. Wind power forecasting using lstm incorporating fourier transformation based denoising technique. In *20th International Workshop on Large-Scale Integration of Wind Power into Power Systems as well as on Transmission Networks for Offshore Wind Power Plants (WIW 2021)*, volume 2021, pages 94–98, 2021.

[20] Kyoung-jae Kim. Financial time series forecasting using support vector machines. *Neurocomputing*, 55(1-2):307–319, 2003.

[21] 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*, 2021.

[22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[23] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *arXiv preprint arXiv:2001.04451*, 2020.

[24] Diego Krapf, Enzo Marinari, Ralf Metzler, Gleb Oshanin, Xinran Xu, and Alessio Squarcini. Power spectral density of a single brownian trajectory: what one can and cannot learn from it. *New Journal of Physics*, 20(2):023029, 2018.

[25] Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 95–104, 2018.

[26] Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 95–104, 2018.

[27] Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyou Zhou, Wenhu Chen, Yu-Xiang Wang, and Xifeng Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *Advances in neural information processing systems*, 32, 2019.

[28] 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*, 2023.

[29] Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. Scinet: Time series modeling and forecasting with sample convolution and interaction. *Advances in Neural Information Processing Systems*, 35:5816–5828, 2022.

[30] Yeqi Liu, Chuanyang Gong, Ling Yang, and Yingyi Chen. Dstp-rnn: A dual-stage two-phase attention-based recurrent neural network for long-term and multivariate time series prediction. *Expert Systems with Applications*, 143:113082, 2020.

[31] Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. *arXiv preprint arXiv:2310.06625*, 2023.

[32] Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Rethinking the stationarity in time series forecasting. *NeurIPS*, 2022.

[33] MickOhrberg. Brownian noise — Wikipedia, the free encyclopedia, 2024. [Online; accessed 9-April-2024].

[34] Fernando Moreno-Pino, Pablo M Olmos, and Antonio Artés-Rodríguez. Deep autoregressive models with spectral attention. *Pattern Recognition*, 133:109014, 2023.

[35] Brian K Nelson. Time series analysis using autoregressive integrated moving average (arima) models. *Academic emergency medicine*, 5(7):739–744, 1998.

[36] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*, 2022.

[37] Vigneashwara Pandiyan, Rita Drissi-Daoudi, Sergey Shevchik, Giulio Masinelli, Roland Logé, and Kilian Wasmer. Analysis of time, frequency and time-frequency domain features from acoustic emissions during laser powder-bed fusion process. *Procedia CIRP*, 94:392–397, 2020.

[38] Adam Paszke, S. Gross, Francisco Massa, A. Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Z. Lin, N. Gimelshein, L. Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 2019.

[39] Hadi Rezaei, Hamidreza Faaljou, and Gholamreza Mansourfar. Stock price prediction using deep learning and frequency decomposition. *Expert Systems with Applications*, 169:114332, 2021.

[40] Frank Rosenblatt. *The perceptron, a perceiving and recognizing automaton Project Para*. Cornell Aeronautical Laboratory, 1957.

[41] Alaa Sagheer and Mostafa Kotb. Time series forecasting of petroleum production using deep lstm recurrent networks. *Neurocomputing*, 323:203–213, 2019.

[42] Donghwan Song, Adrian Matias Chung Baek, and Namhun Kim. Forecasting stock market indices using padding-based fourier transform denoising and time series deep learning models. *IEEE Access*, 9:83786–83796, 2021.

[43] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27, 2014.

[44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[45] Renzhuo Wan, Shuping Mei, Jun Wang, Min Liu, and Fan Yang. Multivariate temporal convolutional network: A deep neural networks approach for multivariate time series forecasting. *Electronics*, 8(8):876, 2019.

[46] Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun. Transformers in time series: A survey. *arXiv preprint arXiv:2202.07125*, 2022.

[47] Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990.

[48] Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. *ICLR*, 2023.

[49] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in neural information processing systems*, 34:22419–22430, 2021.

[50] Zhijian Xu, Ailing Zeng, and Qiang Xu. FITS: Modeling time series with $10k$ parameters. In *The Twelfth International Conference on Learning Representations*, 2024.

[51] Fuhao Yang, Xin Li, Min Wang, Hongyu Zang, Wei Pang, and Mingzhong Wang. Waveform: Graph enhanced wavelet learning for long sequence forecasting of multivariate time series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 10754–10761, 2023.

[52] 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, 2024.

[53] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 11121–11128, 2023.

[54] Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In *The eleventh international conference on learning representations*, 2022.

[55] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115, 2021.

[56] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International conference on machine learning*, pages 27268–27286. PMLR, 2022.

# A  Details on Datasets, Models, and Training

## A.1  Details on datasets

**Dataset Information.** We conducted experiments on 11 real-world datasets to assess the performance of baseline models and the proposed BSA method. The Weather dataset [49] includes 21 meteorological factors acquired every 10 minutes in 2020 from the Weather Station of the Max Planck Institute for Biogeochemistry. The Traffic dataset [49] records hourly road occupancy rates from 862 sensors on San Francisco Bay Area freeways, covering the period from January 2015 to December 2016. The ECL dataset [49] captures the hourly electricity consumption of 321 clients. The ETT dataset [55] contains seven factors related to electricity transformers, spanning from July 2016 to July 2018. It is divided into four sub-datasets: ETTh1 and ETTh2 are recorded hourly, while ETTm1 and ETTm2 are collected every 15 minutes. The Exchange dataset [49] comprises panel data of daily exchange rates from eight countries, ranging from 1990 to 2016. Illness dataset [49] contains weekly data on influenza-like illness (ILI) cases recorded by the U.S. Centers for Disease Control and Prevention (CDC) from 2002 to 2021. The dataset tracks the proportion of ILI patients relative to the overall number of patients seen during that period. PEMS03 dataset [29] is a sub-dataset of the PEMS dataset, which includes public traffic network data from California, recorded at 5-minute intervals. The EnergyData dataset [6] comprises hourly end-use measurements gathered from 454 residential properties and 140 commercial establishments located in the Pacific Northwest. All these public datasets were downloaded from the referenced sources in March 2024.

**Dataset split.** We adhere to the data processing protocol and train-validation-test split used in Times-Net [48], where the training, validation, and test datasets are sequentially separated in chronological order. Our paper's data split ratios for train, validation, and test set are as follows: (0.7, 0.1, 0.2) for Weather, Traffic, ECL, Exchange, Ilnness, PEMS03, and EnergyData datasets and (0.6, 0.2, 0.2) for ETT. The details of datasets are provided in Table 6.

**Forecasting setting.** Following the approach in iTransformer [31], the look-back window length is set to {96}, while the forecast lengths are {96, 192, 336, 720} for the Weather, Traffic, ECL, ETT, PEMS03, and EnergyData datasets. Forecasting lengths is only {96} for the Exchange dataset since the dataset is too short, which causes biased best model selection with the validation set. For Illness dataset, the look-back window length is {36} and the forecasting lengths are {24, 36, 48, 60}.

Table 6: Summary of Datasets. *Channel* denotes the number of time series variables (channels) for each dataset. *Pre-Train* is the number of training epochs for baseline model saturation. *Finetune* is the number of fine-tuning epochs for our method. *Data Split* means the number of time steps in (Train, Validation, Test) data split respectively. *Sampling Rate* denotes how often the data samples are collected. *Type* is to show the domain in which the data is acquired. We have indicated cases where the data split count matches exactly with that used in TimesNet [48] by marking them with an asterisk (*). For the ETT and the Exchange dataset, the length of the downloaded data differed from the data length reported in TimesNet.

| Dataset | Channel | Pre-Train | Finetune | Data Split | Sampling Rate | Type |
|---|---|---|---|---|---|---|
| Weather* | 21 | 40 | 20 | (36792, 5271, 10540) | 10min | Weather |
| Traffic* | 862 | 20 | 20 | (12185, 1757, 3509) | Hourly | Transportation |
| ECL* | 321 | 30 | 20 | (18317, 2633, 5261) | Hourly | Electricity |
| ETTh1, ETTh2 | 7 | 30 | 30 | (10357, 3485, 3485) | Hourly | Electricity |
| ETTm1, ETTm2 | 7 | 30 | 20 | (41713, 13937, 13937) | 15min | Electricity |
| Exchange | 8 | 30 | 30 | (5216, 761, 1518) | Daily | Economy |
| Illness | 7 | 30 | 20 | (692, 120, 226) | Weekly | Health |
| PEMS03 | 358 | 30 | 20 | (18250, 2623, 5242) | 5min | Transportation |
| EnergyData | 28 | 30 | 20 | (13719, 1975, 3948) | Hourly | Energy |

## A.2  Details on Model Implementations

**Model configurations.** In this section, we discuss the configurations of the baseline models and specify where the BSA module was integrated into each baseline model. For each dataset, the base model structure configuration was directly **replicated from the Time-Series-Library** [48] scripts

when available. Where configurations were not provided, we adjusted them to align closely with the available examples.

1. **DLinear** [53]: The only hyperparameter for this baseline model is the (moving_average = 25) for the series decomposition module from Autoformer [49]. We set the Individual to True so that there are separate linear models for each number of input variables. The BSA module is implemented at the very beginning of the forward pass right before series decomposition. The BSA module is implemented for each channel, which is the number of input variables for this case.

2. **iTransformer** [31]: The hyperparameters for this baseline model are as follows:
(d_model = 512, d_ff = 512, dropout = 0.1, num_heads = 8, encoder_layers = 3, activation = GELU [11]) for the Weather, ECL, PEMS03, and EnergyData datasets,
(d_model = 512, d_ff = 512, dropout = 0.1, num_heads = 8, encoder_layers = 4, activation = GELU) for the Traffic dataset,
(d_model = 128, d_ff = 128, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = GELU) for the ETT and Exchange datasets.
The BSA module is implemented at the beginning part of the forward pass right after normalization from the Non-stationary Transformer [32] and right before the input embedding. The BSA module is implemented for each channel, which is the number of input variables for this case.

3. **Crossformer** [54]: The hyperparameters for this baseline model are as follows:
(seg_len = 12, win_size = 2, factor = 3, d_model = 32, d_ff = 32, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = RELU [1]) for the Weather and EnergyData dataset,
(seg_len = 12, win_size = 2, factor = 3, d_model = 128, d_ff = 128, dropout = 0.1, num_heads = 2, encoder_layers = 2, activation = RELU) for the Traffic dataset,
(seg_len = 12, win_size = 2, factor = 3, d_model = 256, d_ff = 512, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = RELU) for the ECL and PEMS03 dataset,
(seg_len = 12, win_size = 2, factor = 3, d_model = 512, d_ff = 2048, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = RELU) for the ETTh1 and ETTh2 datsets,
(seg_len = 12, win_size = 2, factor = 1, d_model = 512, d_ff = 2048, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = RELU) for the ETTm1 and ETTm2 datasets,
(seg_len = 12, win_size = 2, factor = 3, d_model = 64, d_ff = 64, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = RELU) for the Exchange dataset.
The BSA module is implemented at the very beginning of the forward pass, right before the input embedding. The BSA module is implemented for each channel, which is the number of input variables for this case.

4. **PatchTST** [36]: The hyperparameters for this baseline model are as follows:
(patch_len = 16, stride = 8, d_model = 512, d_ff = 2048, dropout = 0.1, num_heads = 4, encoder_layers = 2, activation = GELU [11]) for the Weather and EnergyData dataset,
(patch_len = 16, stride = 8, d_model = 512, d_ff = 512, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = GELU) for the Traffic dataset,
(patch_len = 16, stride = 8, d_model = 512, d_ff = 2048, dropout = 0.1, num_heads = 8, encoder_layers = 2, activation = GELU) for the ECL and PEMS03 dataset,
(patch_len = 16, stride = 8, d_model = 512, d_ff = 2048, dropout = 0.1, num_heads = 8, encoder_layers = 1, activation = GELU) for the ETTh1 datset,
(patch_len = 16, stride = 8, d_model = 512, d_ff = 2048, dropout = 0.1, num_heads = 4, encoder_layers = 3, activation = GELU) for the ETTh2, ETTm1, and ETTm2 datsets.
The BSA module is implemented at the beginning part of the forward pass right after normalization from the Non-stationary Transformer [32] and right before the patch embedding. The BSA module is implemented for each channel, which is the number of input variables for this case.

5. **TimesNet** [48]: The hyperparameters for this baseline model are as follows:
(top_k = 5, num_kernels = 6, embed = 'timeF', freq = 'h', d_model = 32, d_ff = 32, dropout = 0.1, encoder_layers = 2) for the Weather, ETTh2, ETTm2, and EnergyData datasets,
(top_k = 5, num_kernels = 6, embed = 'timeF', freq = 'h', d_model = 512, d_ff = 512, dropout = 0.1, encoder_layers = 2) for the Traffic dataset,
(top_k = 5, num_kernels = 6, embed = 'timeF', freq = 'h', d_model = 256, d_ff = 512, dropout = 0.1, encoder_layers = 2) for the ECL and PEMS03 dataset,
(top_k = 5, num_kernels = 6, embed = 'timeF', freq = 'h', d_model = 16, d_ff = 32, dropout = 0.1, encoder_layers = 2) for the ETTh1 and ETTm1 datasets,

(top_k = 5, num_kernels = 6, embed = 'timeF', freq = 'h', d_model = 64, d_ff = 64, dropout = 0.1, encoder_layers = 2) for the Exchange dataset.
The BSA module is implemented at the beginning part of the forward pass right after normalization from the Non-stationary Transformer [32] and right before the input embedding. The BSA module is implemented for each channel, which is the number of input variables for this case.

6. **FreTS** [52]: The hyperparameters for this baseline model are as follows:
(embed_size = 128, hidden_size = 256, sparsity_threshold = 0.01, scale = 0.02) for all datasets. Based on the paper, the channel-independent strategy is selected.
The BSA module is implemented at the very beginning of the forward pass, right before the token embedding. The BSA module is implemented for each channel, which is the number of input variables for this case.

7. **RLinear** [28]: The only hyperparameter for this baseline model is the (dropout = 0.1). We set the Individual to True so that there are separate linear models for each number of input variables. The BSA module is implemented at the beginning part of the forward pass right after normalization from the RevIN [21] and right before the linear layer. The BSA module is implemented for each channel, which is the number of input variables for this case.

### A.3   Details on Training

**Pre-training and Finetuning configurations.** To show how our BSA module performs when added to the baseline models, we saturated the models by greedy hyperparameter search.

The hyperparameter search space for the base model is as follows: the possible learning rate is (0.03, 0.01, 0.003, 0.001, 0.0003), and the weight decay is (0.01, 0.003, 0.001, 0.0003, 0.0001, 0.00003).

The hyperparameter search space for BSA finetuning is as follows: the possible learning rate for the SA-Matrix in the BSA module is (0.08, 0.05, 0.03, 0.01, 0.003, 0.001), learning rate for the rest of the model, i.e. original modules, is (0.01, 0.003, 0.001, 0.0003, 0.0001, 0.00003, 0.00001), learning rate for smoothing factor $\alpha_k$ is (none, 0.03, 0.01, 0.003, 0.001, 0.0001, 0.00001), initialization for smoothing factor $\alpha_k$ is ([0.9, 0.99, 0.999], [0.9, 0.99, 0.999, 0.999], [0.9, 0.95, 0.992, 0.999], [0.8, 0.96, 0.992, 0.9984, 0.99968]).

**Model selection** Hyperparameter search is conducted based on the validation set. While models trained using conventional sample shuffling evenly represent the entire time series distribution from which the dataset is sampled, BSA learns data in chronological order. Consequently, the final model tends to favor the distribution of the later part of the data. This can be seen as a mild version of catastrophic forgetting, commonly occurring in continual learning. To mitigate this effect, we assigned higher weights to the later samples during the validation process. The weights are represented by $0.5 + 0.5 \times \sin\left(\frac{\pi}{2} \times \frac{val\_idx}{val\_len}\right)$, continuously increasing from 0.5 for the first sample to 1.0 for the last sample.

**Optimization** The whole code is implemented in PyTorch [38]. Each experiment was conducted on a single NVIDIA GeForce RTX 3090Ti or NVIDIA A40 or NVIDIA L40 GPU. The default batch size for baseline model saturation is 64, while for our method—which involves fine-tuning after integrating the BSA module—it is 256. If a baseline model is too heavy and results in GPU memory overflow, the batch size is adjusted to fit within the available memory. We used the ADAM [22] optimizer and L2 loss (MSE loss) for the model optimization. The baseline saturation training epoch is set to 40 epochs for the Weather dataset, 20 epochs for the Traffic dataset, and 30 epochs for the rest of the datasets. The finetuning epoch is set to 30 epochs for the ETTh1, ETTh2, and Exchange datasets and 20 epochs for the rest of the datasets.

## B   Real world dataset experiments

### B.1   Full experiment results

The experiments were conducted on a total of 11 public real-world datasets and 7 forecasting models. The average results (MSE, MAE) of the experiments conducted with three random seeds 0, 1, and 2 are shown in Table 7, and the standard deviations are shown in Table 8. The values reported in Table 1 are labeled as Avg in Table 7. For the Exchange dataset, we only report experiments with a prediction

| (a) Pre-Denoising | (b) Post-denoising | (c) Brownian Noise [33] |

Figure 8: The illustration compares the pre-denoising and post-denoising results alongside Brownian noise. For (a) and (b), the y-axis represents amplitude, and the x-axis is on a negative log scale, indicating that moving toward the right corresponds to lower frequencies. Conversely, in (c), the y-axis represents decibels, and the x-axis is on a log scale, where moving toward the left corresponds to lower frequencies.

length of {96}. Using prediction lengths of {192, 336, 720} causes improper best model selection due to an insufficient number of validation samples. The empty results indicate that the training is too heavy, and the results are not yet available. We expect to have the results ready by the rebuttal.

## B.2 SA-Matrix and FFT visualization

To investigate how the SA-Matrix of BSA predominantly attends to specific frequency bands, we employed the heatmaps, the Gaussian kernel density estimate graphs, and the FFT graphs as shown in Figure 9, 10, 11, 12, 15 and 13.

The 2D heatmap depicts the learnable parameters of the SA-Matrix, defined as *sa-matrix*$\in \mathbb{R}^{(2K+1)\times D}$, where $K$ represents the number of smoothing factors and $D$ represents the dimension of the input features. Consequently, the y-axis of the 2D heatmap matches the length of $D$, and the x-axis corresponds to $2K + 1$, symmetrically arranged around zero. On the x-axis, positive values indicate the low-frequency regions, while negative values represent the high-frequency regions.

The Gaussian kernel density estimate graphs intuitively reveal which frequency bands the SA-Matrix predominantly attends to. The $K$ smoothing factors were modified according to Equation 8 and symmetrically arranged around zero, serving as the data points for kernel density estimation (KDE). The weight values from each row of the SA-Matrix were converted to probabilities using the softmax function, and the resulting outputs established a mapping of these weight values to data points necessary for Gaussian KDE. This mapping facilitated the construction of the Gaussian KDE graph. Subsequently, the overall probability density of the SA-Matrix was estimated by calculating the mean across columns. Consequently, the y-axis of the graphs denotes the attention level of the SA-Matrix, whereas the x-axis indicates the scalar indices, encompassing the range of the smoothing factors. To better represent the variation in weight values across frequency bands, we adjusted the bandwidth of the KDE function to 0.4, based on the product of the Gaussian kernel's covariance factor and the standard deviation of the sampled weights in the matrix.

$$\alpha_k^{'} = \log_{10}\left(\frac{1}{1 - \sigma(\alpha_k)}\right) \tag{8}$$

The FFT graph is depicted with a red line. FFT analysis was conducted for each variable to determine if the fine-tuned SA-Matrix aligns with the frequencies exhibited by the dataset. To compare the low-frequency parts, the negative log-scale was used on the x-axis, showing that progression to the right indicates decreasing frequencies. The Gaussian filter with a fixed standard deviation of 5 was utilized to smooth the signal's amplitude.

The denoising phase was executed on the FFT output to analyze the correlation between the frequency tendencies of the dataset and the manipulated weights in the SA-Matrix. Figure 8(c) depicts Brownian noise, with the x-axis denoting frequency and the y-axis in decibels (dB). Both axes are configured

18

Table 7: Full average results of three random seeds for prediction lengths $S \in \{96, 192, 336, 720\}$ with a fixed look-back window $T = \{96\}$. The average values across all prediction lengths are reported in Table 1. The Exchange dataset is only done with prediction lengths of $S \in \{96\}$ due to its short data length. The setting for the Illness dataset is $T = \{36\}$, $S \in \{24, 36, 48, 60\}$.

| | | DLinear | | | | RLinear | | | | FreTS | | | | TimesNet | | | | iTransformer | | | | Crossformer | | | | PatchTST | | | |
| | | base | | BSA | | base | | BSA | | base | | BSA | | base | | BSA | | base | | BSA | | base | | BSA | | base | | BSA | |
| Dataset | H | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Weather | 96 | 0.1630 | 0.2363 | 0.1519 | 0.2203 | 0.1656 | 0.2108 | 0.1594 | 0.2065 | 0.1670 | 0.2235 | 0.1587 | 0.2214 | 0.1673 | 0.2157 | 0.1647 | 0.2135 | 0.1706 | 0.2096 | 0.1598 | 0.2037 | 0.1475 | 0.2151 | 0.1447 | 0.2140 | 0.1742 | 0.2149 | 0.1633 | 0.2096 |
| | 192 | 0.2091 | 0.2832 | 0.1912 | 0.2601 | 0.2119 | 0.2512 | 0.2013 | 0.2444 | 0.2076 | 0.2628 | 0.1988 | 0.2611 | 0.2160 | 0.2565 | 0.2097 | 0.2521 | 0.2203 | 0.2544 | 0.2052 | 0.2466 | 0.2004 | 0.2665 | 0.1946 | 0.2615 | 0.2194 | 0.2548 | 0.2030 | 0.2462 |
| | 336 | 0.2634 | 0.3271 | 0.2375 | 0.2993 | 0.2679 | 0.2911 | 0.2541 | 0.2823 | 0.2559 | 0.3001 | 0.2479 | 0.3005 | 0.2923 | 0.3082 | 0.2879 | 0.3065 | 0.2762 | 0.2952 | 0.2526 | 0.2856 | 0.2595 | 0.3086 | 0.2460 | 0.2996 | 0.2745 | 0.2951 | 0.2540 | 0.2850 |
| | 720 | 0.3421 | 0.3868 | 0.2979 | 0.3463 | 0.3465 | 0.3412 | 0.3315 | 0.3338 | 0.3293 | 0.3567 | 0.3367 | 0.3861 | 0.3628 | 0.3558 | 0.3441 | 0.3447 | 0.3551 | 0.3473 | 0.3252 | 0.3371 | 0.3527 | 0.3783 | 0.3216 | 0.3596 | 0.3532 | 0.3456 | 0.3222 | 0.3323 |
| | Avg | 0.2444 | 0.3084 | 0.2196 | 0.2815 | 0.2480 | 0.2736 | 0.2366 | 0.2668 | 0.2400 | 0.3108 | 0.2355 | 0.2923 | 0.2596 | 0.2840 | 0.2516 | 0.2792 | 0.2556 | 0.2766 | 0.2357 | 0.2682 | 0.2400 | 0.2921 | 0.2267 | 0.2837 | 0.2553 | 0.2776 | 0.2356 | 0.2683 |
| Traffic | 96 | 0.6891 | 0.4286 | 0.6679 | 0.4212 | 0.6889 | 0.4153 | 0.6510 | 0.3922 | 0.5089 | 0.3140 | 0.4790 | 0.3060 | 0.5911 | 0.3200 | 0.5930 | 0.3189 | 0.3955 | 0.2684 | 0.3928 | 0.2725 | 0.5217 | 0.2713 | 0.5143 | 0.2664 | 0.4409 | 0.2841 | 0.4405 | 0.2845 |
| | 192 | 0.6752 | 0.4228 | 0.6517 | 0.4187 | 0.6720 | 0.4158 | 0.6228 | 0.3827 | 0.5088 | 0.3120 | 0.4876 | 0.3079 | 0.6074 | 0.3272 | 0.6107 | 0.3288 | 0.4177 | 0.2790 | 0.4170 | 0.2811 | 0.5537 | 0.2873 | 0.5540 | 0.2887 | 0.4551 | 0.2899 | 0.4544 | 0.2903 |
| | 336 | 0.7190 | 0.4431 | 0.6821 | 0.4323 | 0.7142 | 0.4405 | 0.6523 | 0.4098 | 0.5208 | 0.3205 | 0.4954 | 0.3205 | 0.6238 | 0.3348 | 0.6286 | 0.3378 | 0.4339 | 0.2872 | 0.4331 | 0.2899 | 0.5830 | 0.2923 | 0.5507 | 0.2844 | 0.4700 | 0.2963 | 0.4691 | 0.2970 |
| | 720 | 0.8634 | 0.5108 | 0.7607 | 0.4695 | 0.8152 | 0.4945 | 0.7355 | 0.4691 | 0.5681 | 0.3340 | 0.5461 | 0.3326 | 0.6573 | 0.3504 | 0.6625 | 0.3537 | 0.4635 | 0.3028 | 0.4701 | 0.3103 | 0.6024 | 0.3153 | 0.5968 | 0.3125 | 0.5027 | 0.3144 | 0.5049 | 0.3157 |
| | Avg | 0.7367 | 0.4513 | 0.6906 | 0.4354 | 0.7226 | 0.4415 | 0.6654 | 0.4135 | 0.5267 | 0.3201 | 0.5020 | 0.3167 | 0.6199 | 0.3331 | 0.6237 | 0.3348 | 0.4276 | 0.2843 | 0.4283 | 0.2885 | 0.5652 | 0.2916 | 0.5539 | 0.2880 | 0.4672 | 0.2962 | 0.4672 | 0.2969 |
| ECL | 96 | 0.1955 | 0.2891 | 0.1839 | 0.2815 | 0.1990 | 0.2871 | 0.1833 | 0.2730 | 0.1700 | 0.2591 | 0.1571 | 0.2543 | 0.1688 | 0.2720 | 0.1687 | 0.2719 | 0.1499 | 0.2415 | 0.1427 | 0.2390 | 0.1407 | 0.2373 | 0.1412 | 0.2393 | 0.1713 | 0.2614 | 0.1593 | 0.2589 |
| | 192 | 0.2037 | 0.3031 | 0.1956 | 0.2968 | 0.2068 | 0.2990 | 0.1917 | 0.2824 | 0.1773 | 0.2674 | 0.1700 | 0.2659 | 0.1947 | 0.2928 | 0.1925 | 0.2907 | 0.1662 | 0.2574 | 0.1632 | 0.2573 | 0.1596 | 0.2562 | 0.1635 | 0.2612 | 0.1819 | 0.2713 | 0.1746 | 0.2714 |
| | 336 | 0.2265 | 0.3286 | 0.2184 | 0.3218 | 0.2298 | 0.3224 | 0.2120 | 0.3030 | 0.1929 | 0.2844 | 0.1860 | 0.2832 | 0.2150 | 0.3115 | 0.2013 | 0.3003 | 0.1825 | 0.2765 | 0.1768 | 0.2740 | 0.1959 | 0.2896 | 0.1979 | 0.2926 | 0.1986 | 0.2887 | 0.1914 | 0.2883 |
| | 720 | 0.2805 | 0.3773 | 0.2684 | 0.3656 | 0.2875 | 0.3695 | 0.2629 | 0.3451 | 0.2407 | 0.3275 | 0.2559 | 0.3429 | 0.2328 | 0.3282 | 0.2316 | 0.3272 | 0.2259 | 0.3139 | 0.2212 | 0.3083 | 0.2315 | 0.3211 | 0.2419 | 0.3288 | 0.2415 | 0.3229 | 0.2324 | 0.3216 |
| | Avg | 0.2266 | 0.3245 | 0.2166 | 0.3164 | 0.2308 | 0.3195 | 0.2125 | 0.3009 | 0.1952 | 0.2846 | 0.1923 | 0.2866 | 0.2028 | 0.3011 | 0.1985 | 0.2976 | 0.1811 | 0.2723 | 0.1760 | 0.2697 | 0.1819 | 0.2760 | 0.1861 | 0.2805 | 0.1983 | 0.2861 | 0.1894 | 0.2850 |
| ETTH1 | 96 | 0.4329 | 0.4381 | 0.4286 | 0.4370 | 0.4430 | 0.4440 | 0.4319 | 0.4421 | 0.4175 | 0.4415 | 0.4163 | 0.4407 | 0.4616 | 0.4679 | 0.4524 | 0.4595 | 0.4316 | 0.4453 | 0.4277 | 0.4426 | 0.4230 | 0.4519 | 0.4182 | 0.4497 | 0.4277 | 0.4458 | 0.4288 | 0.4480 |
| | 192 | 0.4864 | 0.4770 | 0.4816 | 0.4767 | 0.5014 | 0.4843 | 0.4987 | 0.4899 | 0.4643 | 0.4758 | 0.4608 | 0.4740 | 0.5284 | 0.5177 | 0.5184 | 0.5095 | 0.4906 | 0.4857 | 0.4809 | 0.4812 | 0.4678 | 0.4872 | 0.4662 | 0.4867 | 0.4832 | 0.4858 | 0.4817 | 0.4867 |
| | 336 | 0.5361 | 0.5165 | 0.5343 | 0.5169 | 0.5556 | 0.5213 | 0.5535 | 0.5219 | 0.5202 | 0.5137 | 0.5183 | 0.5127 | 0.5861 | 0.5568 | 0.5824 | 0.5547 | 0.5451 | 0.5217 | 0.5377 | 0.5212 | 0.5143 | 0.5238 | 0.5168 | 0.5269 | 0.5420 | 0.5223 | 0.5514 | 0.5316 |
| | 720 | 0.6621 | 0.5983 | 0.6615 | 0.5975 | 0.7047 | 0.6130 | 0.6945 | 0.6148 | 0.6493 | 0.5942 | 0.6476 | 0.5937 | 0.7267 | 0.6387 | 0.7240 | 0.6378 | 0.6991 | 0.6174 | 0.6979 | 0.6204 | 0.6645 | 0.6169 | 0.6801 | 0.6231 | 0.6878 | 0.6104 | 0.6934 | 0.6128 |
| | Avg | 0.5294 | 0.5075 | 0.5265 | 0.5070 | 0.5512 | 0.5156 | 0.5446 | 0.5172 | 0.5128 | 0.5063 | 0.5107 | 0.5053 | 0.5757 | 0.5453 | 0.5693 | 0.5404 | 0.5416 | 0.5175 | 0.5360 | 0.5163 | 0.5174 | 0.5199 | 0.5203 | 0.5216 | 0.5352 | 0.5161 | 0.5388 | 0.5198 |
| ETTH2 | 96 | 0.2658 | 0.3509 | 0.2597 | 0.3433 | 0.2291 | 0.3167 | 0.2279 | 0.3154 | 0.2480 | 0.3386 | 0.2427 | 0.3341 | 0.2720 | 0.3516 | 0.2638 | 0.3451 | 0.2397 | 0.3278 | 0.2347 | 0.3236 | 0.2611 | 0.3560 | 0.2607 | 0.3546 | 0.2371 | 0.3327 | 0.2336 | 0.3197 |
| | 192 | 0.3168 | 0.3813 | 0.3126 | 0.3759 | 0.2786 | 0.3514 | 0.2717 | 0.3514 | 0.2884 | 0.3705 | 0.2874 | 0.3661 | 0.3741 | 0.4146 | 0.3825 | 0.4183 | 0.2950 | 0.3641 | 0.2901 | 0.3624 | 0.3023 | 0.3846 | 0.3016 | 0.3831 | 0.2835 | 0.3578 | 0.2821 | 0.3546 |
| | 336 | 0.3862 | 0.4298 | 0.3691 | 0.4157 | 0.3186 | 0.3786 | 0.3192 | 0.3832 | 0.3080 | 0.3906 | 0.3016 | 0.3837 | 0.3906 | 0.4306 | 0.4086 | 0.4329 | 0.3316 | 0.3880 | 0.3271 | 0.3884 | 0.3239 | 0.4034 | 0.3210 | 0.3993 | 0.3339 | 0.3896 | 0.3296 | 0.3878 |
| | 720 | 0.4265 | 0.4604 | 0.4222 | 0.4558 | 0.4083 | 0.4343 | 0.3968 | 0.4302 | 0.3956 | 0.4704 | 0.3983 | 0.4660 | 0.5056 | 0.5095 | 0.4583 | 0.4740 | 0.4160 | 0.4389 | 0.4144 | 0.4391 | 0.4061 | 0.4655 | 0.3835 | 0.4477 | 0.4143 | 0.4326 | 0.4119 | 0.4331 |
| | Avg | 0.3488 | 0.4056 | 0.3409 | 0.3977 | 0.3087 | 0.3702 | 0.3039 | 0.3701 | 0.3098 | 0.3925 | 0.3075 | 0.3875 | 0.3856 | 0.4266 | 0.3783 | 0.4176 | 0.3206 | 0.3797 | 0.3166 | 0.3784 | 0.3234 | 0.4024 | 0.3167 | 0.3962 | 0.3172 | 0.3782 | 0.3143 | 0.3738 |
| ETTm1 | 96 | 0.3833 | 0.3965 | 0.3568 | 0.3966 | 0.3755 | 0.3990 | 0.3704 | 0.3971 | 0.3755 | 0.3990 | 0.3585 | 0.3955 | 0.4288 | 0.4322 | 0.4073 | 0.4194 | 0.3671 | 0.3942 | 0.3820 | 0.3984 | 0.4279 | 0.4284 | 0.4185 | 0.4294 | 0.3865 | 0.4004 | 0.3770 | 0.3963 |
| | 192 | 0.4488 | 0.4255 | 0.4107 | 0.4261 | 0.4283 | 0.4298 | 0.4283 | 0.4300 | 0.4283 | 0.4298 | 0.4074 | 0.4256 | 0.4700 | 0.4617 | 0.4640 | 0.4584 | 0.4271 | 0.4318 | 0.4230 | 0.4285 | 0.4617 | 0.4555 | 0.4320 | 0.4511 | 0.4402 | 0.4338 | 0.4338 | 0.4311 |
| | 336 | 0.5225 | 0.4609 | 0.4666 | 0.4594 | 0.4884 | 0.4652 | 0.5081 | 0.4718 | 0.4884 | 0.4652 | 0.4526 | 0.4549 | 0.5259 | 0.4906 | 0.5214 | 0.4864 | 0.4962 | 0.4696 | 0.4925 | 0.4667 | 0.5010 | 0.4799 | 0.5020 | 0.4861 | 0.4974 | 0.4696 | 0.4836 | 0.4696 |
| | 720 | 0.5975 | 0.5053 | 0.5314 | 0.5056 | 0.5466 | 0.5073 | 0.5916 | 0.5291 | 0.5466 | 0.5073 | 0.5036 | 0.4929 | 0.6062 | 0.5432 | 0.6036 | 0.5413 | 0.5752 | 0.5195 | 0.5675 | 0.5162 | 0.5490 | 0.5176 | 0.5188 | 0.5101 | 0.5684 | 0.5210 | 0.5356 | 0.5278 |
| | Avg | 0.4880 | 0.4470 | 0.4414 | 0.4469 | 0.4597 | 0.4503 | 0.4746 | 0.4570 | 0.4597 | 0.4503 | 0.4305 | 0.4422 | 0.5077 | 0.4819 | 0.4991 | 0.4764 | 0.4664 | 0.4538 | 0.4663 | 0.4524 | 0.4849 | 0.4704 | 0.4678 | 0.4692 | 0.4731 | 0.4562 | 0.4575 | 0.4562 |
| ETTm2 | 96 | 0.1653 | 0.2697 | 0.1530 | 0.2584 | 0.1556 | 0.2575 | 0.1503 | 0.2560 | 0.1564 | 0.2617 | 0.1589 | 0.2628 | 0.1562 | 0.2612 | 0.1531 | 0.2589 | 0.1569 | 0.2636 | 0.1531 | 0.2607 | 0.1867 | 0.3050 | 0.1830 | 0.2993 | 0.1547 | 0.2582 | 0.1516 | 0.2553 |
| | 192 | 0.2076 | 0.3045 | 0.1890 | 0.2886 | 0.1904 | 0.2875 | 0.1831 | 0.2847 | 0.1920 | 0.2934 | 0.1977 | 0.2956 | 0.1961 | 0.2934 | 0.1890 | 0.2901 | 0.1942 | 0.2954 | 0.1911 | 0.2923 | 0.2126 | 0.3219 | 0.2151 | 0.3186 | 0.1906 | 0.2914 | 0.1862 | 0.2873 |
| | 336 | 0.2588 | 0.3451 | 0.2517 | 0.3353 | 0.2309 | 0.3178 | 0.2204 | 0.3131 | 0.2462 | 0.3284 | 0.2552 | 0.3321 | 0.2653 | 0.3402 | 0.2300 | 0.3213 | 0.2343 | 0.3253 | 0.2313 | 0.3216 | 0.2610 | 0.3559 | 0.2714 | 0.3571 | 0.2317 | 0.3213 | 0.2265 | 0.3179 |
| | 720 | 0.3585 | 0.4092 | 0.2885 | 0.3592 | 0.3048 | 0.3641 | 0.2865 | 0.3565 | 0.3402 | 0.4018 | 0.3509 | 0.4133 | 0.3701 | 0.4079 | 0.3035 | 0.3692 | 0.3098 | 0.3735 | 0.3039 | 0.3692 | 0.3103 | 0.3867 | 0.3113 | 0.3849 | 0.3047 | 0.3696 | 0.2924 | 0.3625 |
| | Avg | 0.2475 | 0.3321 | 0.2205 | 0.3104 | 0.2204 | 0.3067 | 0.2101 | 0.3026 | 0.2337 | 0.3213 | 0.2407 | 0.3259 | 0.2469 | 0.3257 | 0.2189 | 0.3099 | 0.2238 | 0.3145 | 0.2199 | 0.3110 | 0.2426 | 0.3424 | 0.2452 | 0.3399 | 0.2204 | 0.3101 | 0.2142 | 0.3058 |
| Exchange | 96 | 0.0880 | 0.2110 | 0.0810 | 0.2070 | 0.0852 | 0.2036 | 0.0851 | 0.2035 | 0.0953 | 0.2318 | 0.0893 | 0.2193 | 0.0943 | 0.2187 | 0.0941 | 0.2186 | 0.0905 | 0.2115 | 0.0898 | 0.2106 | 0.2199 | 0.3451 | 0.2188 | 0.3399 | 0.0867 | 0.2035 | 0.0856 | 0.2023 |
| PEMS03 | 96 | 0.4405 | 0.5097 | 0.4121 | 0.4875 | 1.0366 | 0.7815 | 0.6789 | 0.6359 | 0.2316 | 0.3347 | 0.1854 | 0.2905 | 0.2066 | 0.2971 | 0.2113 | 0.2997 | 0.2213 | 0.3226 | 0.1569 | 0.2689 | 0.1353 | 0.2402 | 0.1303 | 0.2363 | 0.3346 | 0.3983 | 0.2099 | 0.3214 |
| | 192 | 0.4658 | 0.5259 | 0.3876 | 0.4677 | 1.0995 | 0.7937 | 0.7033 | 0.6326 | 0.2617 | 0.3574 | 0.2133 | 0.3134 | 0.2396 | 0.3199 | 0.2427 | 0.3234 | 0.2687 | 0.3566 | 0.1865 | 0.2949 | 0.2413 | 0.3368 | 0.2150 | 0.3168 | 0.3818 | 0.4295 | 0.2370 | 0.3393 |
| | 336 | 0.3973 | 0.4669 | 0.3461 | 0.4317 | 0.8227 | 0.6434 | 0.5861 | 0.5408 | 0.2528 | 0.3456 | 0.2380 | 0.3286 | 0.2169 | 0.2982 | 0.2199 | 0.3022 | 0.2479 | 0.3303 | 0.1979 | 0.2957 | 0.1987 | 0.2849 | 0.1870 | 0.2748 | 0.3303 | 0.3884 | 0.2390 | 0.3327 |
| | 720 | 0.4420 | 0.5048 | 0.3922 | 0.4703 | 1.0101 | 0.7432 | 0.7289 | 0.6184 | 0.2915 | 0.3765 | 0.2789 | 0.3619 | 0.2644 | 0.3290 | 0.2638 | 0.3300 | 0.3093 | 0.3716 | 0.2488 | 0.3330 | 0.2714 | 0.3321 | 0.2730 | 0.3337 | 0.3979 | 0.4347 | 0.2901 | 0.3724 |
| | Avg | 0.4364 | 0.5018 | 0.3845 | 0.4643 | 0.9922 | 0.7405 | 0.6743 | 0.6069 | 0.2594 | 0.3535 | 0.2289 | 0.3236 | 0.2319 | 0.3110 | 0.2344 | 0.3138 | 0.2618 | 0.3453 | 0.1975 | 0.2981 | 0.2117 | 0.2985 | 0.2013 | 0.2904 | 0.3612 | 0.4127 | 0.2440 | 0.3415 |
| EnergyData | 96 | 0.4860 | 0.4887 | 0.4746 | 0.4810 | 0.4779 | 0.4663 | 0.4639 | 0.4607 | 0.5467 | 0.5252 | 0.5305 | 0.5224 | 0.5018 | 0.4891 | 0.4915 | 0.4838 | 0.4898 | 0.4734 | 0.4515 | 0.4523 | 0.7204 | 0.6206 | 0.7145 | 0.6102 | 0.4962 | 0.4770 | 0.4558 | 0.4555 |
| | 192 | 0.6560 | 0.5950 | 0.6405 | 0.5857 | 0.6306 | 0.5542 | 0.6137 | 0.5470 | 0.7655 | 0.6553 | 0.7221 | 0.6237 | 0.6794 | 0.5867 | 0.6766 | 0.5848 | 0.6340 | 0.5611 | 0.5992 | 0.5432 | 0.9027 | 0.7175 | 0.9405 | 0.7301 | 0.6467 | 0.5651 | 0.6069 | 0.5464 |
| | 336 | 0.9041 | 0.7142 | 0.8885 | 0.7073 | 0.8493 | 0.6547 | 0.8289 | 0.6448 | 1.0241 | 0.7765 | 0.9195 | 0.7216 | 0.9190 | 0.6897 | 0.8985 | 0.6811 | 0.8653 | 0.6645 | 0.8216 | 0.6444 | 1.2777 | 0.8696 | 1.2763 | 0.8683 | 0.8710 | 0.6661 | 0.8268 | 0.6470 |
| | 720 | 1.4487 | 0.9339 | 1.3465 | 0.8940 | 1.3075 | 0.8455 | 1.2808 | 0.8358 | 1.5692 | 0.9824 | 1.5876 | 0.9734 | 1.4197 | 0.8840 | 1.3784 | 0.8684 | 1.3438 | 0.8587 | 1.2712 | 0.8358 | 1.8131 | 1.0623 | 1.8847 | 1.0815 | 1.3411 | 0.8583 | 1.2609 | 0.8350 |
| | Avg | 0.8737 | 0.6829 | 0.8375 | 0.6670 | 0.8163 | 0.6302 | 0.7968 | 0.6221 | 0.9764 | 0.7349 | 0.9399 | 0.7103 | 0.8800 | 0.6624 | 0.8612 | 0.6545 | 0.8332 | 0.6395 | 0.7859 | 0.6189 | 1.1785 | 0.8175 | 1.2040 | 0.8225 | 0.8388 | 0.6416 | 0.7876 | 0.6210 |
| Illness | 24 | 2.7745 | 1.0734 | 2.7651 | 1.1114 | 2.7842 | 1.0319 | 2.7533 | 1.0241 | 1.6104 | 0.8490 | 1.5862 | 0.8380 | 3.5866 | 0.9677 | 3.5957 | 0.9673 | 3.3026 | 1.0037 | 3.2067 | 0.9905 | 4.7314 | 1.4804 | 4.3844 | 1.4091 | 1.9154 | 0.8445 | 1.9200 | 0.8442 |
| | 36 | 2.7486 | 1.0943 | 2.7382 | 1.0941 | 2.5405 | 0.9835 | 2.5383 | 0.9824 | 1.6043 | 0.8560 | 1.6153 | 0.8560 | 3.0292 | 0.9940 | 3.0613 | 0.9901 | 2.0613 | 0.9119 | 2.0367 | 0.9119 | 5.0147 | 1.5197 | 4.9364 | 1.4996 | 2.2098 | 0.9099 | 2.1997 | 0.8882 |
| | 48 | 2.9586 | 1.1552 | 2.8371 | 1.1311 | 2.6409 | 1.0239 | 2.6358 | 1.0220 | 1.5957 | 0.8470 | 1.5674 | 0.8320 | 2.1801 | 0.9824 | 2.1578 | 0.8862 | 2.2795 | 0.9328 | 2.5946 | 0.9328 | 5.0649 | 1.5484 | 4.9638 | 1.5227 | 2.0023 | 0.8758 | 1.9519 | 0.8660 |
| | 60 | 2.9136 | 1.1205 | 2.8481 | 1.1174 | 2.5818 | 1.0139 | 2.4921 | 0.9914 | 1.6015 | 0.8510 | 1.5939 | 0.8450 | 2.0794 | 0.9420 | 2.0628 | 0.9087 | 2.1722 | 0.9392 | 2.5402 | 0.9392 | 4.9448 | 1.5217 | 4.9531 | 1.5234 | 1.9354 | 0.8915 | 1.8238 | 0.8639 |
| | Avg | 2.8488 | 1.1109 | 2.7971 | 1.1135 | 2.6369 | 1.0133 | 2.6049 | 1.0050 | 1.6030 | 0.8508 | 1.5907 | 0.8428 | 2.7188 | 0.9543 | 2.7196 | 0.9381 | 2.4539 | 0.9469 | 2.5402 | 0.9571 | 4.9956 | 1.5175 | 4.8094 | 1.4887 | 2.0157 | 0.8804 | 1.9738 | 0.8656 |

19

Table 8: Full standard deviation results of three random seeds with prediction lengths $S \in \{96, 192, 336, 720\}$ and a fixed lookback length $T = \{96\}$. The Exchange dataset is only done with prediction lengths of $S \in \{96\}$ due to its short data length. The setting for the Illness dataset is $T = \{36\}$, $S \in \{24, 36, 48, 60\}$.

| Dataset | Len | Dlinear base MSE | Dlinear base MAE | Dlinear BSA MSE | Dlinear BSA MAE | RLinear base MSE | RLinear base MAE | RLinear BSA MSE | RLinear BSA MAE | FreTS base MSE | FreTS base MAE | FreTS BSA MSE | FreTS BSA MAE | TimesNet base MSE | TimesNet base MAE | TimesNet BSA MSE | TimesNet BSA MAE | iTransformer base MSE | iTransformer base MAE | iTransformer BSA MSE | iTransformer BSA MAE | Crossformer base MSE | Crossformer base MAE | Crossformer BSA MSE | Crossformer BSA MAE | PatchTST base MSE | PatchTST base MAE | PatchTST BSA MSE | PatchTST BSA MAE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Weather | 96 | 0.0002 | 0.0005 | 0.0000 | 0.0000 | 0.0002 | 0.0003 | 0.0000 | 0.0000 | 0.0001 | 0.0002 | 0.0003 | 0.0002 | 0.0008 | 0.0002 | 0.0015 | 0.0011 | 0.0015 | 0.0017 | 0.0008 | 0.0011 | 0.0024 | 0.0012 | 0.0029 | 0.0027 | 0.0004 | 0.0005 | 0.0012 | 0.0013 |
| Weather | 192 | 0.0003 | 0.0010 | 0.0001 | 0.0002 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0008 | 0.0002 | 0.0002 | 0.0002 | 0.0019 | 0.0014 | 0.0040 | 0.0019 | 0.0014 | 0.0013 | 0.0019 | 0.0015 | 0.0014 | 0.0022 | 0.0018 | 0.0019 | 0.0000 | 0.0007 | 0.0007 | 0.0003 |
| Weather | 336 | 0.0008 | 0.0014 | 0.0002 | 0.0002 | 0.0000 | 0.0001 | 0.0000 | 0.0000 | 0.0006 | 0.0009 | 0.0002 | 0.0004 | 0.0010 | 0.0002 | 0.0034 | 0.0026 | 0.0003 | 0.0005 | 0.0008 | 0.0010 | 0.0026 | 0.0028 | 0.0019 | 0.0016 | 0.0001 | 0.0004 | 0.0007 | 0.0001 |
| Weather | 720 | 0.0064 | 0.0082 | 0.0003 | 0.0005 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0043 | 0.0075 | 0.0058 | 0.0059 | 0.0072 | 0.0059 | 0.0032 | 0.0033 | 0.0007 | 0.0004 | 0.0013 | 0.0012 | 0.0027 | 0.0023 | 0.0018 | 0.0019 | 0.0005 | 0.0004 | 0.0006 | 0.0006 |
| Weather | Avg | 0.0019 | 0.0028 | 0.0002 | 0.0002 | 0.0001 | 0.0001 | 0.0001 | 0.0001 | 0.0013 | 0.0022 | 0.0017 | 0.0017 | 0.0027 | 0.0019 | 0.0032 | 0.0022 | 0.0010 | 0.0010 | 0.0012 | 0.0012 | 0.0023 | 0.0021 | 0.0018 | 0.0006 | 0.0003 | 0.0003 | 0.0006 | 0.0006 |
| Traffic | 96 | 0.0000 | 0.0000 | 0.0001 | 0.0000 | 0.0002 | 0.0001 | 0.0001 | 0.0000 | 0.0008 | 0.0002 | 0.0015 | 0.0002 | 0.0060 | 0.0082 | 0.0053 | 0.0023 | 0.0003 | 0.0003 | 0.0012 | 0.0005 | 0.0208 | 0.0088 | 0.0203 | 0.0099 | 0.0021 | 0.0003 | 0.0008 | 0.0001 |
| Traffic | 192 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0001 | 0.0002 | 0.0000 | 0.0000 | 0.0001 | 0.0001 | 0.0004 | 0.0005 | 0.0082 | 0.0059 | 0.0081 | 0.0049 | 0.0011 | 0.0007 | 0.0009 | 0.0011 | 0.0020 | 0.0004 | 0.0018 | 0.0007 | 0.0009 | 0.0009 | 0.0009 | 0.0008 |
| Traffic | 336 | 0.0001 | 0.0000 | 0.0000 | 0.0000 | 0.0001 | 0.0000 | 0.0000 | 0.0000 | 0.0001 | 0.0001 | 0.0008 | 0.0002 | 0.0120 | 0.0033 | 0.0042 | 0.0030 | 0.0009 | 0.0009 | 0.0005 | 0.0013 | 0.0227 | 0.0058 | 0.0027 | 0.0007 | 0.0025 | 0.0007 | 0.0014 | 0.0010 |
| Traffic | 720 | 0.0001 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0005 | 0.0002 | 0.0009 | 0.0002 | 0.0032 | 0.0030 | 0.0022 | 0.0015 | 0.0006 | 0.0003 | 0.0023 | 0.0005 | 0.0005 | 0.0002 | 0.0071 | 0.0029 | 0.0012 | 0.0004 | 0.0014 | 0.0006 |
| Traffic | Avg | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0001 | 0.0001 | 0.0001 | 0.0001 | 0.0005 | 0.0002 | 0.0009 | 0.0002 | 0.0073 | 0.0034 | 0.0050 | 0.0029 | 0.0007 | 0.0006 | 0.0012 | 0.0008 | 0.0115 | 0.0038 | 0.0071 | 0.0029 | 0.0017 | 0.0004 | 0.0011 | 0.0006 |
| ECL | 96 | 0.0001 | 0.0000 | 0.0000 | 0.0002 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0005 | 0.0004 | 0.0002 | 0.0004 | 0.0021 | 0.0013 | 0.0029 | 0.0005 | 0.0009 | 0.0009 | 0.0006 | 0.0005 | 0.0051 | 0.0047 | 0.0052 | 0.0051 | 0.0005 | 0.0011 | 0.0005 | 0.0006 |
| ECL | 192 | 0.0003 | 0.0002 | 0.0000 | 0.0001 | 0.0000 | 0.0000 | 0.0000 | 0.0002 | 0.0004 | 0.0004 | 0.0001 | 0.0004 | 0.0119 | 0.0084 | 0.0076 | 0.0013 | 0.0008 | 0.0005 | 0.0013 | 0.0011 | 0.0057 | 0.0054 | 0.0084 | 0.0082 | 0.0014 | 0.0011 | 0.0004 | 0.0001 |
| ECL | 336 | 0.0001 | 0.0001 | 0.0000 | 0.0001 | 0.0000 | 0.0002 | 0.0002 | 0.0001 | 0.0007 | 0.0004 | 0.0001 | 0.0001 | 0.0167 | 0.0118 | 0.0124 | 0.0029 | 0.0028 | 0.0039 | 0.0020 | 0.0029 | 0.0013 | 0.0015 | 0.0038 | 0.0058 | 0.0007 | 0.0013 | 0.0007 | 0.0004 |
| ECL | 720 | 0.0001 | 0.0001 | 0.0000 | 0.0000 | 0.0003 | 0.0001 | 0.0000 | 0.0000 | 0.0009 | 0.0008 | 0.0034 | 0.0025 | 0.0070 | 0.0060 | 0.0057 | 0.0002 | 0.0038 | 0.0050 | 0.0001 | 0.0002 | 0.0254 | 0.0169 | 0.0204 | 0.0148 | 0.0009 | 0.0014 | 0.0005 | 0.0007 |
| ECL | Avg | 0.0001 | 0.0001 | 0.0000 | 0.0000 | 0.0004 | 0.0001 | 0.0000 | 0.0000 | 0.0006 | 0.0005 | 0.0010 | 0.0008 | 0.0094 | 0.0069 | 0.0071 | 0.0012 | 0.0021 | 0.0026 | 0.0010 | 0.0012 | 0.0094 | 0.0071 | 0.0095 | 0.0085 | 0.0008 | 0.0014 | 0.0005 | 0.0005 |
| ETTH1 | 96 | 0.0002 | 0.0002 | 0.0001 | 0.0002 | 0.0001 | 0.0002 | 0.0001 | 0.0001 | 0.0053 | 0.0019 | 0.0048 | 0.0015 | 0.0050 | 0.0039 | 0.0048 | 0.0015 | 0.0006 | 0.0014 | 0.0002 | 0.0007 | 0.0005 | 0.0006 | 0.0006 | 0.0003 | 0.0004 | 0.0006 | 0.0010 | 0.0012 |
| ETTH1 | 192 | 0.0006 | 0.0004 | 0.0002 | 0.0001 | 0.0000 | 0.0001 | 0.0004 | 0.0001 | 0.0015 | 0.0004 | 0.0008 | 0.0005 | 0.0063 | 0.0067 | 0.0052 | 0.0011 | 0.0013 | 0.0006 | 0.0003 | 0.0006 | 0.0012 | 0.0009 | 0.0013 | 0.0011 | 0.0047 | 0.0033 | 0.0024 | 0.0010 |
| ETTH1 | 336 | 0.0011 | 0.0006 | 0.0003 | 0.0002 | 0.0001 | 0.0001 | 0.0001 | 0.0000 | 0.0029 | 0.0005 | 0.0024 | 0.0007 | 0.0083 | 0.0058 | 0.0069 | 0.0013 | 0.0024 | 0.0009 | 0.0004 | 0.0004 | 0.0023 | 0.0015 | 0.0008 | 0.0009 | 0.0032 | 0.0013 | 0.0035 | 0.0021 |
| ETTH1 | 720 | 0.0007 | 0.0009 | 0.0004 | 0.0006 | 0.0001 | 0.0001 | 0.0000 | 0.0002 | 0.0026 | 0.0023 | 0.0022 | 0.0019 | 0.0120 | 0.0038 | 0.0111 | 0.0071 | 0.0054 | 0.0046 | 0.0038 | 0.0024 | 0.0125 | 0.0071 | 0.0012 | 0.0007 | 0.0009 | 0.0013 | 0.0039 | 0.0021 |
| ETTH1 | Avg | 0.0007 | 0.0005 | 0.0002 | 0.0003 | 0.0001 | 0.0001 | 0.0002 | 0.0001 | 0.0031 | 0.0013 | 0.0025 | 0.0011 | 0.0079 | 0.0051 | 0.0070 | 0.0025 | 0.0024 | 0.0019 | 0.0012 | 0.0010 | 0.0041 | 0.0025 | 0.0010 | 0.0007 | 0.0023 | 0.0013 | 0.0027 | 0.0016 |
| ETTH2 | 96 | 0.0116 | 0.0108 | 0.0095 | 0.0081 | 0.0001 | 0.0001 | 0.0001 | 0.0001 | 0.0097 | 0.0095 | 0.0051 | 0.0053 | 0.0064 | 0.0051 | 0.0041 | 0.0036 | 0.0011 | 0.0009 | 0.0003 | 0.0006 | 0.0028 | 0.0014 | 0.0032 | 0.0013 | 0.0060 | 0.0027 | 0.0038 | 0.0020 |
| ETTH2 | 192 | 0.0152 | 0.0098 | 0.0145 | 0.0088 | 0.0003 | 0.0002 | 0.0003 | 0.0001 | 0.0009 | 0.0070 | 0.0014 | 0.0049 | 0.0268 | 0.0132 | 0.0268 | 0.0154 | 0.0022 | 0.0012 | 0.0007 | 0.0012 | 0.0012 | 0.0007 | 0.0009 | 0.0056 | 0.0014 | 0.0015 | 0.0005 | 0.0005 |
| ETTH2 | 336 | 0.0191 | 0.0065 | 0.0160 | 0.0062 | 0.0000 | 0.0000 | 0.0003 | 0.0003 | 0.0018 | 0.0047 | 0.0057 | 0.0053 | 0.0152 | 0.0098 | 0.0189 | 0.0107 | 0.0006 | 0.0008 | 0.0020 | 0.0008 | 0.0048 | 0.0097 | 0.0038 | 0.0116 | 0.0021 | 0.0006 | 0.0011 | 0.0021 |
| ETTH2 | 720 | 0.0231 | 0.0141 | 0.0222 | 0.0137 | 0.0003 | 0.0003 | 0.0017 | 0.0012 | 0.0242 | 0.0151 | 0.0051 | 0.0008 | 0.0096 | 0.0048 | 0.0338 | 0.0163 | 0.0005 | 0.0010 | 0.0010 | 0.0010 | 0.0367 | 0.0340 | 0.0130 | 0.0048 | 0.0019 | 0.0010 | 0.0007 | 0.0003 |
| ETTH2 | Avg | 0.0173 | 0.0103 | 0.0156 | 0.0092 | 0.0002 | 0.0003 | 0.0000 | 0.0002 | 0.0092 | 0.0090 | 0.0044 | 0.0041 | 0.0145 | 0.0082 | 0.0209 | 0.0115 | 0.0011 | 0.0010 | 0.0010 | 0.0009 | 0.0114 | 0.0115 | 0.0052 | 0.0058 | 0.0028 | 0.0015 | 0.0015 | 0.0012 |
| ETTm1 | 96 | 0.0003 | 0.0001 | 0.0000 | 0.0000 | 0.0003 | 0.0001 | 0.0006 | 0.0006 | 0.0015 | 0.0005 | 0.0007 | 0.0004 | 0.0055 | 0.0018 | 0.0067 | 0.0025 | 0.0006 | 0.0006 | 0.0027 | 0.0015 | 0.0014 | 0.0014 | 0.0047 | 0.0031 | 0.0069 | 0.0032 | 0.0010 | 0.0007 |
| ETTm1 | 192 | 0.0002 | 0.0001 | 0.0000 | 0.0000 | 0.0003 | 0.0002 | 0.0003 | 0.0003 | 0.0008 | 0.0003 | 0.0009 | 0.0002 | 0.0087 | 0.0062 | 0.0075 | 0.0051 | 0.0063 | 0.0043 | 0.0050 | 0.0030 | 0.0079 | 0.0041 | 0.0047 | 0.0016 | 0.0037 | 0.0017 | 0.0025 | 0.0010 |
| ETTm1 | 336 | 0.0006 | 0.0002 | 0.0001 | 0.0001 | 0.0000 | 0.0000 | 0.0002 | 0.0002 | 0.0005 | 0.0007 | 0.0012 | 0.0009 | 0.0052 | 0.0049 | 0.0036 | 0.0027 | 0.0024 | 0.0008 | 0.0029 | 0.0020 | 0.0186 | 0.0125 | 0.0345 | 0.0211 | 0.0023 | 0.0021 | 0.0026 | 0.0028 |
| ETTm1 | 720 | 0.0011 | 0.0006 | 0.0001 | 0.0000 | 0.0006 | 0.0003 | 0.0001 | 0.0000 | 0.0044 | 0.0007 | 0.0011 | 0.0009 | 0.0058 | 0.0029 | 0.0062 | 0.0031 | 0.0032 | 0.0018 | 0.0011 | 0.0009 | 0.0055 | 0.0024 | 0.0035 | 0.0026 | 0.0041 | 0.0034 | 0.0008 | 0.0011 |
| ETTm1 | Avg | 0.0006 | 0.0003 | 0.0001 | 0.0000 | 0.0003 | 0.0002 | 0.0003 | 0.0002 | 0.0018 | 0.0007 | 0.0010 | 0.0006 | 0.0063 | 0.0039 | 0.0060 | 0.0033 | 0.0031 | 0.0019 | 0.0029 | 0.0019 | 0.0083 | 0.0051 | 0.0119 | 0.0071 | 0.0043 | 0.0026 | 0.0017 | 0.0014 |
| ETTm2 | 96 | 0.0038 | 0.0040 | 0.0005 | 0.0004 | 0.0001 | 0.0002 | 0.0001 | 0.0004 | 0.0010 | 0.0009 | 0.0004 | 0.0003 | 0.0015 | 0.0017 | 0.0004 | 0.0005 | 0.0009 | 0.0010 | 0.0006 | 0.0006 | 0.0048 | 0.0083 | 0.0023 | 0.0038 | 0.0004 | 0.0004 | 0.0001 | 0.0001 |
| ETTm2 | 192 | 0.0069 | 0.0032 | 0.0015 | 0.0009 | 0.0000 | 0.0001 | 0.0000 | 0.0009 | 0.0002 | 0.0011 | 0.0013 | 0.0004 | 0.0050 | 0.0025 | 0.0052 | 0.0030 | 0.0005 | 0.0002 | 0.0007 | 0.0005 | 0.0061 | 0.0080 | 0.0008 | 0.0077 | 0.0010 | 0.0006 | 0.0011 | 0.0003 |
| ETTm2 | 336 | 0.0062 | 0.0057 | 0.0041 | 0.0022 | 0.0002 | 0.0002 | 0.0001 | 0.0000 | 0.0052 | 0.0037 | 0.0046 | 0.0030 | 0.0212 | 0.0114 | 0.0181 | 0.0098 | 0.0002 | 0.0008 | 0.0013 | 0.0009 | 0.0017 | 0.0140 | 0.0057 | 0.0039 | 0.0005 | 0.0003 | 0.0004 | 0.0003 |
| ETTm2 | 720 | 0.0238 | 0.0160 | 0.0044 | 0.0025 | 0.0005 | 0.0003 | 0.0000 | 0.0000 | 0.0323 | 0.0112 | 0.0172 | 0.0062 | 0.0236 | 0.0154 | 0.0015 | 0.0011 | 0.0010 | 0.0009 | 0.0009 | 0.0008 | 0.0034 | 0.0038 | 0.0056 | 0.0040 | 0.0004 | 0.0009 | 0.0009 | 0.0004 |
| ETTm2 | Avg | 0.0102 | 0.0072 | 0.0026 | 0.0015 | 0.0002 | 0.0002 | 0.0001 | 0.0000 | 0.0097 | 0.0043 | 0.0059 | 0.0025 | 0.0128 | 0.0078 | 0.0063 | 0.0036 | 0.0007 | 0.0007 | 0.0009 | 0.0007 | 0.0040 | 0.0085 | 0.0036 | 0.0048 | 0.0006 | 0.0006 | 0.0006 | 0.0004 |
| Exchange | 96 | 0.0025 | 0.0023 | 0.0002 | 0.0003 | 0.0017 | 0.0012 | 0.0015 | 0.0015 | 0.0043 | 0.0077 | 0.0016 | 0.0031 | 0.0005 | 0.0004 | 0.0004 | 0.0004 | 0.0020 | 0.0018 | 0.0020 | 0.0019 | 0.0172 | 0.0133 | 0.0161 | 0.0123 | 0.0033 | 0.0033 | 0.0026 | 0.0025 |
| PEMS03 | 96 | 0.0003 | 0.0004 | 0.0004 | 0.0005 | 0.0017 | 0.0020 | 0.0002 | 0.0001 | 0.0009 | 0.0011 | 0.0008 | 0.0008 | 0.0054 | 0.0035 | 0.0019 | 0.0006 | 0.0021 | 0.0011 | 0.0009 | 0.0006 | 0.0034 | 0.0015 | 0.0039 | 0.0037 | 0.0012 | 0.0016 | 0.0006 | 0.0007 |
| PEMS03 | 192 | 0.0010 | 0.0006 | 0.0000 | 0.0000 | 0.0011 | 0.0011 | 0.0024 | 0.0006 | 0.0002 | 0.0001 | 0.0022 | 0.0014 | 0.0034 | 0.0026 | 0.0063 | 0.0013 | 0.0035 | 0.0020 | 0.0019 | 0.0013 | 0.0508 | 0.0400 | 0.0235 | 0.0207 | 0.0001 | 0.0002 | 0.0012 | 0.0010 |
| PEMS03 | 336 | 0.0001 | 0.0004 | 0.0001 | 0.0000 | 0.0012 | 0.0002 | 0.0009 | 0.0006 | 0.0004 | 0.0001 | 0.0019 | 0.0012 | 0.0042 | 0.0030 | 0.0014 | 0.0014 | 0.0022 | 0.0011 | 0.0019 | 0.0014 | 0.0061 | 0.0066 | 0.0121 | 0.0104 | 0.0007 | 0.0003 | 0.0007 | 0.0007 |
| PEMS03 | 720 | 0.0003 | 0.0004 | 0.0002 | 0.0001 | 0.0015 | 0.0008 | 0.0028 | 0.0017 | 0.0001 | 0.0004 | 0.0011 | 0.0007 | 0.0004 | 0.0012 | 0.0024 | 0.0017 | 0.0045 | 0.0015 | 0.0029 | 0.0017 | 0.0612 | 0.0425 | 0.0610 | 0.0458 | 0.0010 | 0.0007 | 0.0014 | 0.0008 |
| PEMS03 | Avg | 0.0004 | 0.0004 | 0.0002 | 0.0001 | 0.0014 | 0.0010 | 0.0016 | 0.0001 | 0.0004 | 0.0004 | 0.0015 | 0.0010 | 0.0034 | 0.0026 | 0.0030 | 0.0012 | 0.0031 | 0.0014 | 0.0019 | 0.0012 | 0.0304 | 0.0227 | 0.0251 | 0.0201 | 0.0008 | 0.0007 | 0.0010 | 0.0008 |
| EnergyData | 96 | 0.0030 | 0.0027 | 0.0002 | 0.0002 | 0.0002 | 0.0001 | 0.0002 | 0.0001 | 0.0078 | 0.0065 | 0.0022 | 0.0017 | 0.0026 | 0.0018 | 0.0023 | 0.0024 | 0.0036 | 0.0027 | 0.0039 | 0.0024 | 0.0485 | 0.0262 | 0.0101 | 0.0068 | 0.0006 | 0.0018 | 0.0013 | 0.0005 |
| EnergyData | 192 | 0.0044 | 0.0030 | 0.0006 | 0.0005 | 0.0003 | 0.0003 | 0.0001 | 0.0000 | 0.0530 | 0.0288 | 0.0157 | 0.0081 | 0.0020 | 0.0010 | 0.0017 | 0.0005 | 0.0028 | 0.0028 | 0.0006 | 0.0005 | 0.0665 | 0.0250 | 0.0780 | 0.0276 | 0.0009 | 0.0016 | 0.0017 | 0.0006 |
| EnergyData | 336 | 0.0078 | 0.0041 | 0.0027 | 0.0010 | 0.0006 | 0.0003 | 0.0003 | 0.0000 | 0.0096 | 0.0041 | 0.0096 | 0.0032 | 0.0071 | 0.0028 | 0.0091 | 0.0015 | 0.0006 | 0.0010 | 0.0030 | 0.0015 | 0.0285 | 0.0051 | 0.0272 | 0.0049 | 0.0056 | 0.0010 | 0.0018 | 0.0008 |
| EnergyData | 720 | 0.0147 | 0.0046 | 0.0100 | 0.0042 | 0.0006 | 0.0002 | 0.0004 | 0.0001 | 0.0135 | 0.0072 | 0.0577 | 0.0232 | 0.0407 | 0.0108 | 0.0159 | 0.0026 | 0.0071 | 0.0022 | 0.0084 | 0.0026 | 0.0638 | 0.0184 | 0.0393 | 0.0120 | 0.0027 | 0.0013 | 0.0048 | 0.0008 |
| EnergyData | Avg | 0.0075 | 0.0036 | 0.0034 | 0.0015 | 0.0004 | 0.0002 | 0.0002 | 0.0001 | 0.0210 | 0.0117 | 0.0213 | 0.0091 | 0.0131 | 0.0041 | 0.0073 | 0.0017 | 0.0035 | 0.0022 | 0.0040 | 0.0017 | 0.0518 | 0.0187 | 0.0386 | 0.0128 | 0.0024 | 0.0014 | 0.0024 | 0.0007 |
| Illness | 24 | 0.0366 | 0.0186 | 0.0018 | 0.0011 | 0.0561 | 0.0168 | 0.0388 | 0.0131 | 0.8467 | 0.1525 | 0.3844 | 0.0673 | 1.3169 | 0.1001 | 1.3169 | 0.1003 | 6.023 | 0.0524 | 0.5739 | 0.0528 | 0.1945 | 0.0299 | 0.2660 | 0.0507 | 0.0887 | 0.0072 | 0.0794 | 0.0098 |
| Illness | 36 | 0.0730 | 0.0477 | 0.0624 | 0.0409 | 0.0414 | 0.0104 | 0.0395 | 0.0096 | 0.3923 | 0.0615 | 0.5034 | 0.0926 | 0.3605 | 0.0246 | 0.4506 | 0.0295 | 0.2656 | 0.0474 | 0.2402 | 0.0386 | 0.0960 | 0.0217 | 0.0774 | 0.0169 | 0.0629 | 0.0305 | 0.0313 | 0.0128 |
| Illness | 48 | 0.2656 | 0.0754 | 0.1651 | 0.0482 | 0.0304 | 0.0065 | 0.0320 | 0.0056 | 0.4186 | 0.0755 | 0.1296 | 0.0311 | 0.0553 | 0.0273 | 0.0506 | 0.0250 | 0.0527 | 0.0039 | 0.1422 | 0.0168 | 0.1000 | 0.0161 | 0.0906 | 0.0263 | 0.0402 | 0.0181 | 0.0223 | 0.0099 |
| Illness | 60 | 0.0427 | 0.0115 | 0.0791 | 0.0162 | 0.0414 | 0.0123 | 0.0264 | 0.0092 | 0.1066 | 0.0203 | 0.2392 | 0.0413 | 0.0543 | 0.0281 | 0.0668 | 0.0305 | 0.0197 | 0.0041 | 0.0565 | 0.0123 | 0.0908 | 0.0181 | 0.0921 | 0.0182 | 0.1070 | 0.0293 | 0.0436 | 0.0114 |
| Illness | Avg | 0.1045 | 0.0383 | 0.0771 | 0.0266 | 0.0424 | 0.0115 | 0.0342 | 0.0094 | 0.4410 | 0.0775 | 0.3142 | 0.0581 | 0.4454 | 0.0450 | 0.4712 | 0.0463 | 0.2351 | 0.0269 | 0.2532 | 0.0301 | 0.1204 | 0.0214 | 0.1315 | 0.0280 | 0.0747 | 0.0213 | 0.0441 | 0.0110 |

Table 9: Full average results of three random seeds with prediction lengths $S \in \{96, 192, 336, 720\}$ and a fixed lookback length $T = \{96\}$ on synthetic data. The original results are from the public datasets ETTh1 and ETTh2, as shown in Table 7. Synthetic datasets were created by adding sine waves with periods of 100, 300, and 1000 to the original ETTh datasets. Gains mean performance improvements by using our method.

| | | iTransformer | | | | | | | | iTransformer | | | | | |
| | | base | | BSA | | Gain(%) | | | | base | | BSA | | Gain(%) | |
| | | MSE | MAE | MSE | MAE | MSE | MAE | | | MSE | MAE | MSE | MAE | MSE | MAE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ETTh1 original | 96 | 0.4316 | 0.4453 | 0.4277 | 0.4426 | 0.909 | 0.605 | ETTh2 original | 96 | 0.2397 | 0.3278 | 0.2347 | 0.3236 | 2.062 | 1.260 |
| | 192 | 0.4906 | 0.4857 | 0.4809 | 0.4812 | 1.973 | 0.931 | | 192 | 0.2950 | 0.3641 | 0.2901 | 0.3624 | 1.654 | 0.467 |
| | 336 | 0.5451 | 0.5217 | 0.5401 | 0.5243 | 0.908 | -0.504 | | 336 | 0.3316 | 0.3880 | 0.3271 | 0.3884 | 1.370 | -0.103 |
| | 720 | 0.6991 | 0.6174 | 0.6979 | 0.6204 | 0.172 | -0.493 | | 720 | 0.4160 | 0.4389 | 0.4144 | 0.4391 | 0.364 | -0.040 |
| | avg | 0.5416 | 0.5175 | 0.5367 | 0.5171 | 0.991 | 0.135 | | avg | 0.3206 | 0.3797 | 0.3166 | 0.3784 | 1.362 | 0.396 |
| ETTh1 100 | 96 | 0.3108 | 0.3889 | 0.2980 | 0.3800 | 4.111 | 2.286 | ETTh2 100 | 96 | 0.1608 | 0.2738 | 0.1565 | 0.2724 | 2.716 | 0.530 |
| | 192 | 0.3470 | 0.4172 | 0.3298 | 0.4062 | 4.968 | 2.635 | | 192 | 0.1880 | 0.2993 | 0.1839 | 0.3009 | 2.198 | -0.516 |
| | 336 | 0.3825 | 0.4454 | 0.3719 | 0.4409 | 2.774 | 1.001 | | 336 | 0.2098 | 0.3156 | 0.2070 | 0.3153 | 1.335 | 0.105 |
| | 720 | 0.4825 | 0.5162 | 0.4784 | 0.5150 | 0.844 | 0.230 | | 720 | 0.2580 | 0.3526 | 0.2558 | 0.3522 | 0.855 | 0.137 |
| | avg | 0.3807 | 0.4419 | 0.3695 | 0.4355 | 3.174 | 1.538 | | avg | 0.2042 | 0.3103 | 0.2008 | 0.3102 | 1.776 | 0.064 |
| ETTh1 300 | 96 | 0.5350 | 0.5291 | 0.3674 | 0.4334 | 31.330 | 18.081 | ETTh2 300 | 96 | 0.2646 | 0.3663 | 0.1711 | 0.2957 | 35.319 | 19.272 |
| | 192 | 0.6278 | 0.5873 | 0.4156 | 0.4729 | 33.799 | 19.476 | | 192 | 0.3219 | 0.4107 | 0.2189 | 0.3356 | 32.010 | 18.290 |
| | 336 | 0.5741 | 0.5608 | 0.4217 | 0.4780 | 26.537 | 14.763 | | 336 | 0.2971 | 0.3921 | 0.2287 | 0.3444 | 23.023 | 12.178 |
| | 720 | 0.7036 | 0.6361 | 0.5273 | 0.5514 | 25.068 | 13.313 | | 720 | 0.3669 | 0.4410 | 0.2804 | 0.3837 | 23.582 | 12.999 |
| | avg | 0.6101 | 0.5783 | 0.4330 | 0.4839 | 29.183 | 16.408 | | avg | 0.3126 | 0.4025 | 0.2248 | 0.3398 | 28.483 | 15.685 |
| ETTh1 1000 | 96 | 0.3370 | 0.4082 | 0.3295 | 0.4036 | 2.225 | 1.135 | ETTh2 1000 | 96 | 0.1880 | 0.3003 | 0.1799 | 0.2927 | 4.305 | 2.529 |
| | 192 | 0.5019 | 0.5197 | 0.4392 | 0.4833 | 12.487 | 6.992 | | 192 | 0.3283 | 0.4108 | 0.2742 | 0.3745 | 16.489 | 8.842 |
| | 336 | 0.8177 | 0.6785 | 0.5441 | 0.5591 | 33.455 | 17.591 | | 336 | 0.6335 | 0.5958 | 0.4320 | 0.4856 | 31.820 | 18.492 |
| | 720 | 1.2900 | 0.8721 | 0.6741 | 0.6411 | 47.742 | 26.484 | | 720 | 1.1044 | 0.7933 | 0.4554 | 0.5046 | 58.760 | 36.389 |
| | avg | 0.7366 | 0.6196 | 0.4967 | 0.5218 | 23.978 | 13.050 | | avg | 0.5636 | 0.5251 | 0.3354 | 0.4144 | 27.843 | 16.563 |

on a log scale. Converting the y-axis from decibels to amplitude is akin to implementing linear scaling, which reveals a steeper rise in noise at lower frequencies. This noise generally arises in the low-frequency bands of most real-world systems, such as electronic systems and environmental sciences [24]. Figure 8 (a) shows that noise escalates within the low-frequency regions across variables, leading to the assumption that unique Brownian noise originates in each system. Consequently, a linear approximation was used for the denoising function, which was performed manually. The outcomes of this process are presented in Figure 8 (b).

## C   Synthetic dataset experiments

We constructed synthetic datasets by adding sine waves with periods of 100, 300, and 1000 to each channel (variable) of the ETTh1 and ETTh2 datasets. The scale of the sine wave is set to the standard deviation value of each channel. The phases of the sine waves for each channel are randomly sampled from a uniform distribution from 0 to $2\pi$ to make decorrelated sine waves. The full results are reported in Table 9.

### C.1   SA-Matrix and FFT visualization

As illustrated in Figure 14, the LUFL, MULL, OT, and LULL channels are organized into columns, with each graph positioned within rows that correspond to the periods added to the dataset. Figure 14 (a) shows the SA-Matrix from the iTransformer model fine-tuned with a 720 prediction length on the ETTh1 data that does not include added sine waves. From (b) to (d), the models were trained with synthetic datasets increasingly augmented with sine waves at periods of 100, 300, and 1000. As the number of added periods increases, there is a shift in the SA-Matrix weight distribution toward low-frequency patterns in each channel.

## D   Analysis and Ablation Studies

### D.1   BSA insertion site analysis

Table 10 shows the full results for different BSA module insertion sites in the iTransformer [31] and DLinear [53] models on the Weather Dataset. Table 2 is from the average value of Table 10. The BSA module insert position can be found in Figure 7.

Table 10: Full results for different BSA module locations for iTransformer [31] and DLinear [53] on the Weather dataset. Positions 1 and 4 are the only available options for the DLinear Model.

| | | iTransformer | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | base | | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | | 7 | |
| | | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| weather | 96 | 0.1706 | 0.2096 | 0.1598 | 0.2037 | 0.1589 | 0.2027 | 0.1538 | 0.1990 | 0.1667 | 0.2077 | 0.1545 | 0.1999 | 0.1676 | 0.2080 | 0.1582 | 0.2011 |
| | 192 | 0.2203 | 0.2544 | 0.2052 | 0.2466 | 0.2034 | 0.2463 | 0.1996 | 0.2441 | 0.2149 | 0.2526 | 0.2003 | 0.2451 | 0.2206 | 0.2548 | 0.2110 | 0.2504 |
| | 336 | 0.2762 | 0.2952 | 0.2526 | 0.2856 | 0.2510 | 0.2846 | 0.2528 | 0.2837 | 0.2719 | 0.2940 | 0.2519 | 0.2855 | 0.2756 | 0.2948 | 0.2534 | 0.2854 |
| | 720 | 0.3551 | 0.3473 | 0.3252 | 0.3371 | 0.3277 | 0.3386 | 0.3255 | 0.3349 | 0.3497 | 0.3465 | 0.3238 | 0.3380 | 0.3513 | 0.3472 | 0.3434 | 0.3440 |
| | Avg | 0.2556 | 0.2766 | 0.2357 | 0.2682 | 0.2352 | 0.2681 | 0.2329 | 0.2654 | 0.2508 | 0.2752 | 0.2326 | 0.2671 | 0.2538 | 0.2762 | 0.2415 | 0.2702 |

| | | Dlinear | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | base | | 1 | | 2 | | 3 | | 4 | | 5 | | 6 | | 7 | |
| | | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| weather | 96 | 0.1630 | 0.2363 | 0.1565 | 0.2306 | - | - | - | - | 0.1519 | 0.2203 | - | - | - | - | - | - |
| | 192 | 0.2091 | 0.2832 | 0.1985 | 0.2724 | - | - | - | - | 0.1912 | 0.2601 | - | - | - | - | - | - |
| | 336 | 0.2634 | 0.3271 | 0.2590 | 0.3251 | - | - | - | - | 0.2375 | 0.2993 | - | - | - | - | - | - |
| | 720 | 0.3421 | 0.3868 | 0.3167 | 0.3697 | - | - | - | - | 0.2979 | 0.3463 | - | - | - | - | - | - |
| | Avg | 0.2444 | 0.3084 | 0.2327 | 0.2994 | - | - | - | - | 0.2196 | 0.2815 | - | - | - | - | - | - |

## D.2 BSA variable input lengths analysis

Table 11: Full result for the performance of the base model and BSA according to changes in Input length. Experiments were conducted on Weather and PEMS03 data using DLinear, RLinear, and iTransformer models.

| | | Dlinear | | | | RLinear | | | | iTransformer | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | base | | BSA | | base | | BSA | | base | | BSA | |
| | | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| Weather 48 | 96 | 0.2187 | 0.2950 | 0.1616 | 0.2332 | 0.1898 | 0.2322 | 0.1783 | 0.2233 | 0.2054 | 0.2271 | 0.1713 | 0.2124 |
| | 192 | 0.2477 | 0.3203 | 0.1994 | 0.2715 | 0.2315 | 0.2654 | 0.2206 | 0.2583 | 0.2490 | 0.2695 | 0.2178 | 0.2541 |
| | 336 | 0.3086 | 0.3696 | 0.2460 | 0.3094 | 0.2897 | 0.3064 | 0.2768 | 0.2980 | 0.3075 | 0.3099 | 0.2716 | 0.2947 |
| | 720 | 0.3817 | 0.4204 | 0.3085 | 0.3585 | 0.3699 | 0.3566 | 0.3579 | 0.3499 | 0.3831 | 0.3595 | 0.3459 | 0.3442 |
| | Avg | 0.2892 | 0.3513 | 0.2289 | 0.2932 | 0.2702 | 0.2902 | 0.2584 | 0.2824 | 0.2862 | 0.2915 | 0.2517 | 0.2763 |
| Weather 96 | 96 | 0.1630 | 0.2363 | 0.1519 | 0.2203 | 0.1656 | 0.2108 | 0.1594 | 0.2065 | 0.1706 | 0.2096 | 0.1598 | 0.2037 |
| | 192 | 0.2091 | 0.2832 | 0.1912 | 0.2601 | 0.2119 | 0.2512 | 0.2013 | 0.2444 | 0.2203 | 0.2544 | 0.2052 | 0.2466 |
| | 336 | 0.2634 | 0.3271 | 0.2375 | 0.2993 | 0.2679 | 0.2911 | 0.2541 | 0.2823 | 0.2762 | 0.2952 | 0.2526 | 0.2856 |
| | 720 | 0.3421 | 0.3868 | 0.2979 | 0.3463 | 0.3465 | 0.3412 | 0.3315 | 0.3338 | 0.3551 | 0.3473 | 0.3252 | 0.3371 |
| | Avg | 0.2444 | 0.3084 | 0.2196 | 0.2815 | 0.2480 | 0.2736 | 0.2366 | 0.2667 | 0.2556 | 0.2766 | 0.2357 | 0.2682 |
| Weather 192 | 96 | 0.1514 | 0.2208 | 0.1466 | 0.2138 | 0.1522 | 0.1982 | 0.1471 | 0.1960 | 0.1650 | 0.2100 | 0.1521 | 0.2008 |
| | 192 | 0.1952 | 0.2660 | 0.1856 | 0.2530 | 0.1963 | 0.2383 | 0.1898 | 0.2360 | 0.2098 | 0.2510 | 0.2023 | 0.2450 |
| | 336 | 0.2492 | 0.3100 | 0.2315 | 0.2914 | 0.2502 | 0.2785 | 0.2393 | 0.2733 | 0.2625 | 0.2903 | 0.2428 | 0.2810 |
| | 720 | 0.3324 | 0.3771 | 0.2958 | 0.3436 | 0.3285 | 0.3313 | 0.3111 | 0.3242 | 0.3356 | 0.3410 | 0.3089 | 0.3291 |
| | Avg | 0.2321 | 0.2934 | 0.2149 | 0.2754 | 0.2318 | 0.2616 | 0.2218 | 0.2574 | 0.2432 | 0.2731 | 0.2265 | 0.2640 |
| PEMS03 48 | 96 | 0.5159 | 0.5565 | 0.4340 | 0.5021 | 1.2342 | 0.8375 | 1.0040 | 0.7564 | 0.4214 | 0.4518 | 0.2934 | 0.3715 |
| | 192 | 0.5673 | 0.5921 | 0.4646 | 0.5287 | 1.7909 | 1.0790 | 1.3785 | 0.9286 | 0.4010 | 0.4454 | 0.2551 | 0.3513 |
| | 336 | 0.5056 | 0.5437 | 0.4018 | 0.4767 | 1.2898 | 0.8454 | 1.0455 | 0.7490 | 0.4196 | 0.4485 | 0.3084 | 0.3801 |
| | 720 | 0.5524 | 0.5811 | 0.4398 | 0.5047 | 1.5496 | 0.9620 | 1.3303 | 0.8777 | 0.4699 | 0.4848 | 0.3507 | 0.4100 |
| | Avg | 0.5353 | 0.5684 | 0.4351 | 0.5031 | 1.4661 | 0.9310 | 1.1896 | 0.8279 | 0.4280 | 0.4576 | 0.3019 | 0.3782 |
| PEMS03 96 | 96 | 0.4405 | 0.5097 | 0.4121 | 0.4875 | 1.0366 | 0.7815 | 0.6789 | 0.6359 | 0.2213 | 0.3226 | 0.1569 | 0.2689 |
| | 192 | 0.4658 | 0.5259 | 0.3876 | 0.4677 | 1.0995 | 0.7937 | 0.7128 | 0.6326 | 0.2687 | 0.3566 | 0.1865 | 0.2949 |
| | 336 | 0.3973 | 0.4669 | 0.3461 | 0.4317 | 0.8227 | 0.6434 | 0.5861 | 0.5408 | 0.2479 | 0.3303 | 0.1979 | 0.2957 |
| | 720 | 0.4420 | 0.5048 | 0.3922 | 0.4703 | 1.0101 | 0.7432 | 0.7194 | 0.6184 | 0.3093 | 0.3716 | 0.2488 | 0.3330 |
| | Avg | 0.4364 | 0.5018 | 0.3845 | 0.4643 | 0.9922 | 0.7405 | 0.6743 | 0.6069 | 0.2618 | 0.3453 | 0.1975 | 0.2981 |
| PEMS03 192 | 96 | 0.3904 | 0.4570 | 0.2660 | 0.3698 | 0.4078 | 0.4618 | 0.2674 | 0.3609 | 0.1463 | 0.2569 | 0.1218 | 0.2323 |
| | 192 | 0.3156 | 0.3932 | 0.2543 | 0.3552 | 0.3290 | 0.3916 | 0.2581 | 0.3394 | 0.1637 | 0.2641 | 0.1450 | 0.2484 |
| | 336 | 0.3078 | 0.3903 | 0.2650 | 0.3641 | 0.3276 | 0.3890 | 0.2723 | 0.3495 | 0.1738 | 0.2680 | 0.1623 | 0.2596 |
| | 720 | 0.3588 | 0.4370 | 0.3221 | 0.4110 | 0.3991 | 0.4418 | 0.3467 | 0.4038 | 0.2226 | 0.3000 | 0.2091 | 0.2924 |
| | Avg | 0.3432 | 0.4194 | 0.2768 | 0.3750 | 0.3659 | 0.4211 | 0.2861 | 0.3634 | 0.1766 | 0.2723 | 0.1596 | 0.2582 |

Table 11 shows the full results of the performance of the base model and BSA according to changes in Input length. Experiments were conducted on Weather and PEMS03 data using DLinear [53], RLinear [28], and iTransformer [31] models. Table 3 in the main text shows the summarized results.

## D.3 BSA computational cost analysis

Table 12 presents a comprehensive analysis of BSA's computational cost. We measured the model training time (sec/1step), peak memory usage (GB), and the number of parameters (M) for both the base model and the model with BSA applied. The experiments were conducted using the lightweight Weather dataset with 21 channels and the heavy PEMS03 dataset with 358 channels. Tests were

Table 12: Full result for the model training time (sec/1 step), peak memory usage (GB), and the number of parameters (M) for both the base model and the model with BSA applied. The experiments were conducted using the lightweight Weather dataset with 21 channels and the heavy PEMS03 dataset with 358 channels. Tests were performed on Timesnet, iTransformer, Crossformer, and PatchTST models.

| | | TimesNet | | | iTransformer | | | Crossformer | | | PatchTST | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Time (sec/1step) | | base | BSA | gain(%) | base | BSA | gain(%) | base | BSA | gain(%) | base | BSA | gain(%) |
| Weather | 96 | 0.066 | 0.070 | 5.37 | 0.024 | 0.029 | 20.71 | 0.076 | 0.079 | 3.58 | 0.045 | 0.046 | 1.52 |
| | 192 | 0.076 | 0.072 | -4.56 | 0.024 | 0.025 | 3.92 | 0.072 | 0.078 | 8.28 | 0.046 | 0.048 | 3.01 |
| | 336 | 0.081 | 0.078 | -3.30 | 0.025 | 0.027 | 8.05 | 0.077 | 0.080 | 3.66 | 0.050 | 0.052 | 5.12 |
| | 720 | 0.104 | 0.107 | 2.47 | 0.025 | 0.033 | 30.74 | 0.075 | 0.080 | 6.50 | 0.048 | 0.049 | 2.88 |
| | Avg | 0.082 | 0.082 | 0.00 | 0.024 | 0.028 | 15.86 | 0.075 | 0.079 | 5.50 | 0.047 | 0.049 | 3.13 |
| PEMS03 | 96 | 1.160 | 1.147 | -1.06 | 0.077 | 0.080 | 3.96 | 0.163 | 0.162 | -0.95 | 0.878 | 0.877 | -0.10 |
| | 192 | 1.791 | 1.717 | -4.12 | 0.079 | 0.082 | 3.69 | 0.296 | 0.290 | -1.98 | 0.889 | 0.882 | -0.74 |
| | 336 | 2.374 | 2.502 | 5.40 | 0.097 | 0.096 | -1.26 | 0.492 | 0.485 | -1.45 | 0.881 | 0.883 | 0.27 |
| | 720 | 4.397 | 4.425 | 0.63 | 0.115 | 0.118 | 2.57 | 0.965 | 0.929 | -3.78 | 0.908 | 0.902 | -0.57 |
| | Avg | 2.430 | 2.448 | 0.21 | 0.092 | 0.094 | 2.24 | 0.479 | 0.466 | -2.04 | 0.889 | 0.886 | -0.29 |
| Memory (GB) | | base | BSA | gain(%) | base | BSA | gain(%) | base | BSA | gain(%) | base | BSA | gain(%) |
| Weather | 96 | 0.43 | 0.43 | 0.04 | 0.20 | 0.25 | 24.65 | 0.27 | 0.28 | 2.85 | 1.59 | 1.60 | 0.36 |
| | 192 | 0.57 | 0.58 | 1.18 | 0.21 | 0.32 | 54.08 | 0.48 | 0.49 | 1.11 | 1.60 | 1.61 | 0.47 |
| | 336 | 0.79 | 0.81 | 2.15 | 0.22 | 0.29 | 30.93 | 0.87 | 0.88 | 0.81 | 1.63 | 1.63 | 0.30 |
| | 720 | 1.38 | 1.37 | -0.07 | 0.25 | 0.30 | 19.57 | 2.28 | 2.28 | 0.24 | 1.68 | 1.68 | 0.25 |
| | Avg | 0.79 | 0.80 | 0.82 | 0.22 | 0.29 | 32.31 | 0.97 | 0.98 | 1.25 | 1.62 | 1.63 | 0.35 |
| PEMS03 | 96 | 5.79 | 6.01 | 3.66 | 3.57 | 3.68 | 3.04 | 7.23 | 7.25 | 0.32 | 25.40 | 25.52 | 0.46 |
| | 192 | 7.32 | 7.59 | 3.74 | 3.65 | 3.75 | 2.76 | 12.24 | 12.26 | 0.22 | 25.49 | 25.60 | 0.44 |
| | 336 | 10.26 | 10.01 | -2.46 | 3.80 | 3.89 | 2.30 | 19.96 | 19.98 | 0.10 | 25.61 | 25.71 | 0.38 |
| | 720 | 17.07 | 16.46 | -3.58 | 4.18 | 4.23 | 1.27 | 42.34 | 42.36 | 0.05 | 25.95 | 26.02 | 0.24 |
| | Avg | 10.11 | 10.02 | 0.34 | 3.80 | 3.89 | 2.34 | 20.44 | 20.46 | 0.17 | 25.61 | 25.71 | 0.38 |
| Parameters (M) | | base | BSA | gain(%) | base | BSA | gain(%) | base | BSA | gain(%) | base | BSA | gain(%) |
| Weather | 96 | 1.354 | 1.368 | 1.05 | 4.848 | 4.852 | 0.08 | 0.283 | 0.301 | 6.41 | 9.46 | 9.48 | 0.15 |
| | 192 | 1.363 | 1.381 | 1.34 | 4.897 | 4.897 | 0.00 | 0.291 | 0.313 | 7.62 | 10.05 | 10.07 | 0.18 |
| | 336 | 1.377 | 1.391 | 1.03 | 4.971 | 4.971 | 0.00 | 0.302 | 0.317 | 4.67 | 10.94 | 10.96 | 0.17 |
| | 720 | 1.414 | 1.433 | 1.29 | 5.168 | 5.172 | 0.08 | 0.333 | 0.347 | 4.24 | 13.30 | 13.32 | 0.14 |
| | Avg | 1.377 | 1.393 | 1.18 | 4.971 | 4.973 | 0.04 | 0.302 | 0.320 | 5.73 | 10.94 | 10.96 | 0.16 |
| PEMS03 | 96 | 151.6 | 151.9 | 0.16 | 4.834 | 5.076 | 5.00 | 10.69 | 10.94 | 2.26 | 9.46 | 9.71 | 2.55 |
| | 192 | 151.6 | 151.9 | 0.16 | 4.883 | 5.125 | 4.95 | 11.45 | 11.69 | 2.11 | 10.05 | 10.30 | 2.40 |
| | 336 | 151.6 | 151.9 | 0.16 | 4.957 | 5.199 | 4.87 | 12.57 | 12.81 | 1.92 | 10.94 | 11.18 | 2.21 |
| | 720 | 151.7 | 151.9 | 0.16 | 5.154 | 5.396 | 4.69 | 15.58 | 15.82 | 1.55 | 13.30 | 13.54 | 1.82 |
| | Avg | 151.6 | 151.9 | 0.16 | 4.957 | 5.199 | 4.88 | 12.57 | 12.81 | 1.96 | 10.94 | 11.18 | 2.25 |

performed on Timesnet [48], iTransformer [31], Crossformer [54], and PatchTST [36] models. Linear models were excluded from the experiments as they are very lightweight and their training cost is not an issue. Table 4 in the main text shows the summarized results.

### D.4 Pre/Post BSA activation signal visualization

Figure 16 demonstrates the transition between pre-activation $F$ (left) and post-activation $F'$ (right) during the inference process. The figure shows the spectrum of randomly sampled activations from the DLinear model trained with a 720 prediction length on the SWDR (W/m) channel of the weather dataset. The x-axis signifies timesteps while the y-axis indicates frequency. To enhance visibility in the low-frequency region, y-axis values below 0.1 are presented in a log scale, whereas values above 0.1 are depicted in a linear scale. A darkening in the high-frequency region is observed from pre-activation to post-activation, suggesting that the activations may undergo a denoising effect in the high-frequency areas as they pass through the SA-Matrix.

Figure 9: The 2D heatmaps of the SA-Matrix from the DLinear model finetuned with a 720 prediction length on the Weather dataset.

Figure 10: The kernel density estimate graphs of the SA-Matrix from the DLinear model finetuned with a 720 prediction length on the Weather dataset and FFT graphs of the data.

Figure 11: The kernel density estimate graphs of the SA-Matrix from the iTransformer model finetuned with a 720 prediction length on the Weather Data and FFT graphs of the data.

Figure 12: The kernel density estimate graphs of the SA-Matrix from the DLinear model finetuned with a 720 prediction length on the ETTh1 dataset and FFT graphs of the data.



Figure 13: The kernel density estimate graphs of the SA-Matrix from the DLinear model finetuned with a 720 prediction length on the ETTm1 dataset and FFT graphs of the data.

Figure 14: The kernel density estimate graphs of LUFL, MULL, OT, and LULL channels of the SA-Matrix from the iTransformer model with a 720 prediction length on the original and synthetic ETTh1 data. Row (a) represents the original data, while rows (b) - (d) display synthetic datasets. These datasets were generated by adding sine waves with periods of 100, 300, and 1000, respectively.



Figure 15: The kernel density estimate graphs of the SA-Matrix from the iTransformer model finetuned with a 720 prediction length on the ETTm1 dataset and FFT graphs of the data.

Figure 16: Illustration of the SWDR (W/m) channel of the Pre-Activation (left) and Post-Activation (right) through SA-Matrix from the DLinear model trained with a 720 prediction length on the weather dataset.

# NeurIPS Paper Checklist

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

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

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

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

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

IMPORTANT, please:

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

1. **Claims**

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

    Answer:[Yes]

    Justification: The paper's contributions and scope are well-explained throughout the abstract and introduction, and the entire paper is written to provide a consistent and comprehensive conclusion.

    Guidelines:

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

2. **Limitations**

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

    Answer: [Yes]

Justification: Our paper discusses the limitations and future research directions to address them in the conclusion section 5.

Guidelines:

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

3. **Theory Assumptions and Proofs**

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

    Answer: [NA]

    Justification: Our paper proposes a new model architecture and demonstrates its superiority through a practical approach with extensive experiments rather than theoretical proofs.

    Guidelines:

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

4. **Experimental Result Reproducibility**

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

    Answer: [Yes]

    Justification: Our paper provides detailed information to reproduce all experiments. Additionally, the complete experiment code will be provided.

Guidelines:

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

5. **Open access to data and code**

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

   Answer: [Yes]

   Justification: https://github.com/DJLee1208/BSA_2024

   Guidelines:

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

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

6. **Experimental Setting/Details**

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

   Answer: [Yes]

   Justification: Our paper provides detailed information about the experiments, including data splits, hyperparameters, how they were chosen, and the type of optimizer, in the main text and the appendix.

   Guidelines:

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

7. **Experiment Statistical Significance**

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

   Answer: [Yes]

   Justification: In the main results, we used a paired t-test to demonstrate the superiority of our methodology and reported the p-value.

   Guidelines:

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

8. **Experiments Compute Resources**

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

   Answer: [Yes]

   Justification: We reported the types and number of GPUs used in the experiments, as well as the packages used for training.

Guidelines:

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

9. **Code Of Ethics**

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

Answer: [Yes]

Justification: Yes, the research conducted in the paper fully conforms to the NeurIPS Code of Ethics in every respect. We have ensured that all ethical guidelines and standards have been meticulously followed throughout the study.

Guidelines:

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

10. **Broader Impacts**

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

Answer: [Yes]

Justification: Our paper discusses the potential broader impacts of our research in the conclusion.

Guidelines:

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

11. **Safeguards**

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

Answer: [NA] .

Justification: Our model does not pose such risks, and therefore, safeguards for responsible data or model release were not applicable.

Guidelines:

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

12. **Licenses for existing assets**

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

Answer: [Yes]

Justification: Yes, we have properly credited the creators or original owners of assets used in the paper, including publicly available models and packages, with accurate citations.

Guidelines:

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

13. **New Assets**

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

Answer: [NA]

Justification: This paper does not include new assets.

Guidelines:

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

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

14. **Crowdsourcing and Research with Human Subjects**

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

    Answer: [NA]

    Justification: This paper does not involve crowdsourcing.

    Guidelines:

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

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

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

    Answer:[NA]

    Justification: This paper does not involve research with human subjects.

    Guidelines:

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