# Powerformer: A Transformer with Weighted Causal Attention for Time-series Forecasting

Kareem Hegazy[1,2*]        Michael W. Mahoney[1,2,3]        N. Benjamin Erichson[2,3*]

[1]Department of Statistics, University of California Berkeley
[2]International Computer Science Institute, Berkeley
[3]Lawrence Berkeley National Laboratory

## Abstract

Transformers have recently shown strong performance in time-series forecasting, but their all-to-all attention mechanism overlooks the (temporal) causal and often (temporally) local nature of data. We introduce Powerformer, a novel Transformer variant that replaces noncausal attention weights with causal weights that are reweighted according to a smooth heavy-tailed decay. This simple yet effective modification endows the model with an inductive bias favoring temporally local dependencies, while still allowing sufficient flexibility to learn the unique correlation structure of each dataset. Our empirical results demonstrate that Powerformer not only achieves state-of-the-art accuracy on public time-series benchmarks, but also that it offers improved interpretability of attention patterns. Our analyses show that the model's locality bias is amplified during training, demonstrating an interplay between time-series data and power-law-based attention. These findings highlight the importance of domain-specific modifications to the Transformer architecture for time-series forecasting, and they establish Powerformer as a strong, efficient, and principled baseline for future research and real-world applications.

## 1 Introduction

With the recent popularity of Transformer models [36], large language models (LLMs), and foundation models [4], the Transformer architecture gained popularity in time-series forecasting, leading to a number of Transformer based state-of-the-art models: Autoformer [42], FEDformer [48], and PatchTST [27].

The multihead attention (MHA) mechanism is at the core of these models, which is an all-to-all process that enriches each embedding via a weighted combination of other embeddings based on a similarity calculation. For natural language processing (NLP) tasks, each embedding represents a word. Here, an all-to-all embedding mechanism is natural since correlations between words can be only weakly dependent on their ordering and the number of words between them. However, for time-series data, where each embedding is a temporal measurement, the relation between embeddings is both causal and dependent on the time delay. This disparity between the attention mechanism's all-to-all structure motivates the open question: how well-suited is the Transformer's attention mechanism for time-series data? For example, [44] shows that linear models outperform Transformers on common time-series benchmarks. Still, Transformers have seen success in time-series forecasting.

Previous methods sought to impose causal and local implicit biases on Transformer. Before Transformer's widespread adoption, methods such as WaveNet [34] and [12] used causal convolutions. Imparting local and causal implicit biases to MHA requires altering its architecture. Reformer [19] calculates the attention weights via a causal and locality-aware hash map. LogSparse self-attention [20] uses convolutions to encode local variations and feeds them into a causally sparse self-attention mechanism. The work most similar to ours is ETSformer [39], which develops exponential smoothing attention that replaces the similarity-based attention weights with exponential decaying temporal weights. These temporal weights are no longer dependent on a similarity metric like the key-query projection. ETSformer also employs a frequency attention scheme alongside their exponentially weighted attention.

---

*Corresponding authors: Kareem Hegazy (khegazy@berkeley.edu), and N. Benjamin Erichson (erichson@icsi.berkeley.edu).

**a**

X ∈ ℝ^{T×M} — Multivariate Output
X ∈ ℝ^{M×1×T} — Batched Univariate Output
X ∈ ℝ^{M×D×N} — Flatten & Linear Head
X ∈ ℝ^{M×D×N} — Transformer Encoder
X ∈ ℝ^{M×D×N} — Projection Embedding
X ∈ ℝ^{M×P×N} — Norm & Patching
X ∈ ℝ^{M×1×L} — Batched Univariate Input
X ∈ ℝ^{L×M} — Multivariate Input

Add & Norm
Feed Forward
Add & Norm
**Weighted Causal Multi-head Attention**
Weighted Averaging V
Softmax of $\mathbf{S}^{(C,D)}$
Add Mask $\mathbf{M}^{(C,D)}$
Score $\mathbf{S} = \mathbf{KQ}^T$
K   Q   V

**b** Scores: $\mathbf{S} = \mathbf{KQ}^T$     Attention Weights
$\Sigma \Big[ \quad \Big] = $

**c** Scores: $\mathbf{S} = \mathbf{KQ}^T$     Mask: $\mathbf{M}^{(C,D)}$     Weighted Causal Attention Weights
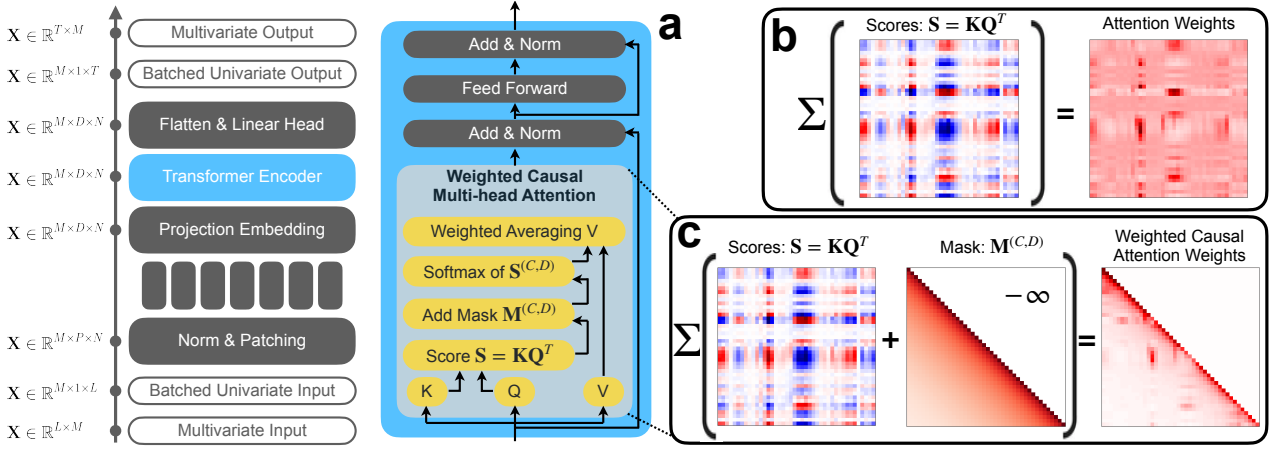$\Sigma \Big[ \quad + \quad -\infty \quad \Big] = $

Figure 1: Illustration of *Powerformer* and the Weighted Causal Multihead Attention (WCMHA) architecture, as well as their effects on attention weights. Panel (a) shows the *Powerformer* architecture (left) and the WCMHA (right). Panels (b) and (c) show the attention weights without and with our local-causal mask, respectively. Here, $\Sigma$ corresponds to the softmax function.

When enforcing a locality bias, previous methods imposed temporal length scales on the dynamics and/or the functional form of the temporal correlations. However, the temporal correlations between data points vary in both functional form and time scale between datasets. Temporal convolutions, and their dilations, are examples of imposing such a time scale. Exponentially weighted attention imposes exponentially decaying temporal correlations but lacks the similarity-weighted MHA that captures varying temporal dependencies. Previous successes from adding local and causal biases [39, 19, 20, 34, 12] show their importance, but the application method is crucial.

In this work, we develop *Powerformer* (illustrated in Fig. 1), a Transformer-based model that uses Weighted Causal Multihead Attention (WCMHA) to learn temporal dependencies unique to each dataset. Our model is motivated by the observation that many physical systems exhibit heavy-tailed autocorrelations, e.g., the pairwise correlation strength may decay as a power law distribution, as the time delay grows [9]. Our goal is to impose this power-law temporal dependence as an inductive bias, while allowing *Powerformer* the flexibility to perturb it sufficiently to capture the temporal pairwise correlations unique to each dataset. To do so, we add a temporally decaying mask to the attention mechanism, specifically to the key-query overlap (referred to as the similarity score). The mask decays attention weights and pairwise dependencies to resemble a power law, but its additive nature maintains the key and query similarity score. This maintained similarity score allows the model to learn data-dependent perturbations to the power-law pairwise dependency. The mask acts as a regularizer that guides the attention weight temporal dependence towards a power law, but still maintains sufficient expressivity for MHA to learn pairwise dependencies unique to each dataset. Moreover, *Powerformer's* simplicity promotes interpretability studies that provide evidence of learned heavy-tailed dependencies in the attention weights.

We build *Powerformer* by combining WCMHA, the commonly used encoder-only Transformer architecture, and patching [27]. *Powerformer* outperforms previous state-of-the-art models on typical tasks despite its simplicity compared to other attention mechanisms. This suggests that these biases and our implementation better align Transformer with natural time-series structures.

Our main contributions are the following.

- We develop *Powerformer*, a Transformer-based model for time-series, with WCMHA to impose causal and local implicit biases that shape attention's pairwise correlation to better respect the data's unique temporal structure.

- *Powerformer* is powered by various weighting schemes: power-law decays and Butterworth filters. The former resembles naturally occurring time-series, while the latter acts analogously to a step function.

- Our empirical results demonstrate that *Powerformer* obtains state-of-the-art performance when compared to similar Transformer models

on a range of public times-series benchmarks. Moreover, we provide interpretability studies that show how our implicit biases shape *Powerformer's* learned attention representations.

- *Powerformer* can reduce the attention complexity from quadratic to linear by defining a cutoff time after which the attention weighting removes their contribution.

## 2 Related Work

The importance of time-series data has led to decades of work, resulting in a myriad of methods: early statistical methods [17, 18, 33, 38] (such as ARIMA [5]); multilayer perceptions (MLP) [44, 21, 10, 46] (such as NHITS [7], TSMixer [8], and N-BEATS [28]); convolutional neural networks (CNNs) [15, 12, 3] (such as TimesNet [40] and SCINet [22]); recurrent neural networks (RNNs) [31, 16] (such as DeepAR [30]); Transformers [49, 13, 39, 48, 26, 41, 43, 23, 47, 19, 20, 32, 11] (such as PatchTST [27], Autoformer [42], FEDformer [48], and iTransformer [24]); and state space models (SSMs) [1, 37, 45] (such as MAMBA [14]). Of these methods, we focus on Transformers.

The success of LLMs and the introduction of powerful foundation models, like GPT-2 [29], inspired models built atop them and others looking to improve upon the Transformer architecture. For example, Time-GPT [13] uses the Transformer encoder-decoder structure, and One-Fits-All [49] is built atop a pretrained GPT-2 and fine-tunes the attention layer normalization and temporal embedding. Some models like AutoTimes [25] leverage LLM's in-context learning capabilities and employ textual prompt engineering to feed into LLMs. Other models have adapted the attention mechanism to time-series data: some models strongly rely on, or partially include, Fourier representations when calculating attention weights [48, 39]; AutoFormer [42] leverages temporal autocorrelations as a similarity metric; Reformer [19] uses a locality-sensitive hashing mechanism to reduce the attention's complexity; Informer [47] selects only high attention weight values, increasing prediction capacity; FlowFormer [41] uses flow conservation to achieve linear complexity attention without imposing locality implicit biases; and iTransformer [24] inverts the temporal and embedding dimensions, applying attention along the embedding.

Transformers were developed for discrete inputs, unlike the continuous valued time-series inputs.

Patching along the time domain (PatchTST) shortens the context window which allows for longer lookback windows and significantly improves performance over more complicated Transformer methods [27]. This simple yet powerful improvement has made patching ubiquitous. Other works look to convert time-series into discretized tokens similar to word vectors. Chronos [2] normalizes the input sequences and matches these continuous values to tokens made from discretizing a sufficiently large domain. TOTEM [32] similarly discretizes the input time series, but instead VQVAE [35] encodes the input and matches it with the most similar token.

## 3 Method

We propose Weighted Causal Multihead Attention (WCMHA) as a foundation for *Powerformer*, a Transformer-based architecture designed for time-series forecasting. In this section, we review standard Multihead Attention (MHA) and demonstrate how we introduce locality and causality into it. We then present various weighting functional forms and insights on how WCMHA modifies attention distributions. We then introduce the *Powerformer* architecture that builds upon these principles. Figure 1 provides a high-level illustration of WCMHA and its effect on attention weights.

### 3.1 Weighted Causal Multihead Attention

Transformers typically leverage MHA to compute new embeddings from weighted sums of input embeddings. Let $\mathbf{X} \in \mathbb{R}^{T \times d}$ be a sequence of $T$ input vectors, each of dimension $d$. For each attention head $h$, MHA first projects $\mathbf{X}$ into *query*, *key*, and *value* matrices, respectively denoted by $\mathbf{Q}_h$, $\mathbf{K}_h$, and $\mathbf{V}_h$. Concretely,

$$\mathbf{Q}_h = \mathbf{X}\,\mathbf{W}_h^{(Q)}, \quad \mathbf{K}_h = \mathbf{X}\,\mathbf{W}_h^{(K)}, \quad \mathbf{V}_h = \mathbf{X}\,\mathbf{W}_h^{(V)}, \tag{1}$$

where each projection matrix $\mathbf{W}_h^{(\cdot)} \in \mathbb{R}^{d \times d_k}$ is learnable and $d_k$ is the dimensionality of each head. The parameters $\mathbf{Q}_h$, $\mathbf{K}_h$, and $\mathbf{V}_h$ thus each have shape $T \times d_k$.

The *unnormalized attention similarity scores* $\mathbf{S}_h$ measure how each query interacts with all keys:

$$\mathbf{S}_h = \frac{\mathbf{K}_h\,\mathbf{Q}_h^\top}{\sqrt{d_k}}. \tag{2}$$

The factor of $1/\sqrt{d_k}$ helps stabilize gradients by normalizing the dot products according to the head

dimension. We then apply a softmax function to obtain the *normalized attention weights*:

$$\mathbf{C}_h = \mathrm{Softmax}(\mathbf{S}_h), \qquad (3)$$

where each row of $\mathbf{C}_h$ sums to 1. Finally, each head's output $\tilde{\mathbf{X}}_h$ is computed as a weighted sum of the values:

$$\tilde{\mathbf{X}}_h = \mathbf{C}_h \mathbf{V}_h. \qquad (4)$$

To form the final MHA output, we concatenate the outputs from all $H$ heads and apply a linear projection:

$$\tilde{\mathbf{X}} = [\tilde{\mathbf{X}}_1, \ldots, \tilde{\mathbf{X}}_H]\,\mathbf{W}^{(\mathrm{A})}, \qquad (5)$$

where $\mathbf{W}^{(\mathrm{A})} \in \mathbb{R}^{(H\,d_k)\times d}$ is another learnable matrix. The result $\tilde{\mathbf{X}}$ has the same dimensionality as the input $\mathbf{X}$ (i.e., $T \times d$), allowing it to be passed into subsequent Transformer layers or decoded for downstream tasks.

### 3.1.1 Enforcing Causality

To ensure that information at time $t$ does not leak from future time steps $t' > t$, we apply a causal mask $\mathbf{M}^{(\mathrm{C})}$:

$$\mathbf{M}^{(\mathrm{C})}_{i,j} = \begin{cases} -\infty & j > i, \\ 0 & j \leq i. \end{cases} \qquad (6)$$

This mask is added to the attention score $\mathbf{S}_h$, forcing any attention weight for $j > i$ to be zero after the softmax:

$$\mathbf{S}^{(\mathrm{C})}_h = \mathbf{S}_h + \mathbf{M}^{(\mathrm{C})}, \quad \mathbf{C}^{(\mathrm{C})}_h = \mathrm{Softmax}(\mathbf{S}^{(\mathrm{C})}_h).$$

Hence, $\mathbf{C}_{t,t'} = 0$ for all $t' > t$, ensuring causal structure for time-series data.

### 3.1.2 Incorporating Locality via Decaying Masks

Time-series often exhibit stronger correlations between nearby time steps compared to distant ones [9]. To bias the model toward local dependencies, we introduce a second mask $\mathbf{M}^{(\mathrm{D})}$ that decays attention weights as the time gap $|\Delta t|$ increases:

$$\mathbf{M}^{(\mathrm{D})}_{i,j} = \begin{cases} 0 & j > i, \\ f(t_i - t_j) & j \leq i, \end{cases} \qquad (7)$$

where $f(\Delta t) \leq 0$ is a non-positive function of the time difference. Adding $\mathbf{M}^{(\mathrm{D})}$ to $\mathbf{S}^{(\mathrm{C})}_h$ yields

$$\mathbf{S}^{(\mathrm{C,D})}_h = \mathbf{S}_h + \mathbf{M}^{(\mathrm{C})} + \mathbf{M}^{(\mathrm{D})}. \qquad (8)$$
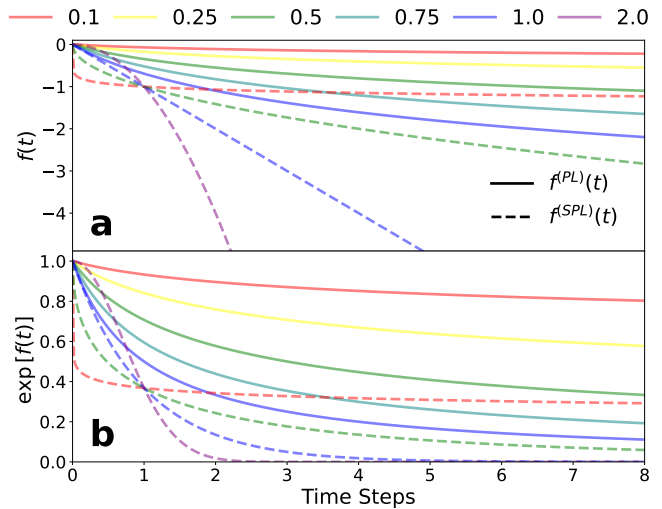


Figure 2: We show the weight power-law (solid line) and similarity power-law (dashed line) masks for varying $\alpha$. Panel (a) shows the contribution added to the attention scores and Panel (b) shows the subsequent effects on the attention weights after applying Softmax.

We thus obtain the *Weighted Causal Attention Weights*

$$\mathbf{C}^{(\mathrm{C,D})}_h = \mathrm{Softmax}(\mathbf{S}^{(\mathrm{C,D})}_h), \qquad (9)$$

where,

$$C^{(\mathrm{C,D})}_{i,j} \propto \begin{cases} 0 & j > i, \\ \exp\big[f(t_i - t_j)\big] & j \leq i. \end{cases} \qquad (10)$$

We refer to this mechanism as WCMHA. By combining causal masking and smoothly decaying weights, WCMHA captures both the temporal ordering and localized dependencies inherent to many real-world processes.

## 3.2 Reweighting Functions for Locality

We explore multiple choices for the reweighting function $f(\Delta t)$:

- **Similarity Power Law:** $f^{(\mathrm{SPL})}(t)(\Delta t) = -(\Delta t)^{\alpha}$.

- **Weight Power Law:** $f^{(\mathrm{PL})}(t)(\Delta t) = -\alpha \log(\Delta t)$.

- **Butterworth Filter:** A frequency-based filter with a smooth transition resembling a step function, parameterized by its cutoff time $t_c$ and order $n$.

Each function has a tunable hyperparameter (e.g., $\alpha$ or $t_c$) that governs the decay rate of attention

as $|\Delta t|$ grows. Figure 1c illustrates the conceptual shift in attention scores when applying these decays. Figures 2 and A1 illustrate the masks' contributions to the attention score and weights for $f^{(\text{PL})}(t)$ and $f^{(\text{SPL})}(t)$, and $f_n^{(\text{BW})}(t)$, respectively. The $f^{(\text{PL})}(t)$ mask imparts a power-law envelope over the attention weight distribution, whereas the $f^{(\text{SPL})}(t)$ envelope is the exponential of a power-law. Therefore, $f^{(\text{SPL})}(t)$, which applies the power-law in the similarity score space, decays faster than $f^{(\text{PL})}(t)$ in weight space (as seen in Fig. 2).

### 3.3 Insights into WCMHA

Before integrating WCMHA into a state-of-the-art Transformer-based model, we investigate how causal and local masks reshape attention distributions in a vanilla Transformer. See Appendix E.3 for details. Figure 3 compares attention score and weight distributions (encoder self-attention, decoder self-attention, and decoder cross-attention) on the Electricity dataset. We observe:

- **Encoder Self-Attention** (modified by WCMHA) shifts toward higher positive values, emphasizing local contexts.

- **Decoder Cross-Attention** (unchanged) retains a heavy-tailed distribution, indicating some long-range dependencies remain relevant.

- **Decoder Self-Attention** (already causal) exhibits minimal distributional change.

These shifts consistently improve forecasting performance on real-world datasets (see Appendix E.3). In many cases, WCMHA significantly outperforms standard MHA, especially for datasets with slowly decaying long-range dependencies. Interestingly, when WCMHA outperforms MHA it does so by much larger margins than when it underperforms. For example, on the ETTh2 dataset, WCMHA improves MSE and MAE by 16% and underperforms by 2%.

### 3.4 Powerformer: Locality, Causality, and Patching

Having shown that WCMHA provides a powerful mechanism for time-series, leveraging its full potential in practical scenarios requires an effective approach to represent temporal data. While incorporating WCMHA into modern Transformer time-series models, we wish to isolate the effects of WCMHA and its local and causal biases. To
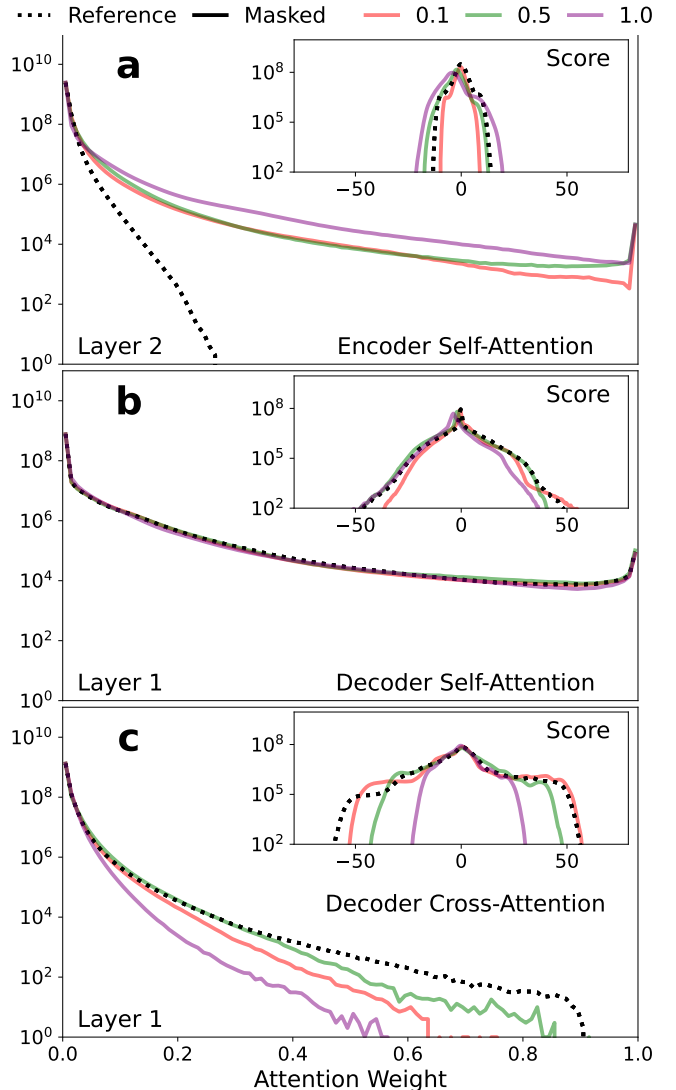


Figure 3: We show the attention score and weight distributions for both the benchmark Transformer (dotted black line) with MHA and our modified Transformer with WCMHA and $f^{(\text{PL})}(t)$ (solid colored lines). Panels (a), (b), and (c) correspond to the last encoder self-attention, decoder self-attention, and decoder cross-attention layers, respectively. The colored lines correspond to different mask decay times ($\alpha$). These results are from the Electricity dataset with a 96 prediction length and 512 input length.

do so requires a simple architecture, and we focus on univariate time-series. Thus, we introduce *Powerformer* (see Fig. 1), a model that integrates Transformer, WCMHA, and PatchTST's [27] successful patching framework.

Patching is a common tool in modern time-series models, and similar to $\mathbf{M}^{(\text{D})}$ it acts as a regularizer to remove high-frequency noise components. We decompose multivariate time-series into univariate

channels before applying *convolutional patching*, similar to PatchTST [27]:

$$\mathbf{X}^{(\text{Patched})}_{s,w} = \mathbf{X} * \mathbf{p}\Big|_s,$$

with patch length $w$ and convolution stride $s$. Patching captures local temporal contexts by converting segments of the raw series into compact embeddings for the Transformer.

Many state-of-the-art models only use the Transformer encoder [27, 49, 32] with a linear decoder head. *Powerformer* combines these insights with our previous Transformer experiments by using only the Transformer encoder with WCMHA and a similar linear readout head. *Powerformer* integrates two modern design choices while injecting causal masks and smoothly decaying attention weights, allowing it to better model nearby past time points.

Our design integrates two complementary strengths:

- **Locality and Causality:** WCMHA introduces both temporal ordering and tailored decay, making it naturally suited for real-world time-series with power-law or exponential correlations.

- **Robust Patch Embeddings:** Patching [27] efficiently captures short-term trends and reduces input length, mitigating quadratic complexity issues that can impede Transformers on long time-series.

*Powerformer* is expressive, computationally efficient, and sufficiently simple to interpret. Appendix C.3 provides further details. Our empirical results confirm that imposing explicit locality and causality boosts forecasting accuracy and the interpretability of learned temporal dependencies.

# 4 Empirical Evaluation

We evaluate *Powerformer*[1], investigating the effects of WCMHA on the attention score and weight distributions, using the mean squared error (MSE) and mean absolute error (MAE) to quantify aggregate performance.

## 4.1 Datasets

We evaluate *Powerformer* on 7 popular publicly-available real-world datasets [47, 42]: ETT (4 subsets), Weather, Electricity, and Traffic. We are

particularly interested in the ETT subsets which measure the electricity transformer oil temperature at 2 different time scales (every 15 minutes and hourly) for two regions. Analyzing the same measurement in different locations and different time scales will illustrate how sensitive WCMHA is to sampling changes and distribution shifts. Section B provides a detailed description of these datasets, available at [42]. We note that Weather, Electricity, and Traffic have much larger timescales, spanning from one to multiple years (Table A1). These 7 datasets are large enough to capture slowly decaying temporal correlations, shown in Appendix B. Illness [47] is traditionally used as a long-term forecasting benchmark. We omit this dataset as its temporal measurements (Table A1) and evaluation look-back window are more than an order of magnitude smaller than the other datasets. This short look-back window is typically the size of our filters.

## 4.2 Evaluation setup

Following previous works [27, 42, 48, 39, 24, 32, 19, 47], we train for multiple prediction and input lengths. The prediction lengths are 96, 192, 336, and 720, while the input sequence lengths are 336 and 512. Each model is trained three times for each setting. Further details on our experimental process and hyperparameters are provided in Appendices E.1 and C.4, respectively.

## 4.3 Empirical Results

*Powerformer* achieves state-of-the-art performance as we compare it against other state-of-the-art models in Table 1. *Powerformer* achieves the best performance in 47 forecasting tasks compared to the next best model (PatchTST) which achieves 17. We note that *Powerformer* outperforms ETSformer in all tasks, further supporting our claim that how and where locally decaying attention weights are applied is crucial. Notably, *Powerformer* outperforms One-Fits-All [49] which is more than twice *Powerformer's* size and is built from pre-trained GPT-2 [29].

*Powerformer* is sensitive to the mask type and aligns $f(t)$ with the dataset's natural pairwise correlation distribution, which we observe as $f^{(\text{PL})}(t)$ and $f^{(\text{SPL})}(t)$ consistently outperforming $f_n^{(\text{BW})}(t)$. See Tables A5, A6, A7, and A8 for full $f^{(\text{PL})}(t)$, $f^{(\text{SPL})}(t)$, $f_1^{(\text{BW})}(t)$, and $f_2^{(\text{BW})}(t)$ evaluation results. The $f_n^{(\text{BW})}(t)$ filters resemble step functions and similarly weight all the values within the win-

---

[1]https://github.com/khegazy/Powerformer

Table 1: We compare *Powerformer's* performance on standard publicly-available time-series datasets to state-of-the-art Transformer models. The best results are bolded and second best are underlined. *Powerformer* and PatchTST are evaluated with look-back windows of 336 and 512, see Section E.1 for selection criteria. TOTEM [32] and iTransformer [24] were evaluated with $T_{\text{seq}} = 96$. All other models were evaluated with multiple look-back windows with the best results selected.

| | | Powerformer | | PatchTST | | One-Fits-All | | TOTEM | | iTransformer | | FEDformer | | ETSformer | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metric | | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 96 | **0.369** | 0.399 | <u>0.370</u> | 0.400 | 0.376 | <u>0.397</u> | 0.380 | **0.394** | 0.386 | 0.405 | 0.376 | 0.415 | 0.494 | 0.479 |
| | 192 | **0.402** | **0.418** | <u>0.413</u> | 0.429 | 0.416 | **0.418** | 0.434 | <u>0.427</u> | 0.441 | 0.436 | 0.423 | 0.446 | 0.538 | 0.504 |
| | 336 | **0.414** | **0.428** | <u>0.422</u> | 0.440 | 0.442 | <u>0.433</u> | 0.490 | 0.459 | 0.487 | 0.458 | 0.444 | 0.462 | 0.574 | 0.521 |
| | 720 | **0.439** | <u>0.459</u> | <u>0.447</u> | 0.468 | 0.477 | **0.456** | 0.539 | 0.513 | 0.503 | 0.491 | 0.469 | 0.492 | 0.562 | 0.535 |
| ETTh2 | 96 | **0.274** | **0.335** | **0.274** | <u>0.336</u> | <u>0.285</u> | 0.342 | 0.293 | 0.338 | 0.297 | 0.349 | 0.332 | 0.374 | 0.340 | 0.391 |
| | 192 | **0.338** | **0.379** | <u>0.339</u> | **0.379** | 0.354 | <u>0.389</u> | 0.375 | 0.390 | 0.380 | 0.400 | 0.407 | 0.446 | 0.430 | 0.439 |
| | 336 | **0.325** | **0.379** | <u>0.331</u> | <u>0.380</u> | 0.373 | 0.407 | 0.422 | 0.431 | 0.428 | 0.432 | 0.400 | 0.447 | 0.485 | 0.479 |
| | 720 | **0.376** | **0.419** | <u>0.379</u> | <u>0.422</u> | 0.406 | 0.441 | 0.610 | 0.567 | 0.427 | 0.445 | 0.412 | 0.469 | 0.500 | 0.497 |
| ETTm1 | 96 | **0.290** | **0.342** | **0.290** | **0.342** | <u>0.292</u> | <u>0.346</u> | 0.320 | 0.347 | 0.334 | 0.368 | 0.326 | 0.390 | 0.375 | 0.398 |
| | 192 | **0.329** | **0.369** | <u>0.332</u> | **0.369** | <u>0.332</u> | <u>0.372</u> | 0.379 | 0.382 | 0.377 | 0.391 | 0.365 | 0.415 | 0.408 | 0.410 |
| | 336 | **0.359** | **0.389** | <u>0.366</u> | <u>0.392</u> | <u>0.366</u> | 0.394 | 0.406 | 0.402 | 0.426 | 0.420 | 0.392 | 0.425 | 0.435 | 0.428 |
| | 720 | **0.412** | **0.421** | 0.420 | <u>0.424</u> | <u>0.417</u> | **0.421** | 0.471 | 0.438 | 0.491 | 0.459 | 0.446 | 0.458 | 0.499 | 0.462 |
| ETTm2 | 96 | **0.162** | **0.252** | <u>0.165</u> | 0.255 | 0.173 | 0.262 | 0.176 | <u>0.253</u> | 0.180 | 0.264 | 0.180 | 0.271 | 0.189 | 0.280 |
| | 192 | <u>0.222</u> | **0.292** | **0.220** | **0.292** | 0.229 | <u>0.301</u> | 0.247 | 0.302 | 0.250 | 0.309 | 0.252 | 0.318 | 0.253 | 0.319 |
| | 336 | **0.275** | **0.327** | <u>0.278</u> | <u>0.329</u> | 0.286 | 0.341 | 0.317 | 0.348 | 0.311 | 0.348 | 0.324 | 0.364 | 0.314 | 0.357 |
| | 720 | **0.351** | **0.380** | <u>0.367</u> | <u>0.385</u> | 0.378 | 0.401 | 0.426 | 0.410 | 0.412 | 0.407 | 0.410 | 0.420 | 0.414 | 0.413 |
| Weather | 96 | **0.147** | **0.198** | <u>0.149</u> | **0.198** | 0.162 | 0.212 | 0.165 | <u>0.208</u> | 0.174 | 0.214 | 0.238 | 0.314 | 0.197 | 0.281 |
| | 192 | **0.191** | **0.239** | <u>0.194</u> | <u>0.241</u> | 0.204 | 0.248 | 0.207 | 0.250 | 0.221 | 0.254 | 0.275 | 0.329 | 0.237 | 0.312 |
| | 336 | **0.243** | **0.279** | <u>0.245</u> | <u>0.282</u> | 0.254 | 0.286 | 0.257 | 0.291 | 0.278 | 0.296 | 0.339 | 0.377 | 0.298 | 0.353 |
| | 720 | **0.310** | **0.329** | <u>0.314</u> | <u>0.334</u> | 0.326 | 0.337 | 0.326 | 0.340 | 0.358 | 0.347 | 0.389 | 0.409 | 0.352 | 0.388 |
| Electricity | 96 | **0.129** | <u>0.223</u> | **0.129** | 0.222 | <u>0.139</u> | 0.238 | 0.178 | 0.263 | 0.148 | 0.240 | 0.186 | 0.302 | 0.187 | 0.304 |
| | 192 | **0.145** | **0.240** | <u>0.147</u> | **0.240** | 0.153 | <u>0.251</u> | 0.187 | 0.272 | 0.162 | 0.253 | 0.197 | 0.311 | 0.199 | 0.315 |
| | 336 | **0.162** | **0.257** | <u>0.163</u> | <u>0.259</u> | 0.169 | 0.266 | 0.199 | 0.285 | 0.178 | 0.269 | 0.213 | 0.328 | 0.212 | 0.329 |
| | 720 | <u>0.198</u> | **0.289** | **0.197** | <u>0.290</u> | 0.206 | 0.297 | 0.236 | 0.318 | 0.225 | 0.317 | 0.233 | 0.344 | 0.233 | 0.345 |
| Traffic | 96 | <u>0.365</u> | <u>0.251</u> | **0.360** | **0.249** | 0.388 | 0.282 | 0.523 | 0.303 | 0.395 | 0.268 | 0.376 | 0.415 | 0.607 | 0.392 |
| | 192 | <u>0.382</u> | <u>0.258</u> | **0.379** | **0.256** | 0.407 | 0.290 | 0.530 | 0.303 | 0.417 | 0.276 | 0.423 | 0.446 | 0.621 | 0.399 |
| | 336 | **0.391** | **0.264** | <u>0.392</u> | **0.264** | 0.412 | 0.294 | 0.549 | 0.311 | 0.433 | <u>0.283</u> | 0.444 | 0.462 | 0.622 | 0.396 |
| | 720 | **0.430** | **0.285** | <u>0.432</u> | <u>0.286</u> | 0.450 | 0.312 | 0.598 | 0.331 | 0.467 | 0.302 | 0.469 | 0.492 | 0.632 | 0.396 |
| 1st/2nd | | **47 / 8** | | <u>17</u> / 33 | | 3 / 13 | | 1 / 3 | | 0 / 1 | | 0 / 0 | | 0 / 0 | |

dow. Alternatively, $f^{(\text{PL})}(t)$ and $f^{(\text{SPL})}(t)$ allow for longer-time couplings but significantly diminish their contribution; see Fig. 1c. The superior performance of $f^{(\text{PL})}(t)$ and $f^{(\text{SPL})}(t)$ over $f_n^{(\text{BW})}(t)$ further supports that the method for enforcing locality is important and that power-law correlation distributions are naturally suited to our datasets. We perform further ablation studies and report detailed results in Appendix F.

In addition to distinguishing between masks that naturally align with the data distribution, *Powerformer* is also sensitive to the pairwise correlation timescale. The ETTh and ETTm datasets test *Powerformer's* ability to adapt the timescale $\alpha$ of $\mathbf{M}^{(\text{D})}$ for datasets sampled from the same distribution but with different temporal resolution. Since ETTh and ETTm are sampled on the hour and

minute timescales, respectively, ETTm requires a slowly decaying mask to access the same information as a quickly decaying mask for ETTh. As expected, *Powerformer* with $f^{(\text{SPL})}(t)$ (Table A6) favors quickly decaying masks on ETTh while selecting slowly decaying masks on ETTm. This illustrates *Powerformer's* sensitivity to both the temporal information timescale, the timescale covering information needed to accurately forecast, and the temporal sampling rate.

To learn the dataset's temporal information timescale, we consider a learnable $\alpha$ instead of treating it as a hyperparameter. Section F.3 outlines method and experimental details, as well as the full results. During training, we observe that $\alpha$ consistently and monotonically drops. This is expected, as it facilitates overfitting by providing the
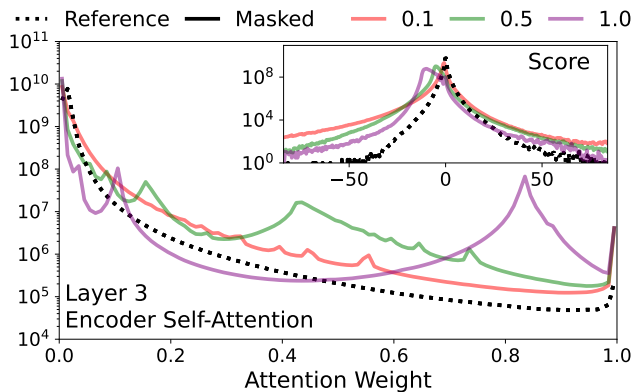
Figure 4: We show *Powerformer's* attention score and weight distributions with MHA (dotted line) and with WCMHA (solid lines) for $f^{(PL)}(t)$. The colored lines correspond to different mask decay times ($\alpha$). These results are for the Weather dataset with a 96 prediction length and 512 input length.

Figure 5: We show the causal and local biases' implicit and explicit effects on *Powerformer's* attention score and weight distributions. The reference (dotted line) has no mask (MHA), the solid line uses WCMHA and has the mask applied, and the dashed-dotted line is the WCMHA distribution calculated before applying the mask. This result are for the Weather dataset with a 96 prediction length and 512 input length.

model with extraneous information. Ultimately, we observe very little difference between learnable and constant $\alpha$ (see Table A9). This indicates that WCMHA acts as a regularizer during training, removing superfluous information to avoid overfitting and improve generalization.

### 4.4 Interpretability

The simplicity of *Powerformer* is a key strength and indicator that causal and local biases are important for time-series Transformer-based models. Previous models employ more complicated attention schemes using autocorrelations [42] or Fourier decompositions [48]; or they induce causal and local implicit biases via various mechanisms [39, 20, 19]. Some models, like FEDformer [48], purposely avoid such implicit biases. The success and simplicity of WCMHA and *Powerformer* over other models and their attention techniques bolsters the importance of prioritizing embeddings with causal correlation structures natural to the dataset. Moreover, *Powerformer's* simplicity, along with the local and causal biases, promotes interpretability and provides insight into the dependencies of each forecast.

By leveraging *Powerformer's* simplicity, we interprate it's behavior through the attention weight and score distributions. As the power-law dampening grows (Fig. 4), an emergent bimodal distribution in attention weights appears. The score distributions (inset in Fig. 4) shift towards lower values, while the positive heavy-tail increases as the
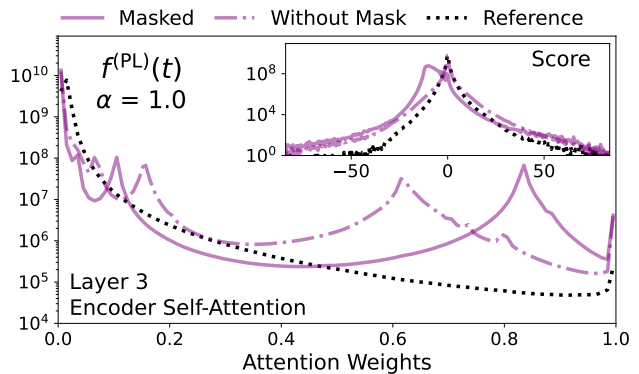
decay constant grows. This shift in scores to larger magnitude values produces an emergent bimodal distribution in the attention weights. This bimodal behavior is well aligned with our implicit biases as the right mode (larger weights) corresponds to our amplified local interactions and the left mode (smaller weights) corresponds to the diminished non-local interactions. The bimodal $\mathbf{C}^{(C,D)}$ distributions in Fig. 4 indicates that $\mathbf{M}^{(C,D)}$ achieves our stated goal of amplifying local couplings, removing causally inconsistent couplings, and diminishing unphysically long-time couplings.

Remarkably, we observe *Powerformer* further enforcing our local bias beyond $\mathbf{M}^{(D)}$, and we consider this a key insight into the local bias' importance. Because our masks are constant, *Powerformer* has the potential to mitigate the effects of $\mathbf{M}^{(D)}$ by learning to reweight the attention scores. Instead, we observe the opposite, *Powerformer* further enforces our locality bias during training by imposing a bimodal attention weight distribution before $\mathbf{M}^{(C,D)}$ is applied. Figure 5 shows the attention score and weight distributions with and without applying $\mathbf{M}^{(C,D)}$. The dashed-dotted line in the attention weight indicates that *Powerformer* imposed a local bias before the locally biasing $\mathbf{M}^{(C,D)}$ mask is applied. This learned effect amplifies the impact of $\mathbf{M}^{(C,D)}$ as seen by the further separation in modes between the solid and dashed-dotted distributions in Fig. 5. Ultimately, *Powerformer* learns embeddings that further strengthen

the effects of our implicit bias. We consider *Powerformer's* learned amplification of our causal and local implicit biases to be strong evidence these biases are crucial in time-series forecasting.

Our local biasing power-law $\mathbf{M}^{(C,D)}$ endows *Powerformer* with an implicit temporal multiscale structure as its embeddings change as the timescale increases. Let us consider the case when $\mathbf{M}^{(D)}$ mask has a long tail (e.g., $f^{(PL)}(t)$ for $\alpha \leq 0.5$ in Fig. 2). The mask weighting in this tail begins to asymptote and similarly scales the attention weights at large $\Delta t$. Summing over many similarly weighted time steps washes out high-frequency information and amplifies dynamic information on the time scale of the summation. For example, if one added a fast oscillating sinusoid to a line and summed this signal over many periods of the sinusoid, the high-frequency oscillations would cancel and one would see a linear signal. Conversely, in temporally local regions the mask quickly changes to amplify local high-frequency components. These local high-frequency components thus play a more prominent role than high-frequency components from longer windows. The mask imposes a temporal multiscale attention structure by capturing low-frequency information from longer-time summations and high frequency from local interactions. In Appendix D, we observe this phenomena as *Powerformer* better predicts higher frequency contributions fast-varying signals than PatchTST.

By removing longer-term pairwise contributions, WCMHA has the potential to improve the attention complexity cost from quadratic to linear. Moreover, since multiple timescales can provide similar results (Tables A5, A6, A7, and A8), such speed improvements may come at a small cost to accuracy. Each mask type and timescale has a cutoff length $\tau$, after which the contributions from data are negligible: e.g., for $f^{(PL)}(t)$ with $\alpha = 1$, points beyond $\tau = 100$ are diminished by two orders of magnitude and may be considered negligible. When WCMHA is only applied within $\tau$, it scales linearly $O(\tau T)$ with input sequence length $T$, instead of quadratically ($O(T^2)$) like MHA. Because $\alpha$ varies between datasets and can be lowered to trade speed for small changes in accuracy, we expect a large variance in potential speedup times, and so we do not provide results here. For example, $f^{(PL)}(t)$ with $\alpha = 1$ provides a factor of 3 and 5 speedup for typical look-back window of 336 and 512, respectively.

WCMHA leverages the strengths of both statistical models and MHA. The induced power-law decay in attention weights, and thus temporal correlations, resembles previous statistical models. However, the WCMHA attention mechanism is not fully overwhelmed by this decay and still amplifies important longer-time correlations, as shown in Fig. 1b and Fig. 1c where the attention mask clearly selects regions of strong correlations near and past the tail of the power-law decay. We believe this balance between power-law temporal correlations and maintaining WCMHA's ability to attend to highly correlated input is key to WCMHA's success.

# 5   Conclusion

We develop *Powerformer* to impose and investigate the importance of causal and local biases for Transformer-based time-series models. *Powerformer* outperforms state-of-the-art time series models on standard tasks often with a simpler architecture. *Powerformer's* simplicity promotes interpretability studies that elucidate the effects of our biases. *Powerformer* strengthens our local bias by learning to decrease attention scores for long-time interactions.

In addition to outperforming previous models on forecasting accuracy, *Powerformer* also has the potential to significantly increase attention efficiency. The decaying nature of $f(t)$ determines a cutoff time after which $\mathbf{M}^{(D)}$ effectively removes contributions. Evaluating within this cutoff time will result in $O(\tau T)$ rather than $O(T^2)$ complexity.

# Acknowledgements

# References

[1] Md Atik Ahamed and Qiang Cheng. Tscmamba: Mamba meets multi-view learning for time series classification, 2024.

[2] Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, Jasper Zschiegner,

Danielle C. Maddix, Hao Wang, Michael W. Mahoney, Kari Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Yuyang Wang. Chronos: Learning the language of time series, 2024.

[3] Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *CoRR*, abs/1803.01271, 2018.

[4] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, S. Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen A. Creel, Jared Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren E. Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas F. Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, O. Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir P. Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Benjamin Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, J. F. Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Robert Reich, Hongyu Ren, Frieda Rong, Yusuf H. Roohani, Camilo Ruiz, Jack Ryan, Christopher R'e, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishna Parasuram Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei A. Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and risks of foundation models. *ArXiv*, 2021.

[5] George E. P. Box and Gwilym M. Jenkins. *Time Series Analysis: Forecasting and Control.* Holden-Day, 1970.

[6] Stephen Butterworth. On the theory of filter amplifiers. *Experimental Wireless and the Wireless Engineer*, 7:536,541, 1930.

[7] Cristian Challu, Kin G. Olivares, Boris N. Oreshkin, Federico Garza Ramirez, Max Mergenthaler Canseco, and Artur Dubrawski. Nhits: Neural hierarchical interpolation for time series forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(6):6989–6997, Jun. 2023.

[8] Si-An Chen, Chun-Liang Li, Sercan O Arik, Nathanael Christian Yoder, and Tomas Pfister. TSMixer: An all-MLP architecture for time series forecast-ing. *Transactions on Machine Learning Research*, 2023.

[9] Aaron Clauset, Cosma Rohilla Shalizi, and M. E. J. Newman. Power-law distributions in empirical data. *SIAM Review*, 51(4):661–703, 2009.

[10] Abhimanyu Das, Weihao Kong, Andrew Leach, Shaan K Mathur, Rajat Sen, and Rose Yu. Long-term forecasting with tiDE: Timeseries dense encoder. *Transactions on Machine Learning Research*, 2023.

[11] Carson Eisenach, Yagna Patel, and Dhruv Madeka. Mqtransformer: Multi-horizon forecasts with context dependent and feedback-aware attention, 2022.

[12] Jean-Yves Franceschi, Aymeric Dieuleveut, and Martin Jaggi. Unsupervised scalable representation learning for multivariate time series. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

[13] Azul Garza, Cristian Challu, and Max Mergenthaler-Canseco. Timegpt-1, 2024.

[14] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024.

[15] Yangdong He and Jiabao Zhao. Temporal convolutional networks for anomaly detection

in time series. *Journal of Physics: Conference Series*, 1213(4):042050, jun 2019.

[16] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, 11 1997.

[17] Charles C. Holt. *Forecasting Seasonals and Trends by Exponentially Weighted Moving Averages*. O.N.R. research memorandum. Defense Technical Information Center, 1957.

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

[19] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020.

[20] 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. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché Buc, Edward A. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems*, pages 5244–5254, 2019.

[21] Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting long-term time series forecasting: An investigation on linear mapping, 2023.

[22] Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. SCINet: Time series modeling and forecasting with sample convolution and interaction. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.

[23] Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and forecasting. In *International Conference on Learning Representations*, 2022.

[24] Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024.

[25] Yong Liu, Guo Qin, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. Autotimes: Autoregressive time series forecasters via large language models, 2024.

[26] Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 9881–9893. Curran Associates, Inc., 2022.

[27] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2023.

[28] Boris N. Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis expansion analysis for interpretable time series forecasting. In *International Conference on Learning Representations*, 2020.

[29] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners, 2019.

[30] David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.

[31] Lifeng Shen, Zhuocong Li, and James Kwok. Timeseries anomaly detection using temporal hierarchical one-class network. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 13016–13026. Curran Associates, Inc., 2020.

[32] Sabera Talukder, Yisong Yue, and Georgia Gkioxari. Totem: Tokenized time series embeddings for general time series analysis, 2024.

[33] Sean J. Taylor and Benjamin Letham. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.

[34] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alexander Graves, Nal Kalchbrenner, Andrew Senior,

and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. In *Arxiv*, 2016.

[35] Aaron van den Oord, Oriol Vinyals, and koray kavukcuoglu. Neural discrete representation learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[37] Zihan Wang, Fanheng Kong, Shi Feng, Ming Wang, Han Zhao, Daling Wang, and Yifei Zhang. Is mamba effective for time series forecasting? *ArXiv*, abs/2403.11144, 2024.

[38] Peter R. Winters. Forecasting sales by exponentially weighted moving averages. *Management Science*, 6(3):324–342, 1960.

[39] Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. Etsformer: Exponential smoothing transformers for time-series forecasting, 2022.

[40] Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *The Eleventh International Conference on Learning Representations*, 2023.

[41] Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Flowformer: Linearizing transformers with conservation flows. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 24226–24242. PMLR, 17–23 Jul 2022.

[42] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.

[43] Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Anomaly transformer: Time series anomaly detection with association discrepancy. In *International Conference on Learning Representations*, 2022.

[44] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(9):11121–11128, Jun. 2023.

[45] Michael Zhang, Khaled Kamal Saab, Michael Poli, Tri Dao, Karan Goel, and Christopher Re. Effectively modeling time series with simple discrete state spaces. In *The Eleventh International Conference on Learning Representations*, 2023.

[46] Tianping Zhang, Yizhuo Zhang, Wei Cao, Jiang Bian, Xiaohan Yi, Shun Zheng, and Jian Li. Less is more: Fast multivariate time series forecasting with light sampling-oriented mlp structures, 2022.

[47] 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. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12):11106–11115, May 2021.

[48] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 27268–27286. PMLR, 17–23 Jul 2022.

[49] Tian Zhou, Peisong Niu, Xue Wang, Liang Sun, and Rong Jin. One fits all: Power general time series analysis by pretrained LM. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

# Appendix

## A    Butterworth Filter

The Butterworth filter [6] is often used in signal processing for low-, high-, and band-pass filters. It is designed to be optimally flat within the passband and sometimes referred to as the "maximally flat filter." The motivation behind it's derivation is to have equal sensitivity to frequency modes within the passband. We leveraged the Butterworth filter to equally weight temporal measurements within the lookback window. The analog Butterworth filter gain is given by

$$f_n^{\mathrm{BF}}(z) = \frac{-1}{\sqrt{1 + \left(\frac{z}{t_c}\right)^{2n}}}, \tag{A1}$$

where $t_c$ is the critical time that sets the width of the filter and the order $n$ determines how fast the gain decays after $t_c$. In this work we use the digital Butterworth filter gain. Figure A1 shows the $f_1^{(\mathrm{BW})}(t)$ and $f_2^{(\mathrm{BW})}(t)$ contributions on the attention score and weights for varying decay constants.



Figure A1: We show the effects of the Butterworth filter masks $\left(f_n^{(\mathrm{BW})}(t)\right)$ for orders 1 (dashed lines) and 2 (solid lines), while varying the critical time. Panel (a) shows the effects on the attention scores and Panel (b) shows the corresponding effects on the attention weights after applying Softmax to $f_n^{(\mathrm{BW})}(t)$.

### A.1    Digital Filter Gain

To transform the analog Butterworth filter into a digital filter, the filter design must be discretized. This discretization is often done using the bilinear transform

$$z = \frac{2}{T}\frac{z-1}{z+1}, \tag{A2}$$

where $T$ is the numerical integration step size used in the trapezoidal rule.

### A.2    Code to Calculate the Gain

For reproducibility, we provide the code used to calculate the digital Butterworth filter gains. We multiply the gain by 5 to scale with the attention key-query overlap values.

```
import numpy as np
import scipy as sp
```

```
def butterworth_filter(scale, order, times):
    b, a = sp.signal.butter(order, 0.8, 'lowpass', analog=False)
    t, decay = sp.signal.freqz(b, a)
    t = scale*t/2
    dc = 5*np.log(np.abs(decay))
    decay_interp = sp.interpolate.interp1d(t, dc)
    return decay_interp(times)
```

# B  Datasets

We evaluate *Powerformer* on 7 real-world datasets that have become standard public benchmarks for time-series forecast tasks. The datasets can be found and downloaded in Ref. [42], or individually at references below.

Table A1: We provide the number of variables and the number of timesteps for each dataset.

| Datasets | Illness | ETTh* | ETTm* | Weather | Electricity | Traffic |
|---|---|---|---|---|---|---|
| Features | 7 | 7 | 7 | 21 | 321 | 862 |
| Timesteps | 966 | 17420 | 69680 | 52696 | 26304 | 17544 |

**ETT**[1] [47] provides 7 measurements of the electrical transformer power load (including oil temperature) between July 2016 and July 2018. This is done over 2 regions in China at 2 different sampling rates (hourly and every 15 minutes), resulting in 4 separate datasets (ETTh1, ETTh2, ETTm1, ETTm2). We show the pairwise correlation dependence in Fig. A5.

**Weather**[2] [42] provides 21 meteorological measurements (air temperature, air pressure, humidity, precipitation, etc.) collected in Germany over the whole of 2020. These measurements are recorded every 10 minutes. We show the pairwise correlation dependence in Fig. A4.

**Electricity**[3] [42] provides the electricity consumption (kWh) of 321 consumers from 2012 to 2014. These measurements are sampled every hour. We show the pairwise correlation dependence in Fig. A3.

**Traffic**[4] [42] provides occupancy rates on San Francisco Bay Area freeways from 826 sensors. This data comes from the California Department of Transportation and is sampled hourly. We show the pairwise correlation dependence in Fig. A2.

# C  Architecture Details

## C.1  Transformer

We use a vanilla encoder-decoder Transformer architecture for all Transformer experiments. The decoder input is the last 48 time-steps in the sequence concatenated by an array of zeros with a length of the prediction window. For our baseline experiment, the decoder self-attention has a causal mask, which is common in Transformer architectures, while the encoder self-attention and decoder cross-attention do not have masks. Our WCMHA variant of the Transformer applies WCMHA to both the encoder and decoder self-attention, but not the decoder cross-attention. We do not apply WCMHA to the decoder cross-attention because this process is already causal as it forecasts future points based on the given sequence. We present the Transformer hyperparameters used in this work in Table A2.
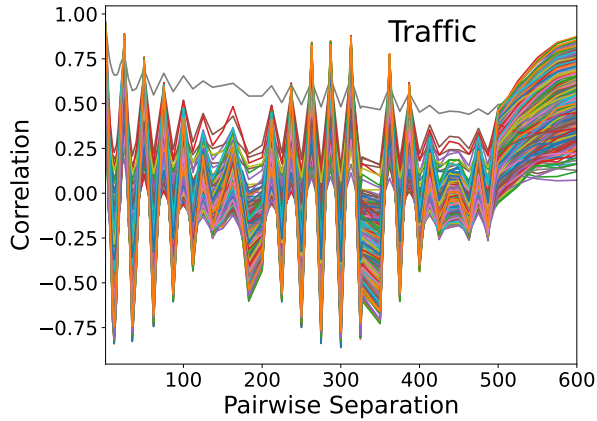
---

[1] https://github.com/zhouhaoyi/ETDataset
[2] https://www.bg-ena.mpg.de/wetter/
[3] https://archive.ics.uci.edu/ dataset/321/electricityloaddiagrams20112014
[4] http://pems.dot.ca.gov

Figure A2: We show the Weather dataset's pairwise correlations as a pairwise separation function. Each colored line represents a different variable in the dataset.

## C.2 PatchTST

Here, we review PatchTST's patching mechanism, as PatchTST is otherwise similar to a traditional Transformer with a linear readout head that simultaneously predicts all future points. For a detailed description of PatchTST, see Ref. [27]. Let us consider an input time-series $\mathbf{X}$ with dimensions $[B, T_{\text{seq}}, D]$, where $B$ is batch size, $T_{\text{seq}}$ is the time interval, and $D$ is the number of input dimensions. The multivariate input is first normalized to zero mean and unit variance along the time axis and flattened into $D$ univariate inputs, resulting in a dimensionality of $[B \times D, T_{\text{seq}}]$. PatchTST batches these time series via a convolution with $N$ filters, a stride $s$, and patch size $p$, resulting in a dimensionality of $[B \times D, P, N]$. Here, $P = (T_{\text{seq}} - p)//s + 1$ is the number of patches, and $//$ is integer division. At this point, the patched data is fed into the transformer where a positional embedding is added and the patched embeddings are further contextualized by vanilla MHA [36]. The linear output head takes the $[B \times D, P, N]$ output of the Transformer, flattens the last two dimensions, and applies a linear matrix multiplication (of shape $[P \times N, T_{\text{pred}}]$) to predict the output time series of length $T_{\text{pred}}$. The linear head returns a matrix of shape $[B \times D, T_{\text{pred}}]$. This matrix is reshaped into the original multivariate dimensionality ($[B, T_{\text{pred}}, D]$) and renormalized. We note that the patching mechanism is performed before applying the Transformer and is independent of it.
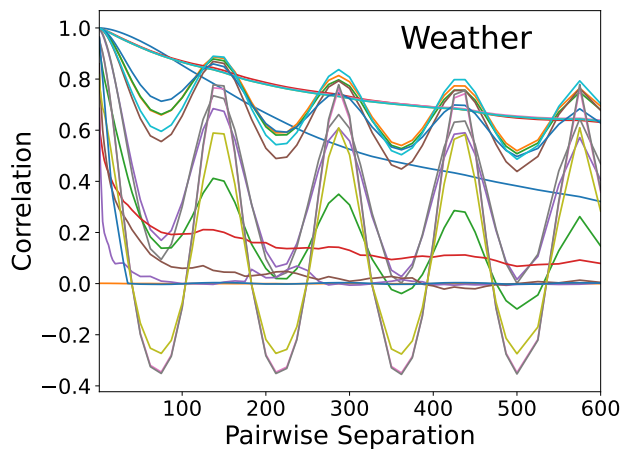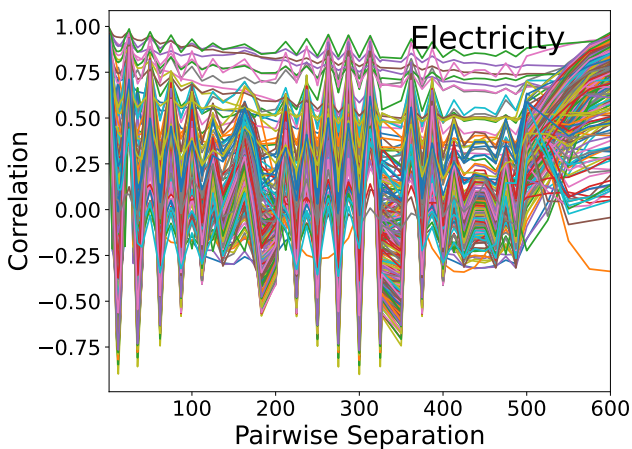


Figure A3: We show the Electricity dataset's pair- Figure A4: We show the Weather dataset's pairwise correlations as a pairwise separation function. correlations as a pairwise separation function. Each Each colored line represents a different variable in colored line represents a different variable in the the dataset. dataset.
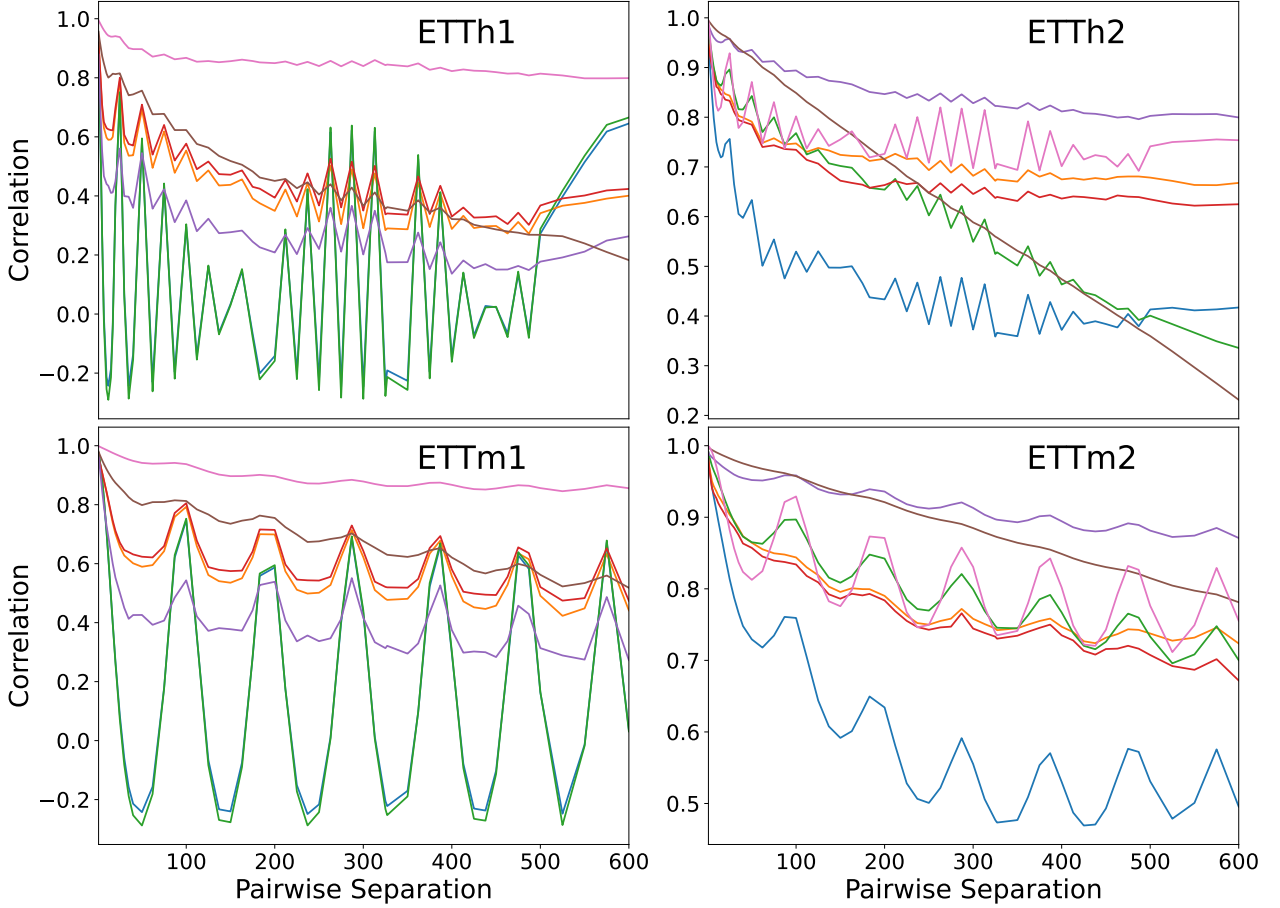
3

Figure A5: We show the ETT datasets' pairwise correlations as a pairwise separation function. Each colored line represents a different variable in the dataset.

## C.3 Powerformer

*Powerformer* forecasts univariate time-series by combining WCMHA, the commonly used encoder-only Transformer architecture, and patching [27]. We selected PatchTST as the base of *Powerformer* since PatchTST already includes decomposes multivariate time-series into univariate channels and performs the patching. Additionally, similar to other modern time-series Transformer-based models [49, 24, 32], PatchTST [27] uses an encoder-only architecture with a linear readout head. *Powerformer* uses a standard Transformer encoder architecture with the primary difference being the replacement of the vanilla MHA by weighted causal WCMHA, described in Eqs. 6-10. We additionally remove the dropout applied to the attention weights because our implicit biases enforce a deterministic temporal dependency on pairwise correlations. Dropout enforces redundancy between pairwise correlations and is therefore in conflict with our implicit biases. The input data comes in with $D$ variates and shape $[B, D, T_{\text{seq}}]$, where $B$ is the batch size and $T_{\text{seq}}$ is the input sequence length. The PatchTST patching mechanism (described in Appendix C.2) returns the patched and normalized univariate embeddings with shape $[B \times D, P, N]$, where $P$ is the number of patches and $N$ is the size of the embeddings. These embeddings are encoded through multiple WCMHA Transformer encoder layers (see Appendix C.4 for hyperparameter details). The encoder output (of shape $[B \times D, P, N]$) is flattened along the last two dimensions ($[B \times D, P \times N]$) in preparation for the linear readout head. The linear readout matrix has the shape $[P \times N, T_{\text{pred}}]$ and after acting on the flattened encoder output produces the univariate forecasts of shape $[B \times D, T_{\text{pred}}]$. Finally, these univariate forecasts are renormalized and reshaped into the multivariate input structure $[B, D, T_{\text{pred}}]$ and returned to the user. Consequently, *Powerformer* is a simple Transformer-based model with causal and local masks added to the attention scores $\mathbf{S}$, and patching applied to the input data. We present the *Powerformer* hyperparameters used in this work in Table A2.

4

## C.4   Model Tuning Parameters

To reproduce our results, we list our hyperparameters used on each dataset for *Powerformer* (Table A2) and Transformer (Table A3). For each dataset and model, we trained with the following 3 random seeds: 2021, 1776, and 1953.

Table A2: We provide *Powerformer* architecture and training parameters used for each dataset.

| Parameters | ETTh1 | ETTh2 | ETTm1 | ETTm2 | Weather | Electricity | Traffic |
|---|---|---|---|---|---|---|---|
| Sequence Length | 336 / 512 | 336 / 512 | 336 / 512 | 336 / 512 | 336 / 512 | 336 / 512 | 336 / 512 |
| Patch Length | 16 | 16 | 16 | 16 | 16 | 16 | 16 |
| Patch Stride | 8 | 8 | 8 | 8 | 8 | 8 | 8 |
| Encoder Layers | 3 | 3 | 3 | 3 | 3 | 3 | 3 |
| Embedding Size | 16 | 16 | 128 | 128 | 128 | 128 | 128 |
| MHA Heads | 4 | 4 | 16 | 16 | 16 | 16 | 16 |
| MHA Feed Forward | 128 | 128 | 256 | 256 | 256 | 256 | 256 |
| Dropout % | 30 | 30 | 20 | 20 | 20 | 20 | 20 |
| Linear Dropout % | 30 | 30 | 20 | 20 | 20 | 20 | 20 |
| Train Epochs | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| Patience | – | – | 20 | 20 | 20 | 10 | 10 |
| Learning Rate | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ |
| Batch Size | 128 | 128 | 128 | 128 | 128 | 32 | 24 |

Table A3: We provide Transformer architecture and training parameters used for each dataset.

| Parameters | ETTh1 | ETTh2 | ETTm1 | ETTm2 | Weather | Electricity | Traffic |
|---|---|---|---|---|---|---|---|
| Sequence Length | 256 | 256 | 256 | 256 | 256 | 256 | 256 |
| Decoder Input Length | 48 | 48 | 48 | 48 | 48 | 48 | 48 |
| Encoder Layers | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| Decoder Layers | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Embedding Size | 512 | 512 | 512 | 512 | 512 | 512 | 512 |
| MHA Heads | 8 | 8 | 8 | 8 | 8 | 8 | 8 |
| MHA Feed Forward | 2048 | 2048 | 2048 | 2048 | 2048 | 2048 | 2048 |
| Dropout % | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| Linear Dropout % | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| Train Epochs | 100 | 100 | 100 | 100 | 100 | 100 | 100 |
| Patience | – | – | – | – | – | – | – |
| Learning Rate | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ | $10^{-4}$ |
| Batch Size | 128 | 128 | 128 | 128 | 128 | 128 | 128 |

# D   Forecasting Visualization

Here, we show *Powerformer* forecasting visualizations and compare various decay time constants against ground truth and PatcTST [27]. We also highlight plots where both *Powerformer* and PatchTST struggle. In these cases *Powerformer* generally outperforms better capture faster-moving features, non-periodic features, and is better able to alter periodic predictions. We believe this is due to *Powerformer's* ability to focus on and amplify local high-frequency features while dampening or removing long-time high-frequency components. WCMHA achieves this through the power-law masks $f^{(\text{PL})}(t)$ and $f^{(\text{SPL})}(t)$, which amplify local correlations and dampen long-time correlations.

# E Experiments

## E.1 Experiment Details

All experiments are evaluated 3 times with different initialization (C.4). For each dataset, we select a single input sequence based on the best-performing input sequence length across all prediction lengths. We choose a single input sequence length for each dataset because the strength of the correlations between time steps is independent of the prediction length, but the required input sequence length to capture the information necessary to forecast is unique for each dataset. For each prediction length, we treat the attention mask as a hyperparameter and select the best mask and time constant pair. We select attention masks for each prediction length because the forecasting timescale will likely emphasize short or long-term correlations for shorter or longer forecasting windows. We make our selections based on the MSE error.

## E.2 Powerformer

We evaluate the weight power-law ($f^{(\mathrm{PL})}(t)$) and similarity power-law ($f^{(\mathrm{SPL})}(t)$) filters on all 7 datasets. For $f^{(\mathrm{PL})}(t)$ we evaluate $\alpha \in \{0.1, 0.25, 0.5, 0.75, 1.0\}$ and for $f^{(\mathrm{SPL})}(t)$ we evaluate $\alpha \in \{0.1, 0.5, 1.0, 2.0\}$. Similar to PatchTST [27] we evaluate the power-law filters with two input lengths: 336 and 512.

We evaluate both Butterworth filters $\left(f_1^{(\mathrm{BW})}(t) \text{ and } f_2^{(\mathrm{BW})}(t)\right)$ on all the ETT datasets and the Weather datasets. For each Butterworth filter, we evaluated $\alpha \in \{2, 5, 10, 15, 20\}$ with an input length of 336. Due to poor performance, we did not evaluate the Butterworth filters on the larger Electricity and Traffic datasets, nor with the longer 512 input length.



Figure A6: We show forecasting results on the Electricity dataset for *Powerformer* with $f^{(\mathrm{PL})}(t)$ and PatchTST. The colored lines represent PatchTST and different $f^{(\mathrm{PL})}(t)$ decay constants. For these forecasts, the sequence length is 512 and the prediction length is 336.

## E.3 Transformer

We train a standard encoder-decoder Transformer architecture both with and without WCMHA. Generic Transformer models often incorporate causality through masking. To account for this, our Transformer benchmark uses masking for the decoder MHA but not for the encoder MHA nor the decoder cross-attention. We investigate WCMHA's effects by replacing the encoder and decoder self-attention, but not the decoder cross-attention as it is already causal. Section C.1 provides more architectural details.

We train Transformer models, with and without WCMHA, using the weight power-law ($f^{(\mathrm{PL})}(t)$) filter on all 7 datasets. We explore the following decay constants: $\alpha \in [0.1, 0.25, 0.5, 0.75, 1.0]$ and provide further experimental detail in sections 4 and E.1.

## E.4 Baseline Selection

We compare *Powerformer* against multiple Transformer-based state-of-the-art methods. We populate Table 1 with values taken from the models' respective papers. All models except TOTEM [32] and iTransformer [24] treat the input length as a hyperparameter and report the best results. PatchTST [27] provides results for two input sequence lengths (336 and 512). In Table 1 we report the best result between these input sequence lengths for each dataset based on the MSE. We employ the same input sequence length selection for *Powerformer*.
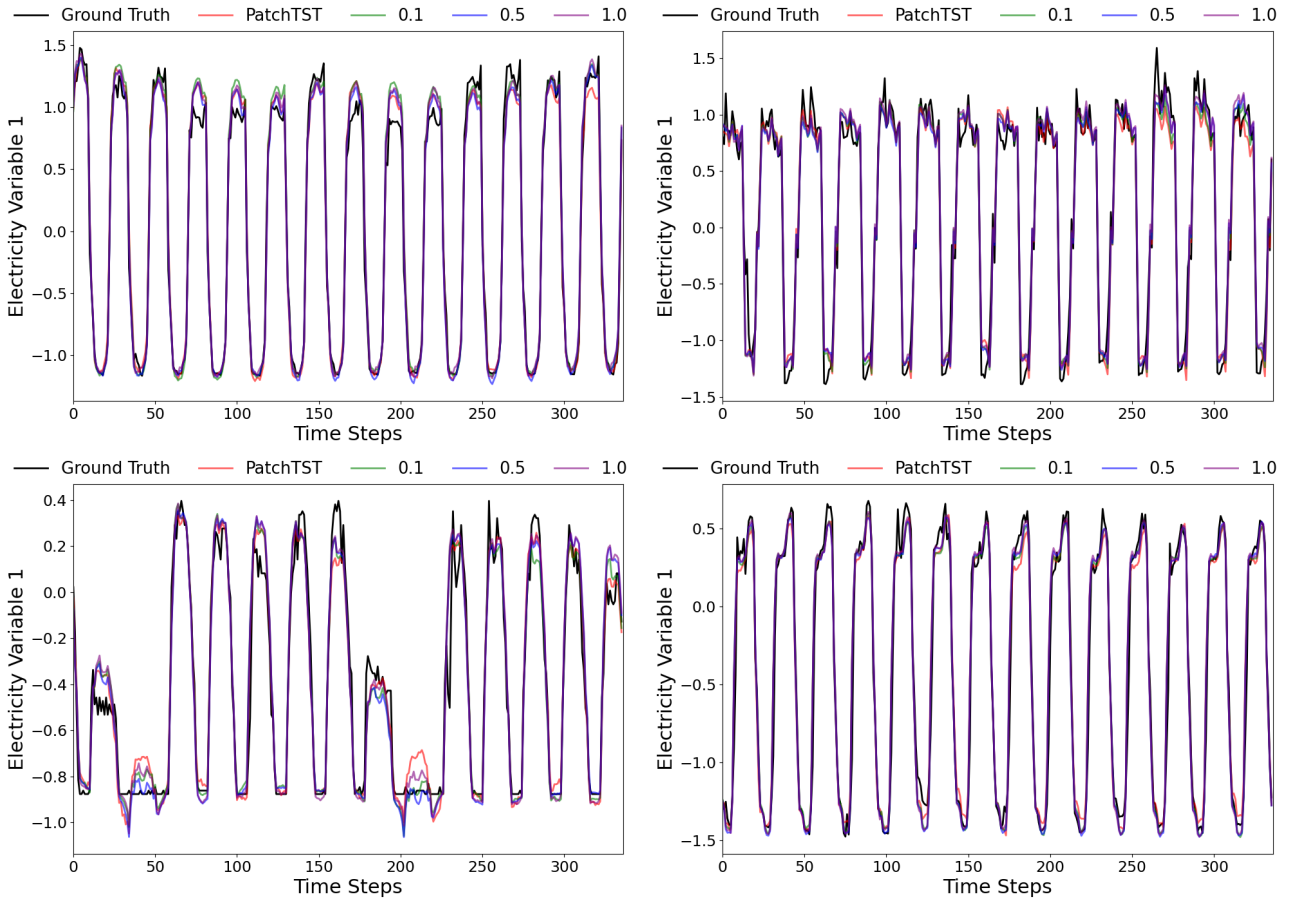


Figure A7: We show forecasting results on the ETTm1 dataset for *Powerformer* with $f^{(\mathrm{PL})}(t)$ and PatchTST. The colored lines represent PatchTST and different $f^{(\mathrm{PL})}(t)$ decay constants. For these forecasts, the sequence length is 512 and the prediction length is 336.

# F    Ablation Studies

Here, we perform ablation experiments on WCMHA in both Transformer and *Powerformer*. For both models, we examine the effects of different masks ($f^{(PL)}(t)$, $f^{(SPL)}(t)$, and $f_n^{(BW)}(t)$) and various time decays $\alpha$. We do not perform the same amount of experiments for each model and mask pairing due to poor performance. In the remainder of this section, we provide detailed tables of the MSE and MAE scores, as well as plots of the attention scores and weights. In the last subsection (F.3), we investigate the effects of making $\alpha$ a learnable parameter.

## F.1    Transformer Results

Here, we present the full results from evaluating Transformer with WCMHA. The architecture is described in Appendix C.1 and the experimental details are provided in Appendices E.1 and E.3.

Table A4 provides the aggregate Transformer performance without WCMHA and with WCMHA for varying decay times ($\alpha$). WCMHA decisively outperforms MHA on 4 or the 7 datasets and ties with the number of best results on ETTh2. Looking further into ETTh2 we see that when WCMHA outperforms MHA it does so by much larger margins (16%) than when it underperforms (2%). Moreover, we observe the same larger margins of overperformance across all of the datasets. This indicates that when the natural pairwise distribution is found we can see significant improvements in MSE and MAE, but the effects of imposing the wrong pairwise distribution are much less significant.

In Figs. A9-A22 we present the attention score and weight distributions with and without applying the local-causal mask. The distributions can have significant deviations between datasets. In general, the WCMHA decoder self-attention shows the smallest deviation from the MHA base case. We note
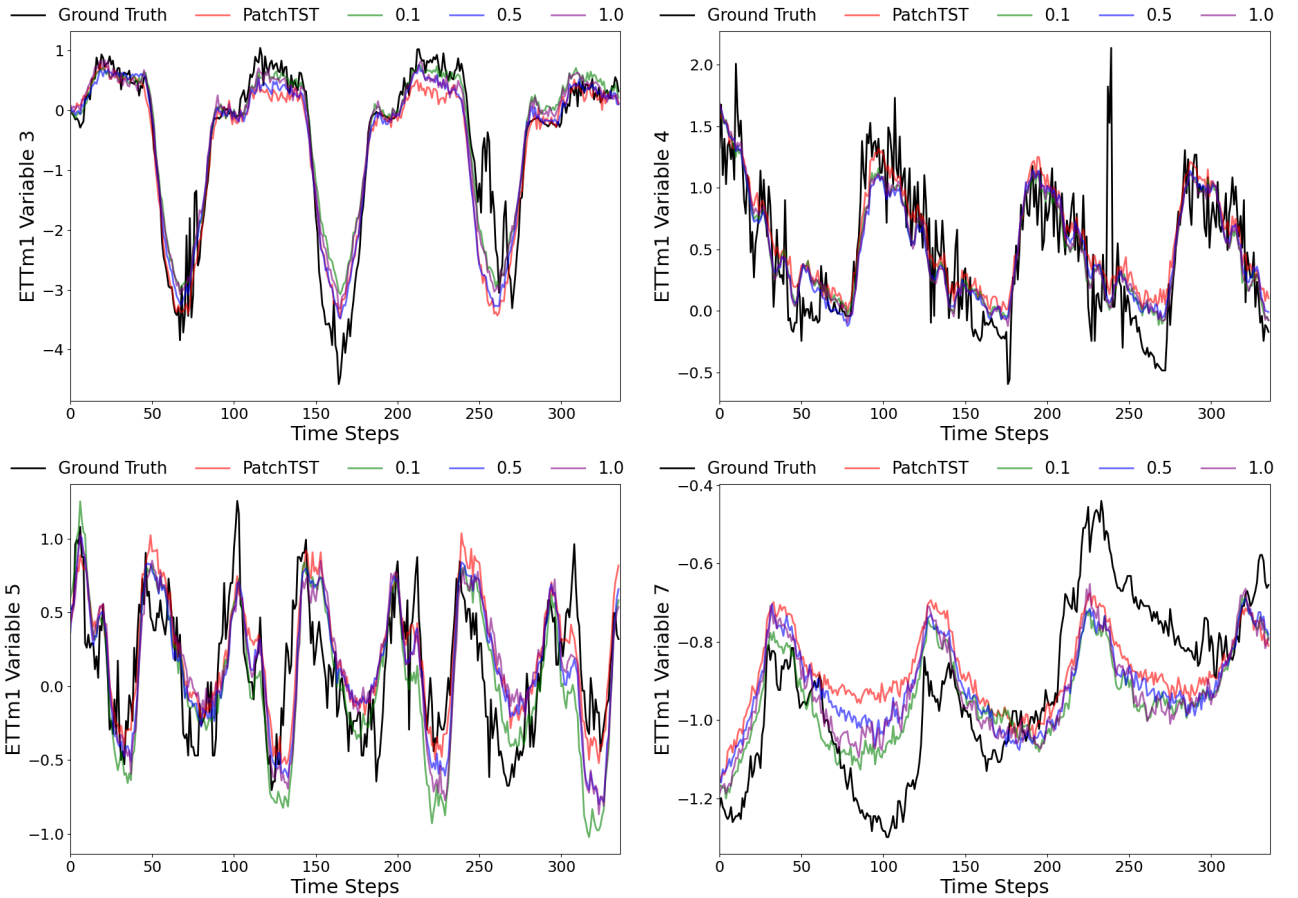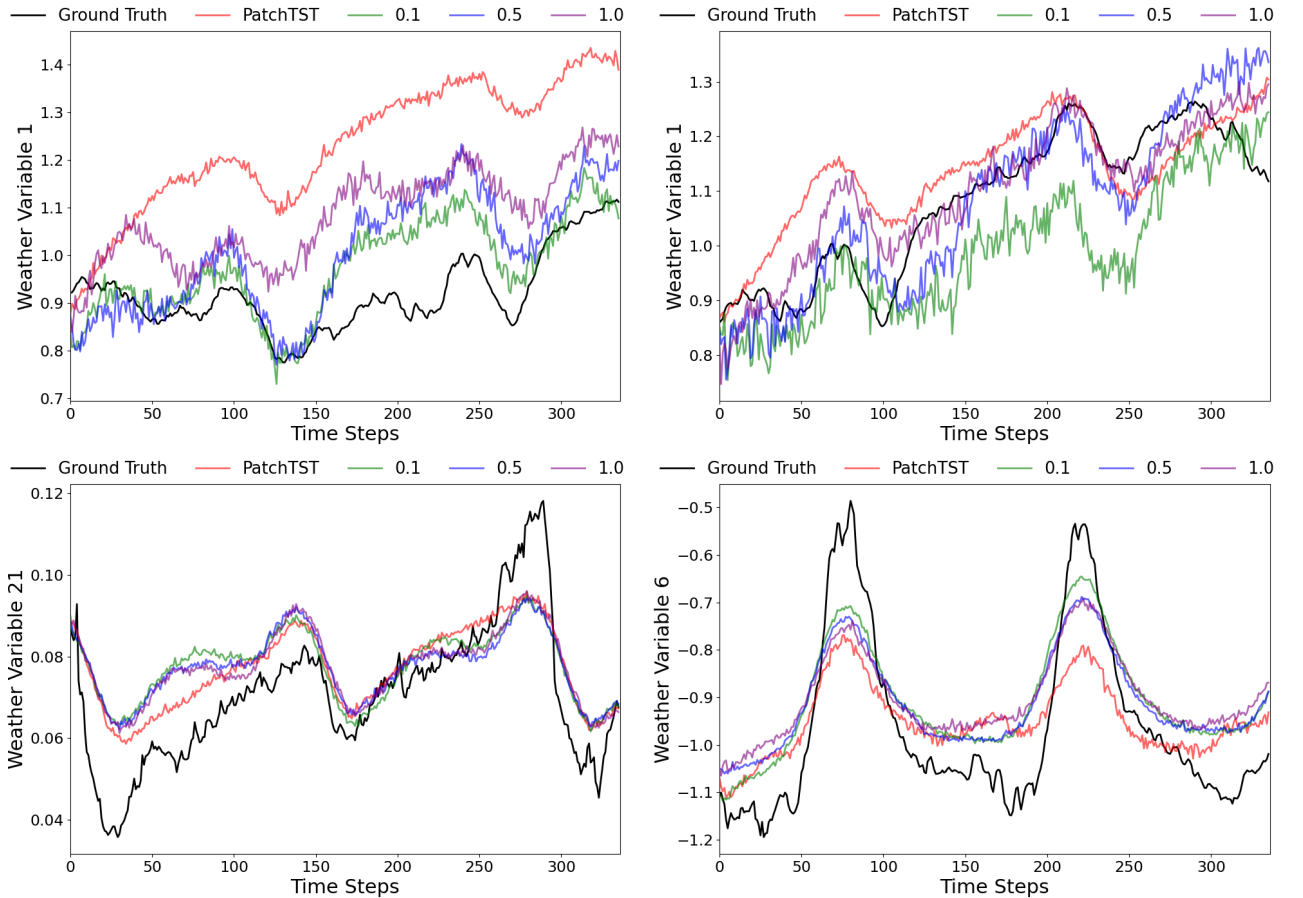


Figure A8: We show forecasting results on the Weather dataset for *Powerformer* with $f^{(PL)}(t)$ and PatchTST. The colored lines represent PatchTST and different $f^{(PL)}(t)$ decay constants. For these forecasts, the sequence length is 512 and the prediction length is 336.

that the MHA base case already has a causal mask. The last WCMHA encoder self-attention layer often shows the largest difference between WCMHA and MHA distributions. This makes sense since the MHA encoder attention is not causal. We do not alter the decoder cross-attention, but we still observe shifts in the distribution, which differ between datasets. Interestingly we do not observe the same bimodal distribution in attention weights as we do for *Powerformer*. This may be due to the Transformer's encoder-decoder structure, whereas *Powerformer* only has an encoder.

Table A4: We provide the forecasting MSE and MAE on the test sets for Transformer with MHA (base case) and for Transformer with WCMHA and the weight power-law mask $f^{(\text{PL})}(t)$. The top row indicates the model or decay constant ($\alpha$), the bold numbers are the best (lowest) performance, and the underlined numbers are the second best.

| | | Base Case | | 0.1 | | 0.25 | | 0.5 | | 0.75 | | 1 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metric | | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 96 | 0.936 | 0.741 | 0.955 | 0.752 | 0.944 | 0.745 | 0.921 | 0.729 | <u>0.896</u> | <u>0.711</u> | **0.874** | **0.693** |
| | 192 | 0.988 | 0.769 | 0.916 | 0.759 | 0.951 | 0.789 | <u>0.887</u> | <u>0.749</u> | 0.895 | **0.745** | **0.876** | <u>0.749</u> |
| | 336 | 1.061 | 0.815 | 1.067 | 0.835 | 1.063 | 0.814 | 1.054 | <u>0.809</u> | **1.038** | 0.836 | <u>1.044</u> | **0.802** |
| | 720 | 1.113 | 0.881 | **1.047** | **0.842** | <u>1.061</u> | <u>0.855</u> | 1.180 | 0.893 | 1.185 | 0.895 | 1.259 | 0.899 |
| ETTh2 | 96 | 1.585 | 1.038 | **1.344** | **0.918** | <u>1.349</u> | <u>0.920</u> | 1.392 | 0.936 | 1.496 | 1.000 | 1.511 | 1.005 |
| | 192 | **1.889** | **1.131** | <u>1.899</u> | <u>1.133</u> | 3.025 | 1.459 | 2.916 | 1.423 | 2.800 | 1.384 | 2.761 | 1.365 |
| | 336 | **2.668** | **1.384** | 2.875 | 1.393 | 2.872 | <u>1.386</u> | 2.879 | 1.412 | <u>2.860</u> | 1.409 | 2.951 | 1.442 |
| | 720 | 3.143 | 1.528 | 3.389 | 1.570 | 3.455 | 1.582 | 3.533 | 1.621 | <u>2.635</u> | <u>1.331</u> | **2.488** | **1.272** |
| ETTm1 | 96 | 0.560 | 0.545 | 0.566 | 0.533 | 0.592 | 0.545 | 0.563 | 0.541 | <u>0.535</u> | <u>0.522</u> | **0.487** | **0.492** |
| | 192 | 0.874 | 0.729 | 0.907 | 0.724 | 0.838 | 0.711 | 0.830 | 0.698 | <u>0.774</u> | <u>0.686</u> | **0.718** | **0.645** |
| | 336 | 0.960 | 0.786 | 0.948 | 0.786 | <u>0.938</u> | <u>0.779</u> | 0.956 | 0.786 | 1.029 | 0.800 | **0.888** | **0.752** |
| | 720 | 0.900 | 0.746 | **0.864** | **0.710** | <u>0.886</u> | <u>0.729</u> | 0.902 | 0.749 | 0.892 | <u>0.729</u> | 1.005 | 0.788 |
| ETTm2 | 96 | **0.597** | **0.614** | <u>0.619</u> | <u>0.628</u> | 0.632 | 0.636 | 0.660 | 0.653 | 0.697 | 0.674 | 0.700 | 0.638 |
| | 192 | **0.849** | **0.749** | <u>0.876</u> | <u>0.761</u> | 0.893 | 0.769 | 0.927 | 0.787 | 0.970 | 0.808 | 1.011 | 0.827 |
| | 336 | **1.227** | **0.909** | <u>1.259</u> | <u>0.920</u> | 1.271 | 0.925 | 1.295 | 0.936 | 1.320 | 0.947 | 1.333 | 0.955 |
| | 720 | **2.176** | **1.237** | 2.214 | <u>1.241</u> | 2.213 | <u>1.241</u> | 2.216 | 1.252 | 2.209 | 1.252 | <u>2.195</u> | 1.250 |
| Weather | 96 | 0.250 | 0.333 | 0.211 | <u>0.297</u> | <u>0.209</u> | 0.302 | **0.207** | **0.291** | 0.216 | 0.305 | 0.228 | 0.316 |
| | 192 | 0.344 | 0.414 | <u>0.295</u> | <u>0.373</u> | **0.287** | **0.367** | 0.297 | 0.375 | 0.307 | 0.383 | 0.313 | 0.390 |
| | 336 | 0.435 | 0.467 | <u>0.403</u> | <u>0.446</u> | 0.416 | 0.455 | 0.420 | 0.458 | 0.410 | 0.452 | **0.384** | **0.436** |
| | 720 | 0.463 | 0.489 | <u>0.420</u> | <u>0.460</u> | 0.435 | 0.471 | 0.424 | 0.464 | 0.426 | 0.468 | **0.403** | **0.454** |
| Electricity | 96 | <u>0.265</u> | 0.360 | 0.271 | 0.364 | **0.264** | **0.355** | 0.266 | <u>0.359</u> | 0.272 | 0.367 | 0.280 | 0.370 |
| | 192 | 0.284 | 0.378 | 0.278 | 0.372 | 0.286 | 0.378 | **0.270** | **0.366** | <u>0.276</u> | <u>0.369</u> | 0.291 | 0.386 |
| | 336 | <u>0.280</u> | 0.371 | 0.289 | 0.381 | 0.290 | 0.381 | 0.288 | 0.381 | 0.281 | <u>0.369</u> | **0.277** | **0.366** |
| | 720 | 0.300 | 0.382 | 0.318 | 0.395 | 0.308 | 0.389 | <u>0.296</u> | <u>0.378</u> | 0.302 | 0.383 | **0.289** | **0.373** |
| Traffic | 96 | 0.662 | 0.365 | 0.656 | 0.359 | 0.653 | <u>0.357</u> | <u>0.652</u> | **0.355** | 0.669 | 0.364 | **0.649** | 0.358 |
| | 192 | **0.655** | **0.350** | 0.675 | 0.368 | 0.669 | 0.361 | 0.669 | 0.365 | <u>0.660</u> | 0.356 | **0.655** | <u>0.355</u> |
| | 336 | **0.651** | **0.350** | <u>0.652</u> | <u>0.353</u> | 0.657 | 0.357 | 0.670 | 0.364 | 0.667 | 0.361 | 0.658 | 0.358 |
| | 720 | **0.660** | **0.354** | 0.669 | 0.363 | 0.676 | 0.360 | <u>0.667</u> | <u>0.355</u> | 0.668 | 0.361 | 0.689 | 0.370 |

## F.2 Powerformer

We present the full *Powerformer* results on all the evaluated masks: $f^{(\text{PL})}(t)$, $f^{(\text{SPL})}(t)$, $f_1^{(\text{BW})}(t)$, and $f_2^{(\text{BW})}(t)$. Architecture details can be found in Appendix C.3, while Appendices E.1 and E.2 outline the experimental details.

### F.2.1 Weight Power-Law Mask

Table A5 presents the aggregate MSE and MAE results for weight power-law mask ($f^{(\text{PL})}(t)$) with varying time decays ($\alpha$), forecasting lengths, and input sequence lengths. The 512 input sequence generally outperforms the shorter 336 input length. As expected, the best performing $\alpha$ varies between datasets, but it also varies as a function of forecast length and input sequence length. Given that some

Figure A9: We present the Transformer attention score (inset) and weight distributions on the ETTh1 dataset with a forecasting length of 96. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\text{PL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{PL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.

of the best-performing decay lengths are $\alpha = 1$ with further experiments one may improve upon these values. However, our $f^{(\text{SPL})}(t)$ mask includes faster-decaying weights, as shown in Fig. 2.

We compare all datasets' attention score and weight distributions both with MHA and WCMHA in Figs A23-A36. We observe similar bimodal distributions for all datasets and note wider $\mathbf{C}^{(\text{C,D})}$ distributions for ETTm1, ETTm2, and Weather.

### F.2.2 Similarity Power-Law Mask

Table A6 presents the aggregate MSE and MAE results for the similarity power-law mask ($f^{(\text{SPL})}(t)$) with varying time decays ($\alpha$), forecast lengths, and input sequence lengths. The 512 input sequence generally outperforms the shorter 336 input length. Compared to $f^{(\text{PL})}(t)$, there are more instances in which a 336 input sequence length outperforms. This is likely due to the much faster decay of $f^{(\text{SPL})}(t)$ compared to $f^{(\text{PL})}(t)$. As expected, the best performing $\alpha$ varies between datasets. However, there is much less variation as a function of forecast length and input sequence length when compared to $f^{(\text{PL})}(t)$.

We compare all datasets' attention score and weight distributions for weighted causal attention in Figs A37-A50. Due to the faster decay of $f^{(\text{SPL})}(t)$ we observe stronger bimodal distributions for all datasets than for $f^{(\text{PL})}(t)$, with wider $\mathbf{C}^{(\text{C,D})}$ distributions for ETTm1, ETTm2, and Weather.
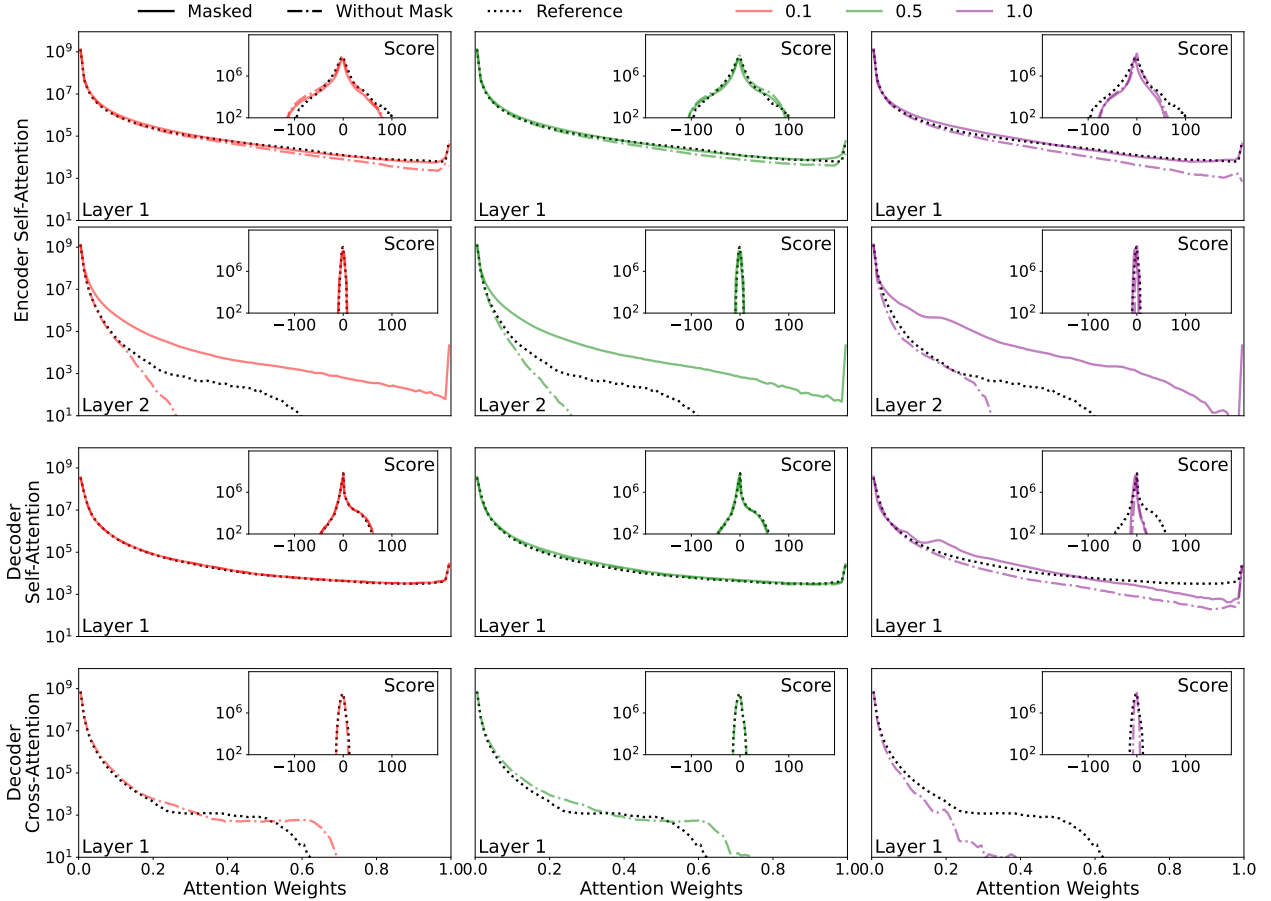
Figure A10: We present the Transformer attention score (inset) and weight distributions on the ETTh1 dataset with a forecasting length of 720. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

### F.2.3 Butterworth Filter: Order 1

Here we evaluate the Butterworth order 1 mask ($f_1^{(\mathrm{BW})}(t)$) on the ETT and Weather datasets with a 336 look-back window. The $f_1^{(\mathrm{BW})}(t)$ decay is not as steep as $f^{(\mathrm{PL})}(t)$ but has a more gradual decay than $f_2^{(\mathrm{BW})}(t)$, as shown in Figs. A1 and 2. This mask offers a middle ground between the step function like $f_2^{(\mathrm{BW})}(t)$ and the faster decaying $f^{(\mathrm{PL})}(t)$. Table A7 presents the results for *Powerformer* with WCMHA and $f_1^{(\mathrm{BW})}(t)$. We observe that $f_1^{(\mathrm{BW})}(t)$ underperforms $f^{(\mathrm{PL})}(t)$, $f^{(\mathrm{SPL})}(t)$ and $f_2^{(\mathrm{BW})}(t)$. The poor performance further supports our hypothesis that the functional form of $f(t)$ is important when imposing a locality bias.

### F.2.4 Butterworth Filter: Order 2

Here we evaluate the Butterworth order 2 mask ($f_2^{(\mathrm{BW})}(t)$) on the ETT and Weather datasets with a 336 look-back window. The $f_2^{(\mathrm{BW})}(t)$ decay is analogous to the step function, but has a smooth decay to 0 at the end of the window, as shown in Figs. A1. We use $f_2^{(\mathrm{BW})}(t)$ as an improved step function to avoid the sharp cutoff. Table A8 presents the results for *Powerformer* with WCMHA and $f_2^{(\mathrm{BW})}(t)$. We observe that $f_2^{(\mathrm{BW})}(t)$ underperforms $f^{(\mathrm{PL})}(t)$ and $f^{(\mathrm{SPL})}(t)$, but outperforms $f_1^{(\mathrm{BW})}(t)$. The poor performance further supports our hypothesis that the functional form of $f(t)$ is important when imposing a locality bias.
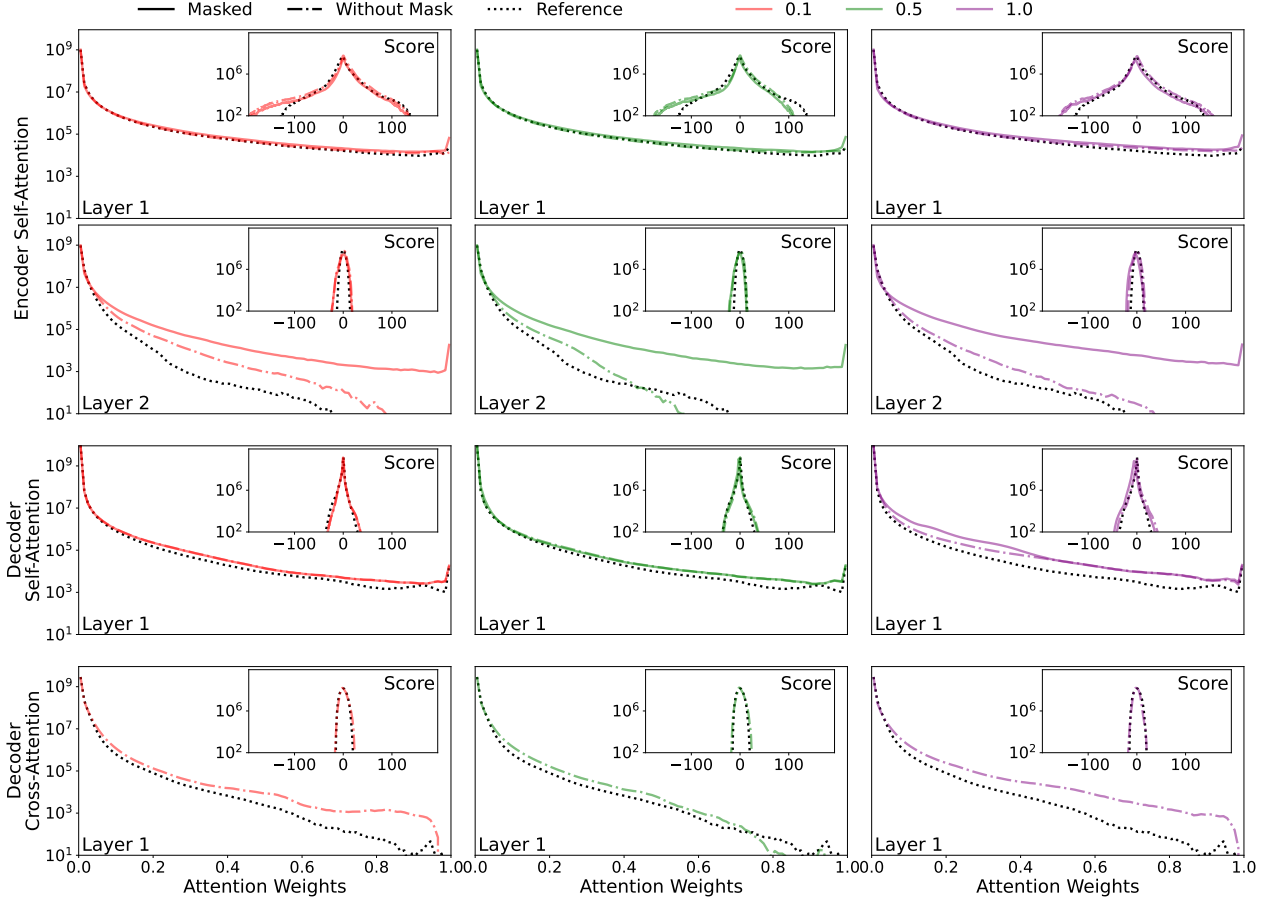
Figure A11: We present the Transformer attention score (inset) and weight distributions on the ETTh2 dataset with a forecasting length of 96. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

## F.3 Learnable Mask Time Decays

Although we treat the decay time $\alpha$ as a hyperparameter, one may consider learning the optimal decay time by making $\alpha$ learnable. We observe that $\alpha$ monotonically decreases and in some cases flips sign, changing our mask from a power-law decay to a power-law growth. To combat this continuous decrease we applied a learning rate scheduler to limit how far $\alpha$ can change from its initialization. To fairly compare with our hyperparameter results, we initialize $\alpha$ at the same times used in our hyperparameter tuning. We present the aggregate MSE and MAE results in Table A9.

We observe very minor differences between the best hyperparameter results and the best results from learning $\alpha$. During training, we still observe a monotonically decreasing $\alpha$ that slows down only because of the learning rate scheduler. This makes sense since the training algorithm gains access to more information as $\alpha$ decreases and facilitates overfitting. This behavior further supports our claim that $\mathbf{M}^{(\mathrm{C,D})}$ and $\alpha$ act as regularizers.

Figure A12: We present the Transformer attention score (inset) and weight distributions on the ETTh2 dataset with a forecasting length of 720. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
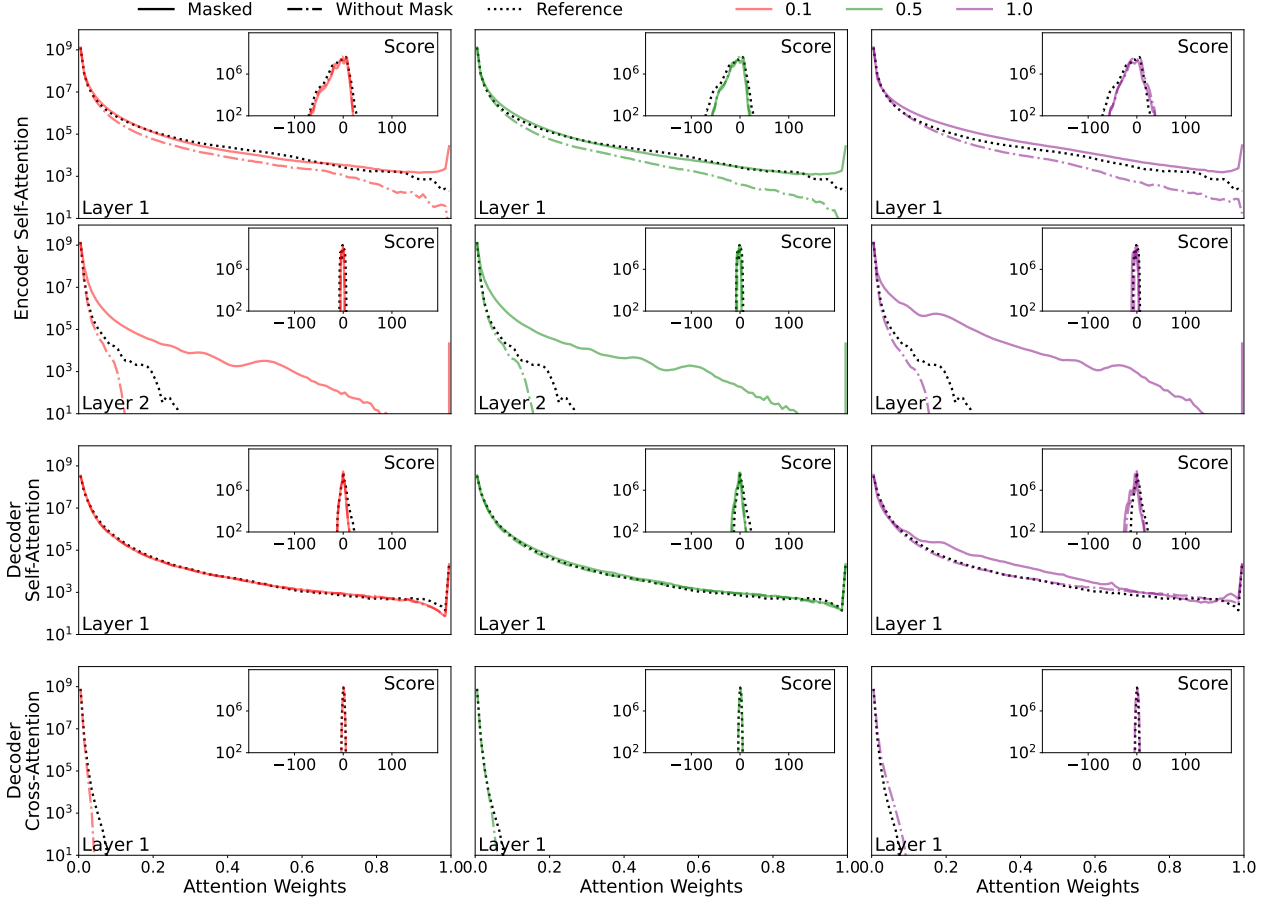
Figure A13: We present the Transformer attention score (inset) and weight distributions on the ETTm1 dataset with a forecasting length of 96. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Figure A14: We present the Transformer attention score (inset) and weight distributions on the ETTm1 dataset with a forecasting length of 720. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Figure A15: We present the Transformer attention score (inset) and weight distributions on the ETTm2 dataset with a forecasting length of 96. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Figure A16: We present the Transformer attention score (inset) and weight distributions on the ETTm2 dataset with a forecasting length of 720. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
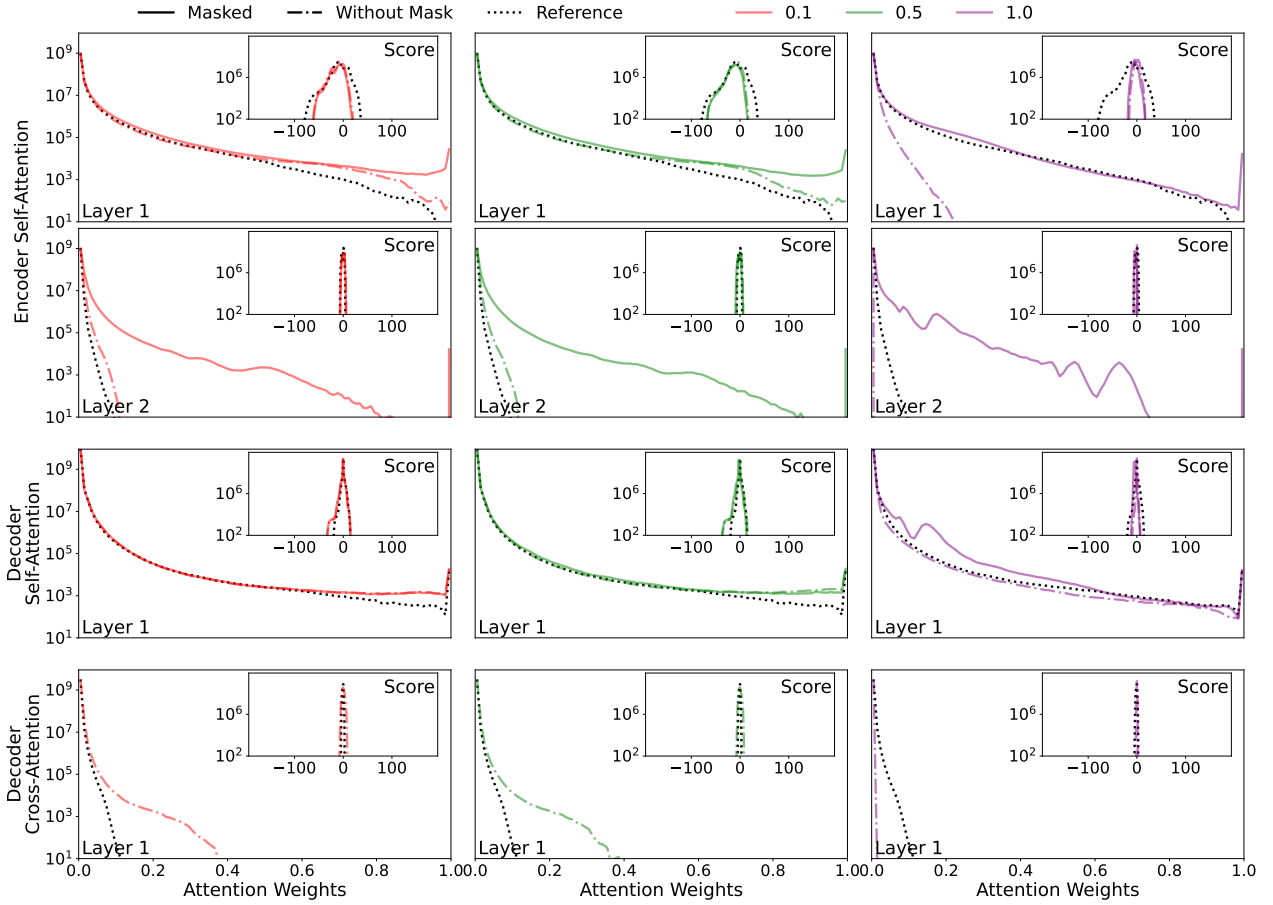
Figure A17: We present the Transformer attention score (inset) and weight distributions on the Weather dataset with a forecasting length of 96. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\text{PL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{PL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.

Figure A18: We present the Transformer attention score (inset) and weight distributions on the Weather dataset with a forecasting length of 720. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
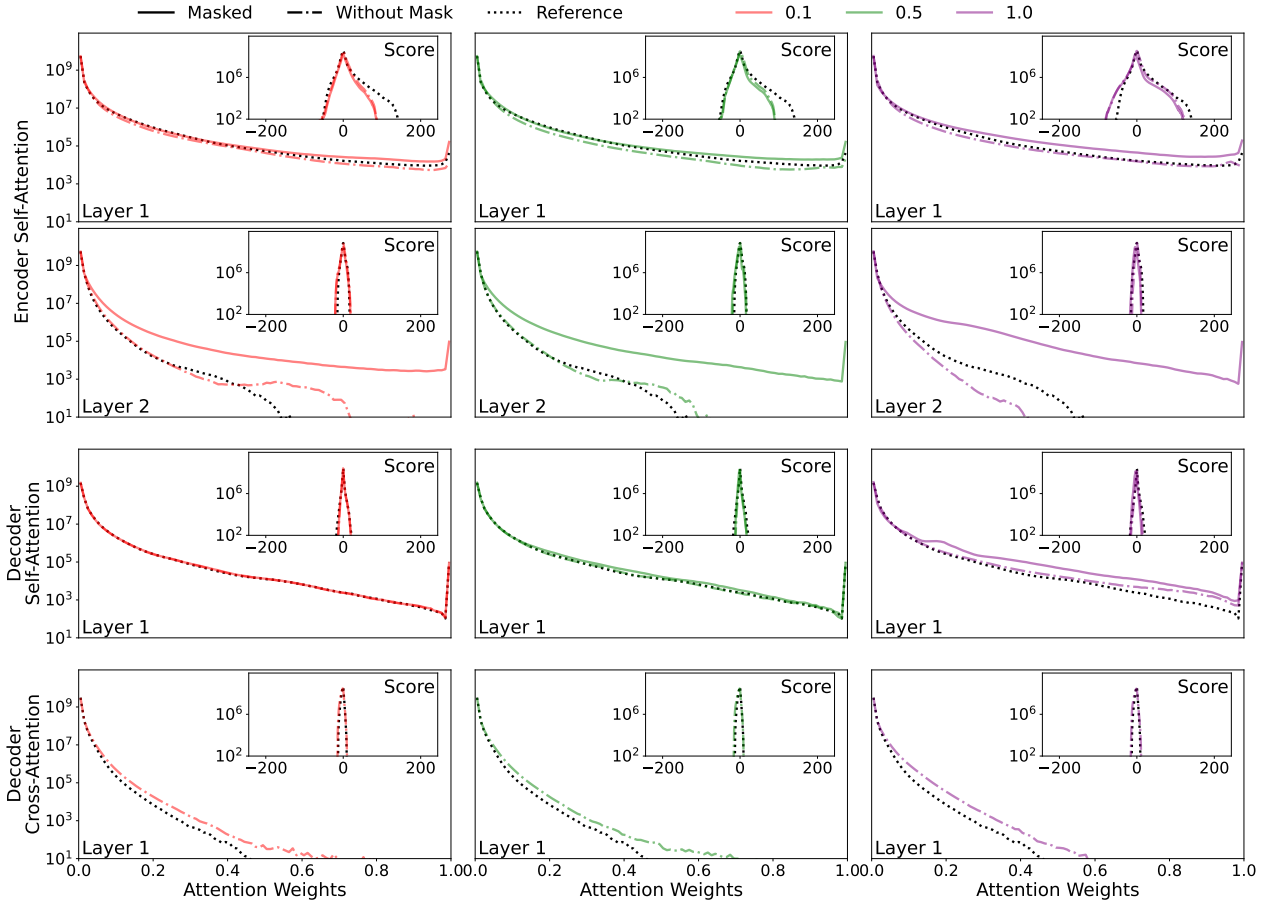
Figure A19: We present the Transformer attention score (inset) and weight distributions on the Electricity dataset with a forecasting length of 96. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
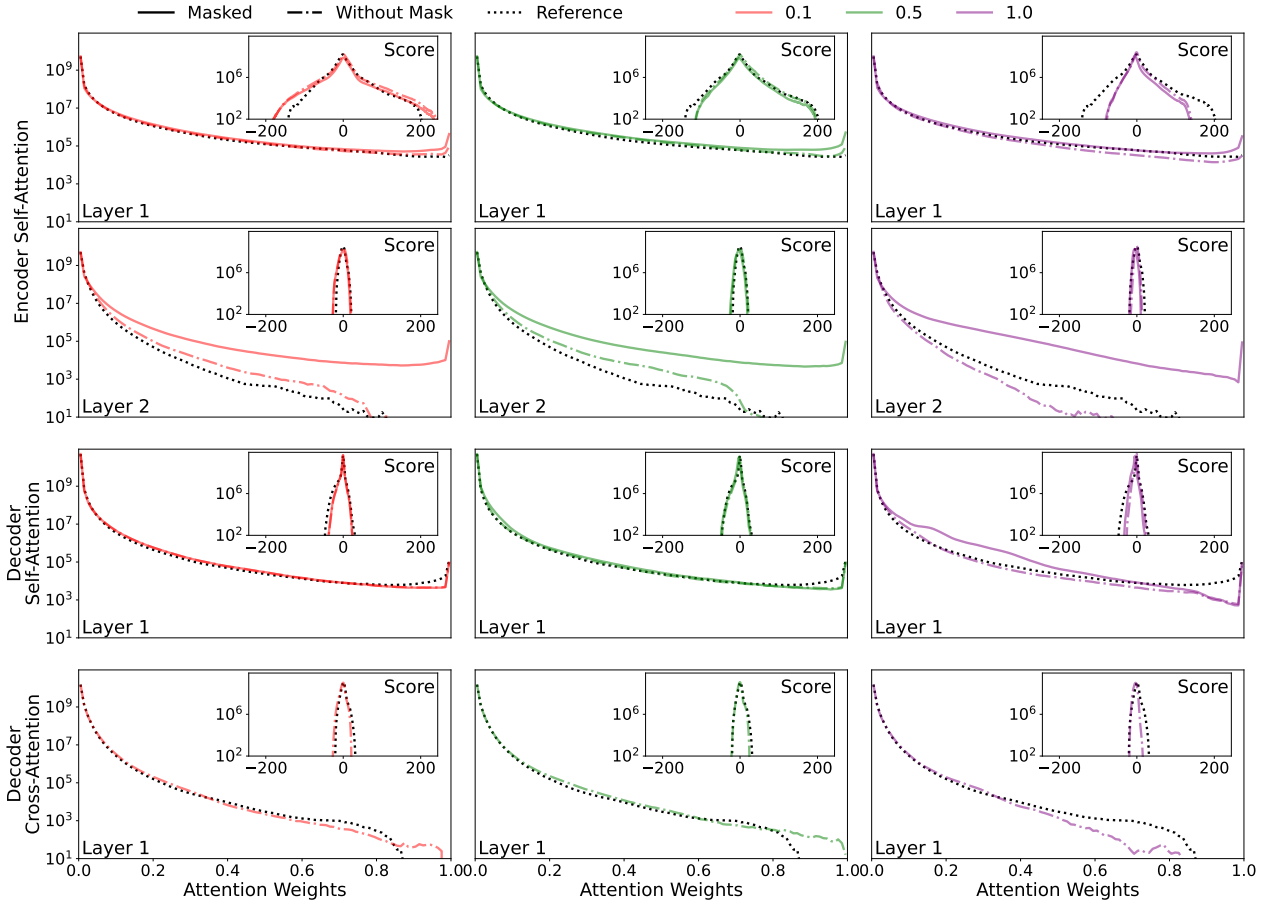
Figure A20: We present the Transformer attention score (inset) and weight distributions on the Electricity dataset with a forecasting length of 720. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
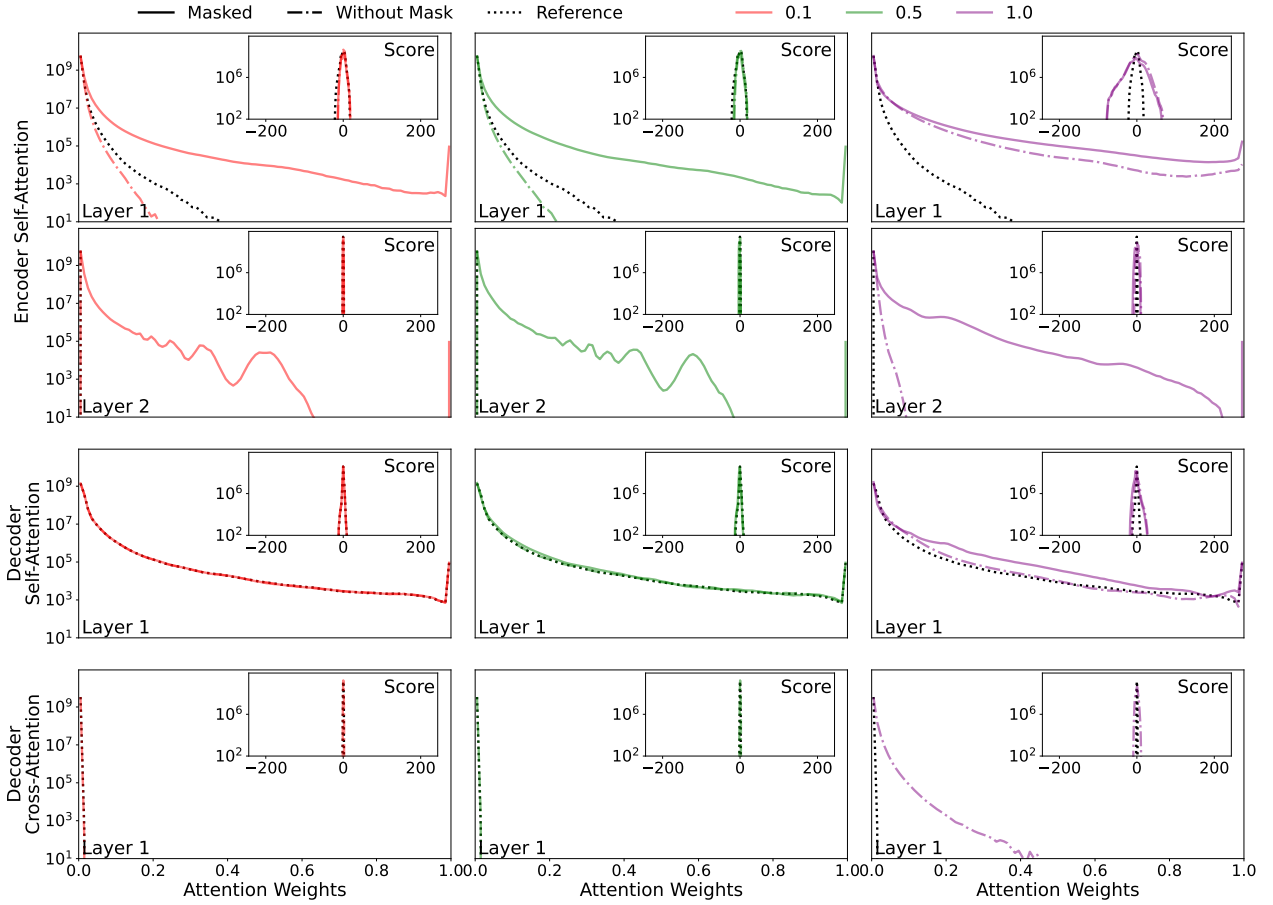
Figure A21: We present the Transformer attention score (inset) and weight distributions on the Traffic dataset with a forecasting length of 96. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
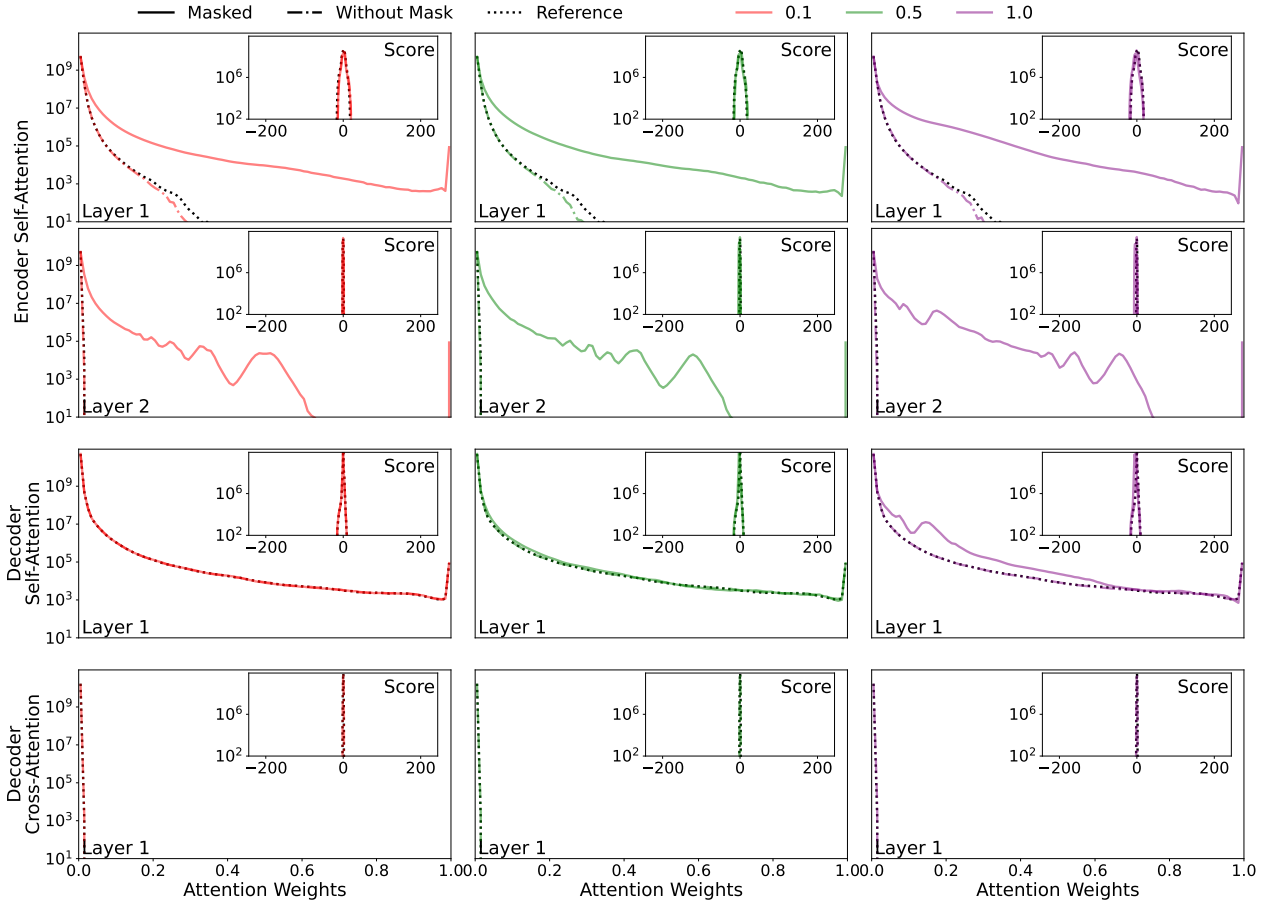
Figure A22: We present the Transformer attention score (inset) and weight distributions on the Traffic dataset with a forecasting length of 720. The dotted line represents the reference MHA Transformer results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Table A5: We compare *Powerformer's* performance with the weight power law mask $f^{(\mathrm{PL})}(t)$ on standard time-series datasets for varying decay lengths. The best results are bolded and second best are underlined.

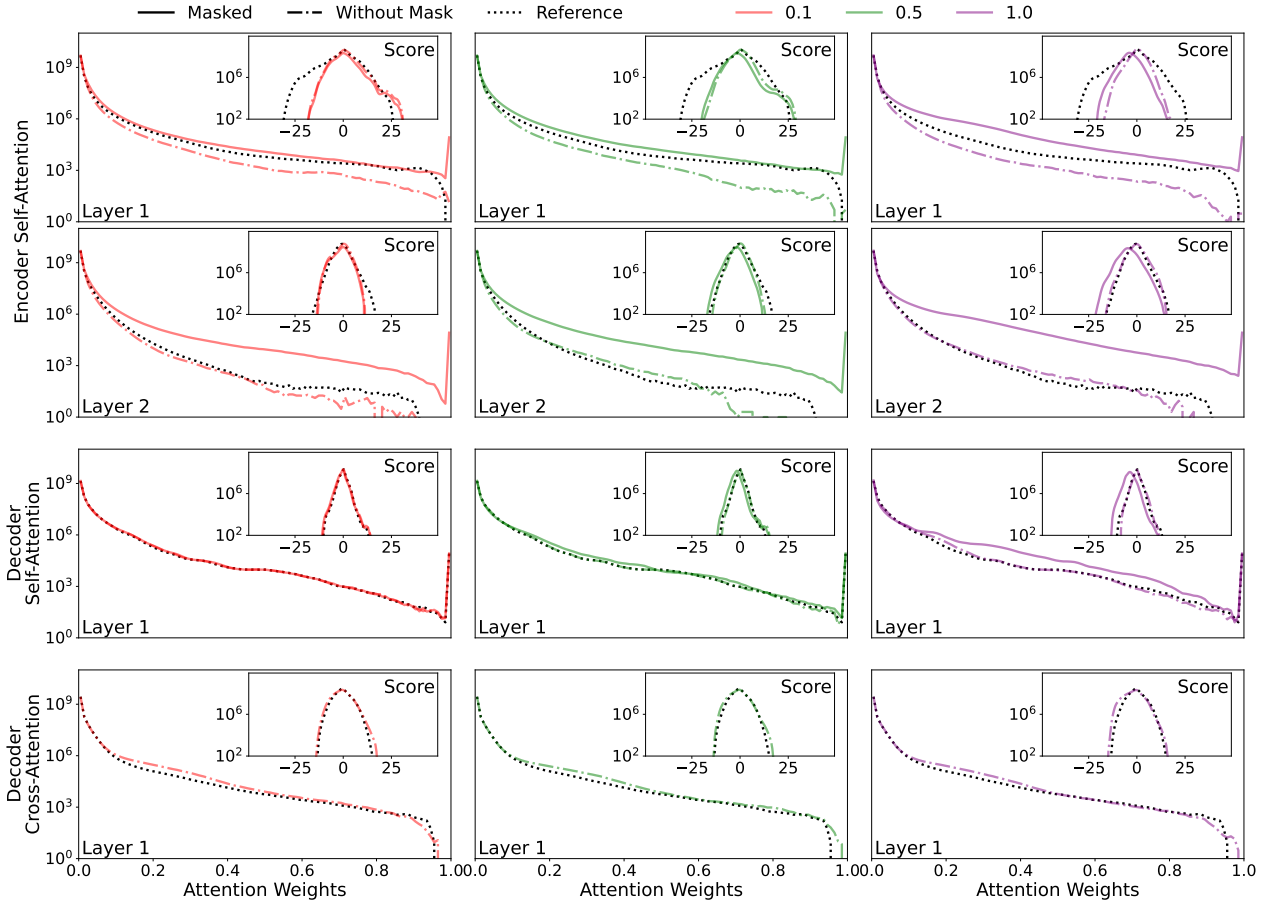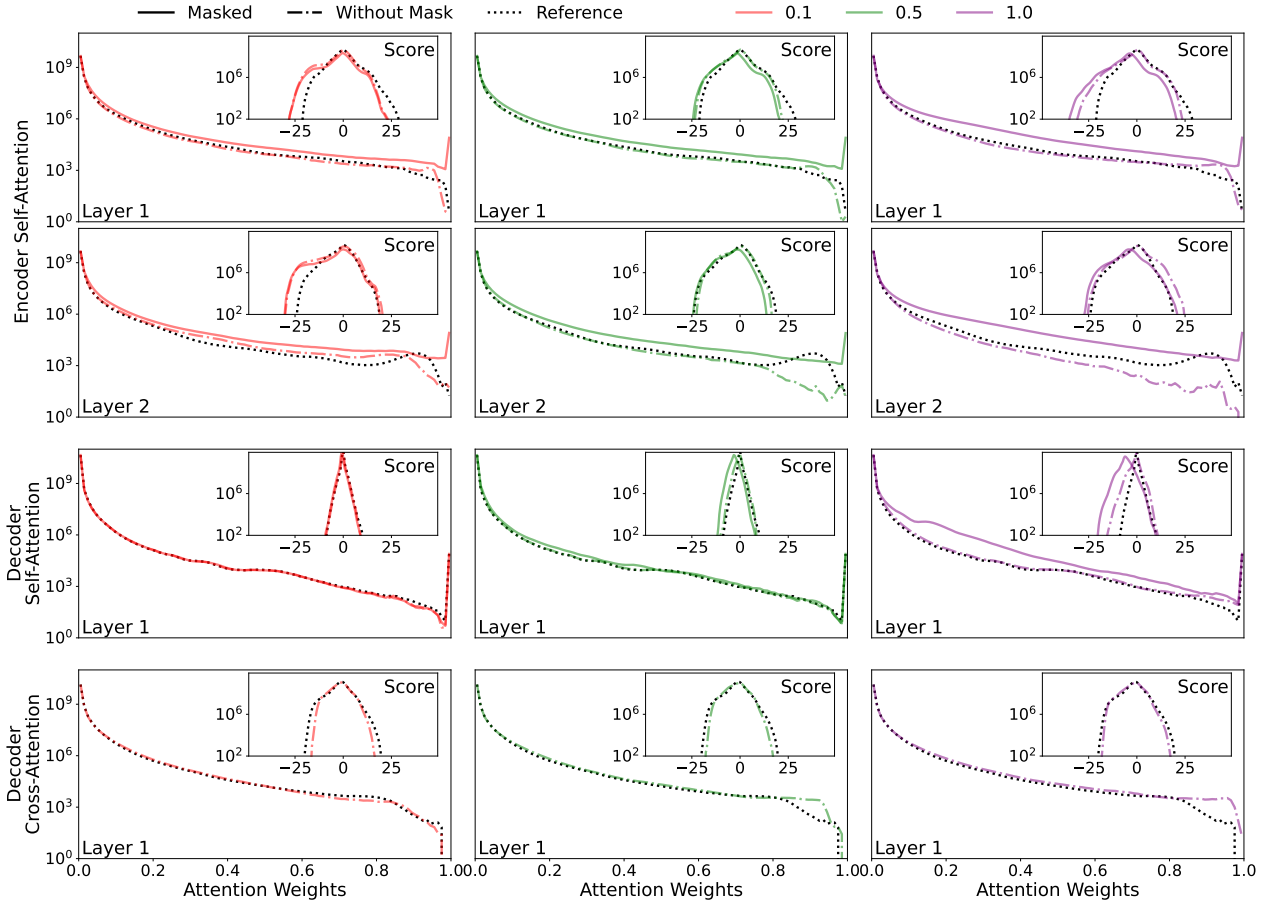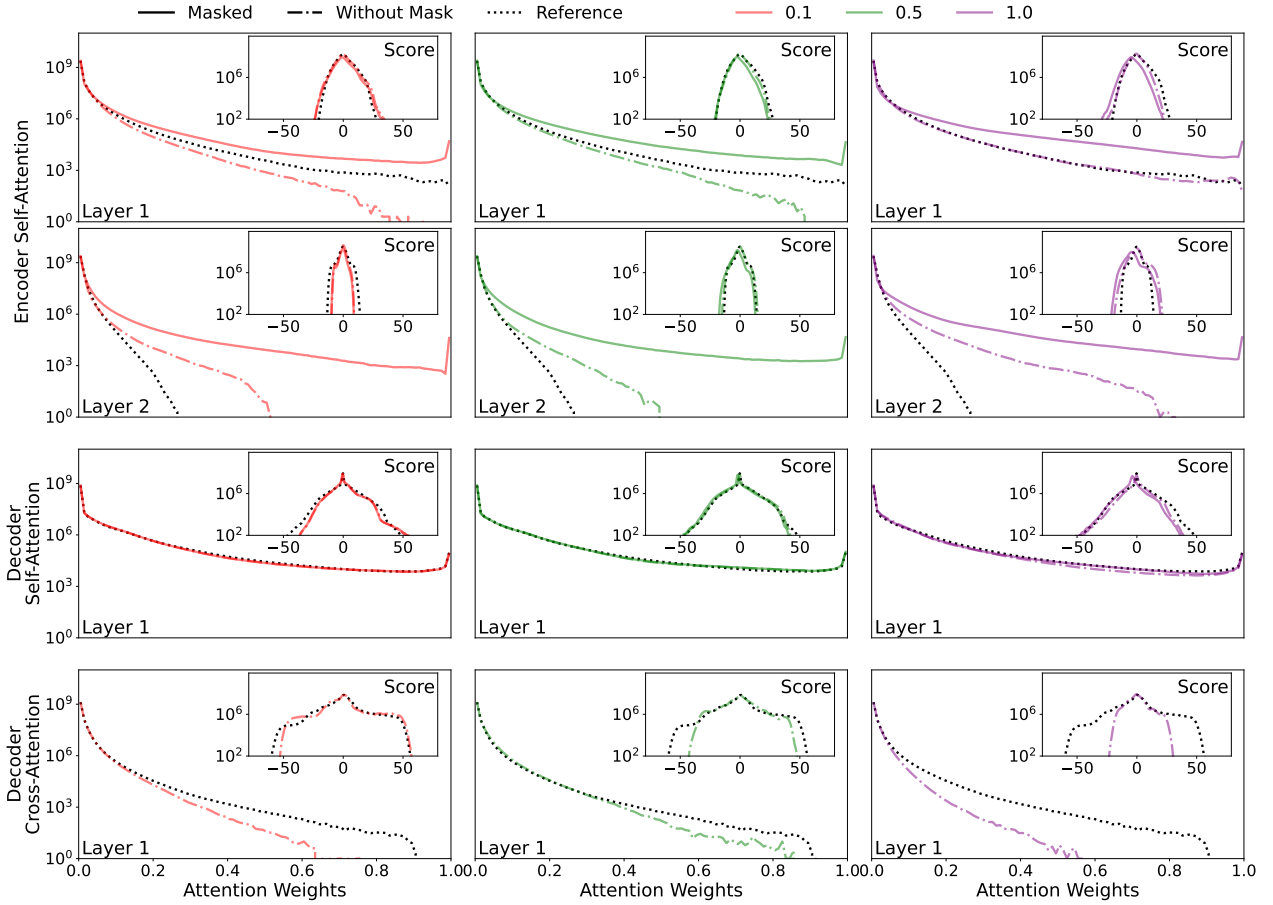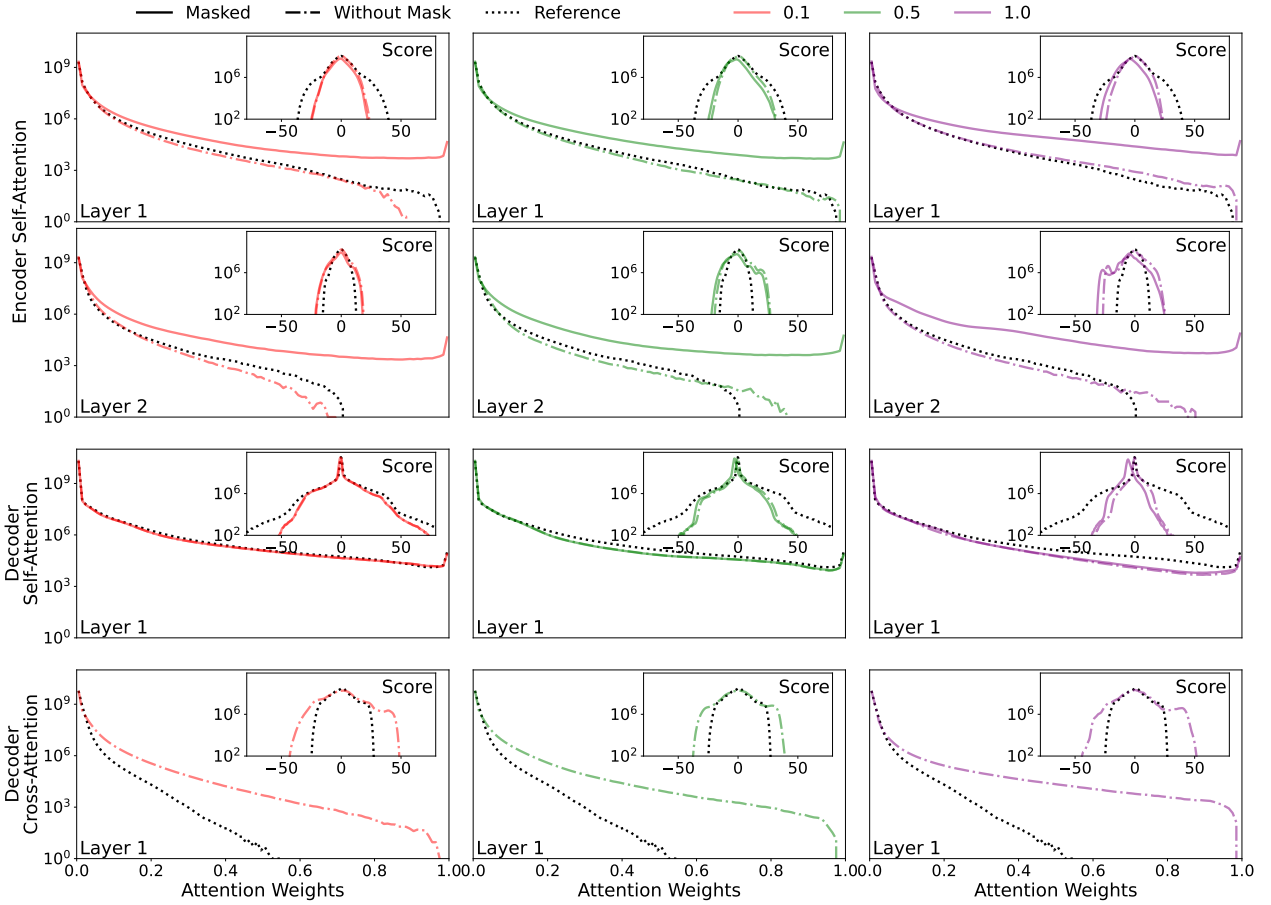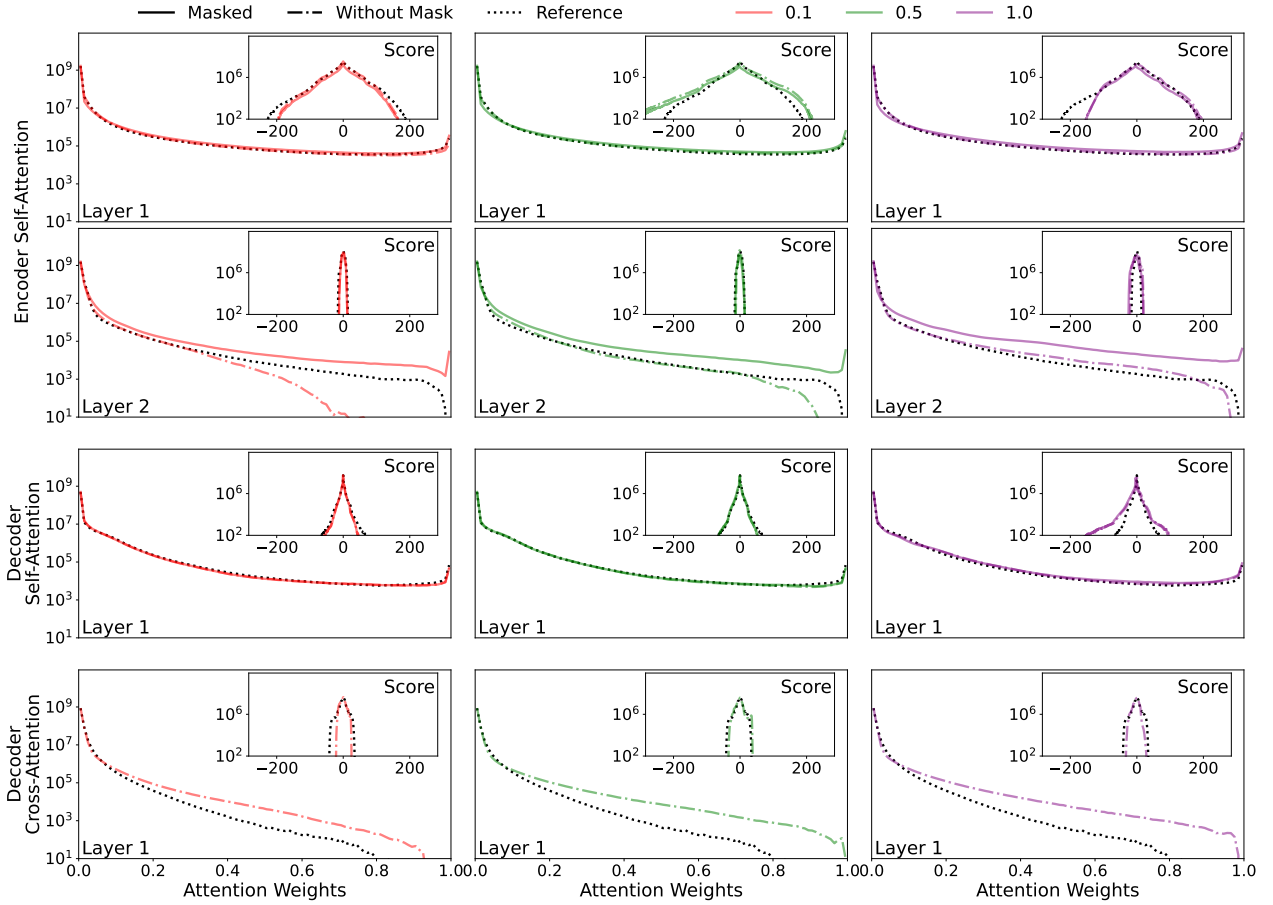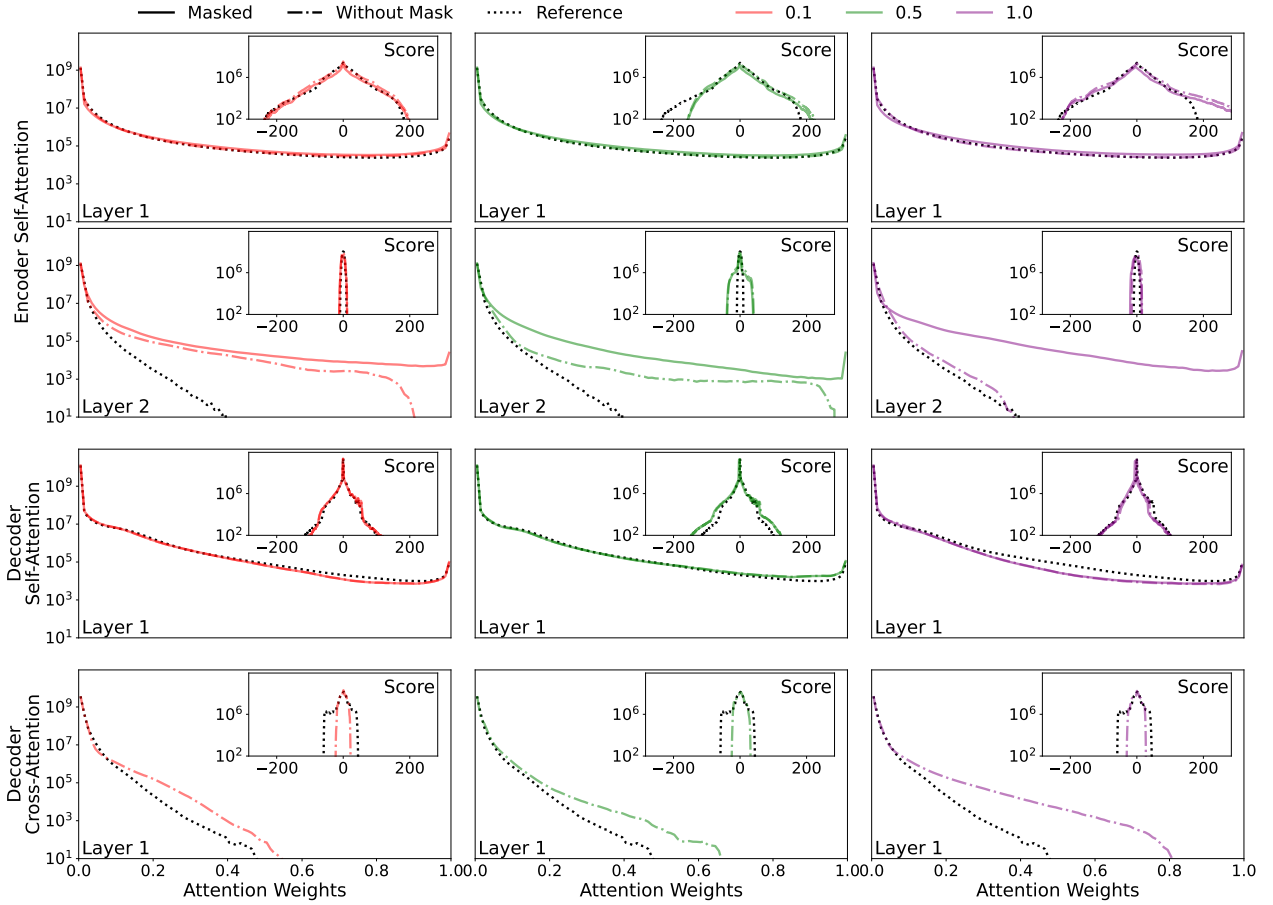| | | Delay | 0.1 | | 0.25 | | 0.5 | | 0.75 | | 1 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Metric | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 336 | 96 | <u>0.377</u> | **0.401** | <u>0.377</u> | **0.401** | <u>0.377</u> | **0.401** | **0.376** | **0.401** | **0.376** | **0.401** |
| | | 192 | **0.413** | **0.421** | **0.413** | **0.421** | **0.413** | **0.421** | **0.413** | **0.421** | <u>0.414</u> | **0.421** |
| | | 336 | **0.424** | **0.430** | <u>0.425</u> | **0.430** | <u>0.425</u> | **0.430** | <u>0.425</u> | **0.430** | <u>0.425</u> | **0.430** |
| | | 720 | **0.437** | 0.455 | **0.437** | 0.455 | **0.437** | 0.455 | **0.437** | <u>0.456</u> | <u>0.438</u> | <u>0.456</u> |
| | 512 | 96 | **0.369** | **0.399** | **0.369** | **0.399** | <u>0.370</u> | **0.399** | <u>0.370</u> | **0.399** | <u>0.370</u> | **0.399** |
| | | 192 | 0.404 | <u>0.420</u> | 0.404 | <u>0.420</u> | **0.402** | **0.418** | <u>0.403</u> | **0.418** | <u>0.403</u> | **0.418** |
| | | 336 | **0.414** | **0.428** | **0.414** | **0.428** | <u>0.415</u> | <u>0.429</u> | <u>0.415</u> | 0.430 | <u>0.416</u> | 0.431 |
| | | 720 | **0.440** | **0.460** | **0.440** | **0.460** | **0.440** | **0.460** | <u>0.443</u> | 0.462 | <u>0.443</u> | **0.460** |
| ETTh2 | 336 | 96 | **0.274** | **0.335** | **0.274** | <u>0.336</u> | <u>0.275</u> | <u>0.336</u> | <u>0.275</u> | **0.335** | <u>0.275</u> | **0.335** |
| | | 192 | 0.341 | <u>0.381</u> | **0.338** | **0.379** | <u>0.340</u> | **0.379** | <u>0.340</u> | **0.379** | 0.341 | **0.379** |
| | | 336 | **0.325** | **0.379** | <u>0.326</u> | <u>0.380</u> | 0.331 | 0.384 | 0.333 | 0.383 | 0.334 | 0.383 |
| | | 720 | <u>0.377</u> | <u>0.420</u> | **0.376** | **0.419** | 0.381 | 0.423 | 0.381 | 0.422 | 0.381 | 0.422 |
| | 512 | 96 | **0.274** | **0.337** | **0.274** | **0.337** | **0.274** | **0.337** | <u>0.275</u> | **0.337** | <u>0.275</u> | <u>0.338</u> |
| | | 192 | **0.340** | **0.381** | **0.340** | **0.381** | <u>0.341</u> | <u>0.382</u> | <u>0.342</u> | <u>0.382</u> | <u>0.342</u> | <u>0.382</u> |
| | | 336 | **0.330** | **0.387** | <u>0.331</u> | **0.387** | <u>0.333</u> | <u>0.388</u> | 0.334 | <u>0.388</u> | 0.334 | <u>0.388</u> |
| | | 720 | **0.383** | <u>0.426</u> | <u>0.384</u> | <u>0.426</u> | **0.383** | <u>0.426</u> | <u>0.384</u> | <u>0.426</u> | **0.383** | **0.425** |
| ETTm1 | 336 | 96 | <u>0.292</u> | <u>0.343</u> | **0.290** | **0.342** | <u>0.292</u> | <u>0.343</u> | 0.293 | 0.344 | <u>0.292</u> | 0.345 |
| | | 192 | <u>0.333</u> | <u>0.371</u> | 0.334 | 0.372 | <u>0.332</u> | <u>0.371</u> | <u>0.332</u> | <u>0.370</u> | **0.330** | **0.369** |
| | | 336 | 0.362 | 0.392 | <u>0.361</u> | 0.391 | 0.362 | <u>0.390</u> | **0.359** | **0.389** | 0.362 | **0.389** |
| | | 720 | <u>0.413</u> | 0.423 | <u>0.413</u> | <u>0.422</u> | <u>0.413</u> | **0.421** | <u>0.413</u> | **0.421** | **0.412** | 0.423 |
| | 512 | 96 | <u>0.291</u> | **0.345** | **0.290** | **0.345** | **0.290** | **0.345** | <u>0.291</u> | 0.346 | 0.292 | 0.346 |
| | | 192 | <u>0.331</u> | <u>0.370</u> | <u>0.330</u> | <u>0.370</u> | **0.329** | **0.368** | 0.332 | <u>0.370</u> | 0.332 | 0.372 |
| | | 336 | <u>0.362</u> | <u>0.391</u> | <u>0.362</u> | **0.390** | **0.361** | 0.392 | <u>0.362</u> | **0.390** | 0.363 | <u>0.391</u> |
| | | 720 | 0.417 | 0.428 | **0.415** | <u>0.425</u> | 0.417 | **0.423** | 0.420 | 0.426 | <u>0.416</u> | 0.430 |
| ETTm2 | 336 | 96 | **0.162** | **0.252** | 0.165 | 0.255 | <u>0.163</u> | <u>0.253</u> | 0.164 | 0.254 | 0.164 | 0.254 |
| | | 192 | **0.222** | **0.292** | **0.222** | <u>0.293</u> | **0.222** | <u>0.293</u> | **0.222** | <u>0.293</u> | **0.222** | <u>0.293</u> |
| | | 336 | **0.277** | **0.329** | **0.277** | **0.329** | **0.277** | **0.329** | **0.277** | **0.329** | <u>0.278</u> | <u>0.330</u> |
| | | 720 | 0.363 | <u>0.381</u> | 0.363 | <u>0.381</u> | <u>0.362</u> | 0.382 | **0.359** | **0.380** | **0.359** | 0.384 |
| | 512 | 96 | 0.165 | **0.255** | 0.165 | **0.255** | 0.165 | 0.256 | 0.165 | 0.257 | **0.164** | **0.255** |
| | | 192 | <u>0.222</u> | 0.295 | 0.224 | 0.297 | 0.224 | 0.296 | 0.223 | 0.296 | **0.221** | **0.294** |
| | | 336 | <u>0.274</u> | <u>0.329</u> | 0.274 | <u>0.329</u> | 0.274 | <u>0.329</u> | **0.272** | **0.327** | <u>0.273</u> | **0.327** |
| | | 720 | <u>0.358</u> | <u>0.383</u> | <u>0.358</u> | <u>0.383</u> | <u>0.358</u> | <u>0.382</u> | <u>0.358</u> | 0.383 | **0.356** | **0.381** |
| Weather | 336 | 96 | **0.151** | <u>0.201</u> | **0.151** | <u>0.201</u> | **0.151** | **0.200** | <u>0.152</u> | **0.200** | 0.154 | 0.202 |
| | | 192 | **0.196** | <u>0.243</u> | **0.196** | <u>0.242</u> | **0.196** | 0.241 | **0.196** | 0.241 | <u>0.197</u> | <u>0.242</u> |
| | | 336 | <u>0.248</u> | 0.283 | **0.247** | <u>0.283</u> | **0.247** | 0.281 | **0.247** | <u>0.282</u> | **0.247** | <u>0.281</u> |
| | | 720 | <u>0.318</u> | 0.334 | 0.318 | 0.334 | 0.318 | 0.334 | 0.317 | **0.333** | 0.317 | 0.333 |
| | 512 | 96 | <u>0.148</u> | 0.198 | **0.147** | **0.197** | **0.147** | 0.199 | <u>0.148</u> | 0.199 | <u>0.148</u> | 0.199 |
| | | 192 | <u>0.193</u> | <u>0.241</u> | 0.193 | <u>0.241</u> | 0.193 | <u>0.241</u> | **0.191** | **0.239** | <u>0.193</u> | <u>0.241</u> |
| | | 336 | **0.243** | <u>0.281</u> | 0.244 | <u>0.281</u> | **0.243** | <u>0.280</u> | 0.243 | 0.279 | 0.244 | <u>0.280</u> |
| | | 720 | **0.311** | <u>0.330</u> | **0.311** | <u>0.330</u> | 0.314 | 0.331 | <u>0.312</u> | <u>0.330</u> | 0.311 | **0.329** |
| Electricity | 336 | 96 | <u>0.131</u> | **0.224** | <u>0.131</u> | **0.224** | **0.130** | **0.224** | <u>0.131</u> | **0.224** | <u>0.131</u> | **0.224** |
| | | 192 | **0.147** | <u>0.241</u> | **0.147** | **0.240** | **0.147** | **0.240** | **0.147** | <u>0.241</u> | <u>0.148</u> | <u>0.241</u> |
| | | 336 | **0.164** | <u>0.259</u> | **0.164** | **0.257** | **0.164** | 0.259 | **0.164** | <u>0.258</u> | **0.164** | <u>0.259</u> |
| | | 720 | **0.201** | **0.291** | **0.201** | <u>0.292</u> | <u>0.202</u> | 0.292 | **0.201** | **0.291** | <u>0.202</u> | 0.293 |
| | 512 | 96 | <u>0.130</u> | <u>0.224</u> | **0.129** | <u>0.224</u> | **0.129** | <u>0.224</u> | **0.129** | **0.223** | **0.129** | **0.223** |
| | | 192 | <u>0.147</u> | <u>0.241</u> | **0.146** | <u>0.241</u> | <u>0.147</u> | **0.240** | **0.146** | **0.240** | <u>0.147</u> | <u>0.241</u> |
| | | 336 | <u>0.163</u> | <u>0.258</u> | **0.162** | **0.257** | **0.162** | <u>0.258</u> | **0.162** | **0.257** | **0.162** | <u>0.258</u> |
| | | 720 | <u>0.199</u> | <u>0.290</u> | **0.198** | <u>0.290</u> | **0.198** | <u>0.290</u> | **0.198** | 0.289 | **0.198** | <u>0.290</u> |
| Traffic | 336 | 96 | <u>0.371</u> | **0.253** | **0.370** | **0.253** | 0.372 | <u>0.254</u> | 0.372 | <u>0.254</u> | 0.373 | 0.255 |
| | | 192 | **0.388** | **0.260** | **0.388** | **0.260** | 0.390 | <u>0.261</u> | 0.390 | <u>0.261</u> | 0.390 | 0.261 |
| | | 336 | **0.400** | <u>0.267</u> | **0.400** | **0.266** | 0.404 | 0.269 | <u>0.402</u> | 0.268 | 0.403 | 0.268 |
| | | 720 | <u>0.435</u> | <u>0.287</u> | **0.434** | 0.286 | <u>0.435</u> | 0.286 | **0.434** | 0.286 | **0.434** | 0.286 |
| | 512 | 96 | **0.365** | 0.252 | **0.365** | 0.252 | <u>0.367</u> | <u>0.253</u> | <u>0.367</u> | <u>0.253</u> | 0.394 | 0.282 |
| | | 192 | **0.382** | **0.258** | 0.391 | 0.269 | <u>0.383</u> | <u>0.259</u> | **0.382** | 0.260 | 0.407 | 0.286 |
| | | 336 | 0.393 | **0.265** | **0.391** | **0.265** | <u>0.392</u> | **0.265** | 0.393 | **0.265** | 0.394 | **0.265** |
| | | 720 | 0.435 | 0.289 | **0.432** | 0.288 | **0.432** | <u>0.287</u> | <u>0.431</u> | <u>0.287</u> | **0.430** | **0.286** |

Figure A23: We present the *Powerformer* attention score (inset) and weight distributions on the ETTh1 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{PL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{PL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.



Figure A24: We present the *Powerformer* attention score (inset) and weight distributions on the ETTh1 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{PL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{PL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.

Figure A25: We present the *Powerformer* attention score (inset) and weight distributions on the ETTh2 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.



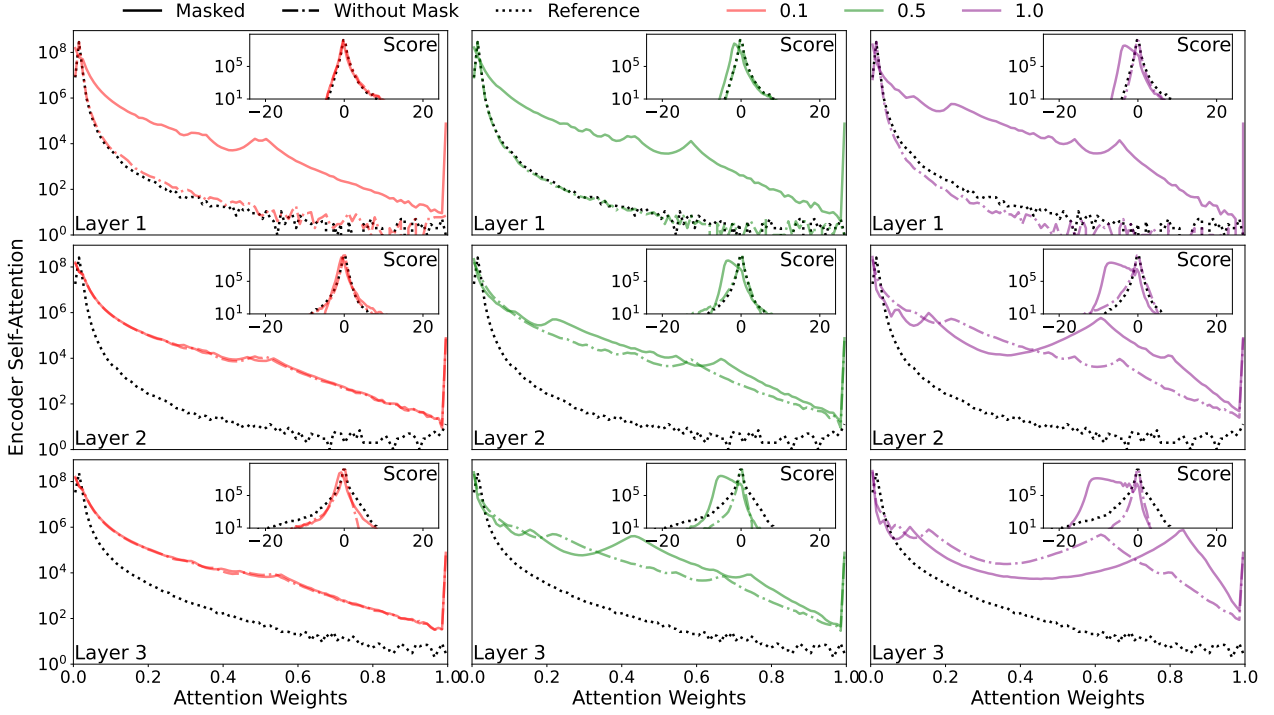Figure A26: We present the *Powerformer* attention score (inset) and weight distributions on the ETTh2 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Figure A27: We present the *Powerformer* attention score (inset) and weight distributions on the ETTm1 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(PL)}(t)$ results before applying $\mathbf{M}^{(C,D)}$, and the solid lines represent WCMHA with $f^{(PL)}(t)$ results after applying $\mathbf{M}^{(C,D)}$.
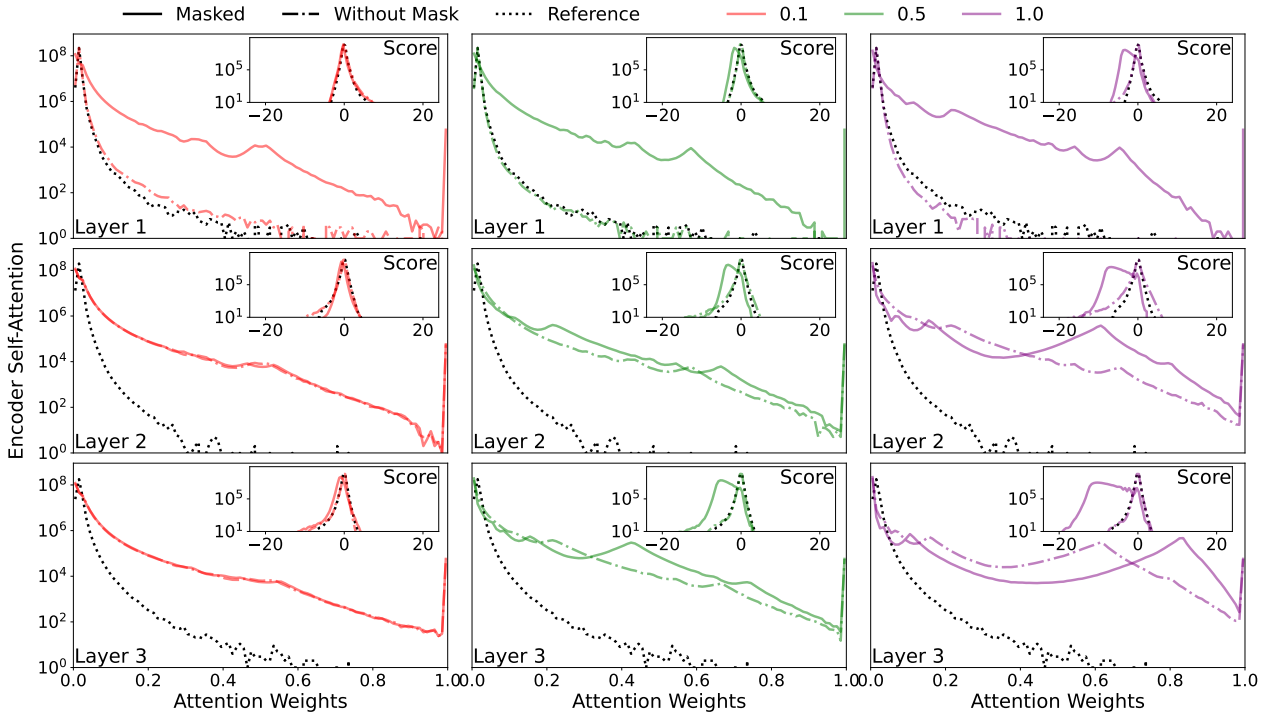


Figure A28: We present the *Powerformer* attention score (inset) and weight distributions on the ETTm1 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(PL)}(t)$ results before applying $\mathbf{M}^{(C,D)}$, and the solid lines represent WCMHA with $f^{(PL)}(t)$ results after applying $\mathbf{M}^{(C,D)}$.
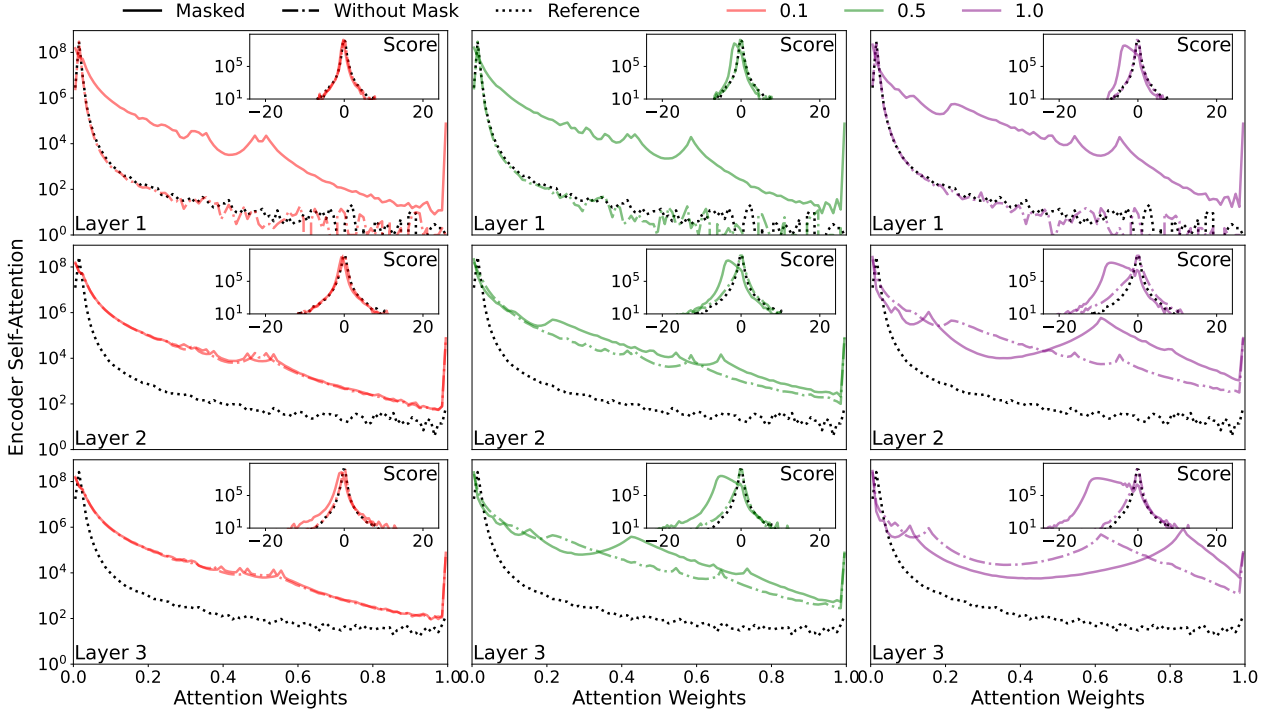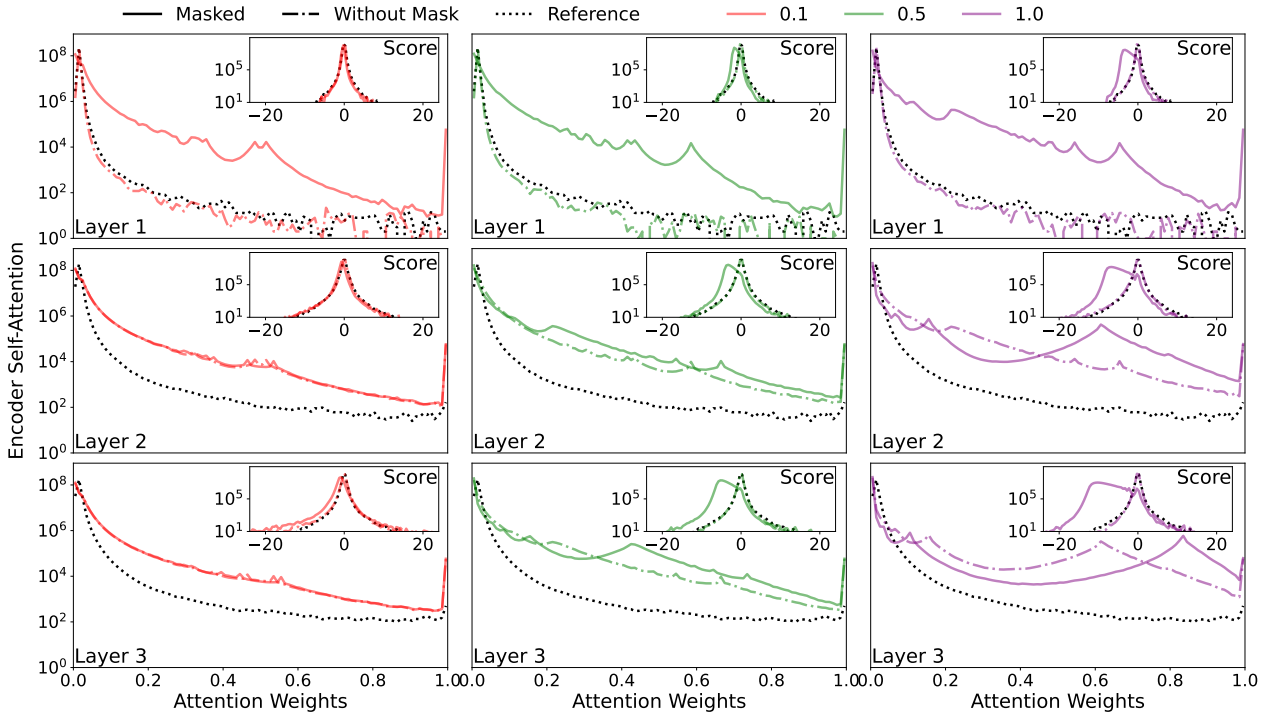
Figure A29: We present the *Powerformer* attention score (inset) and weight distributions on the ETTm2 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{PL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{PL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.



Figure A30: We present the *Powerformer* attention score (inset) and weight distributions on the ETTm2 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{PL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{PL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.
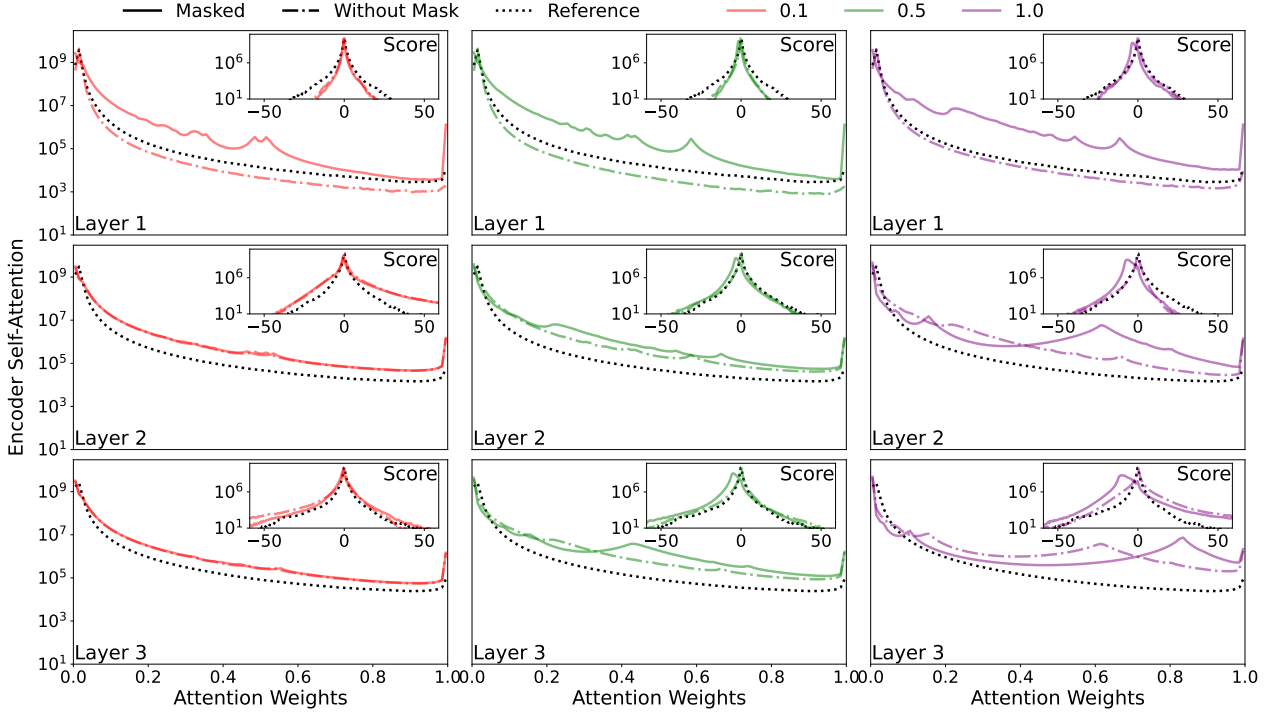
Figure A31: We present the *Powerformer* attention score (inset) and weight distributions on the Weather dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.



Figure A32: We present the *Powerformer* attention score (inset) and weight distributions on the Weather dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
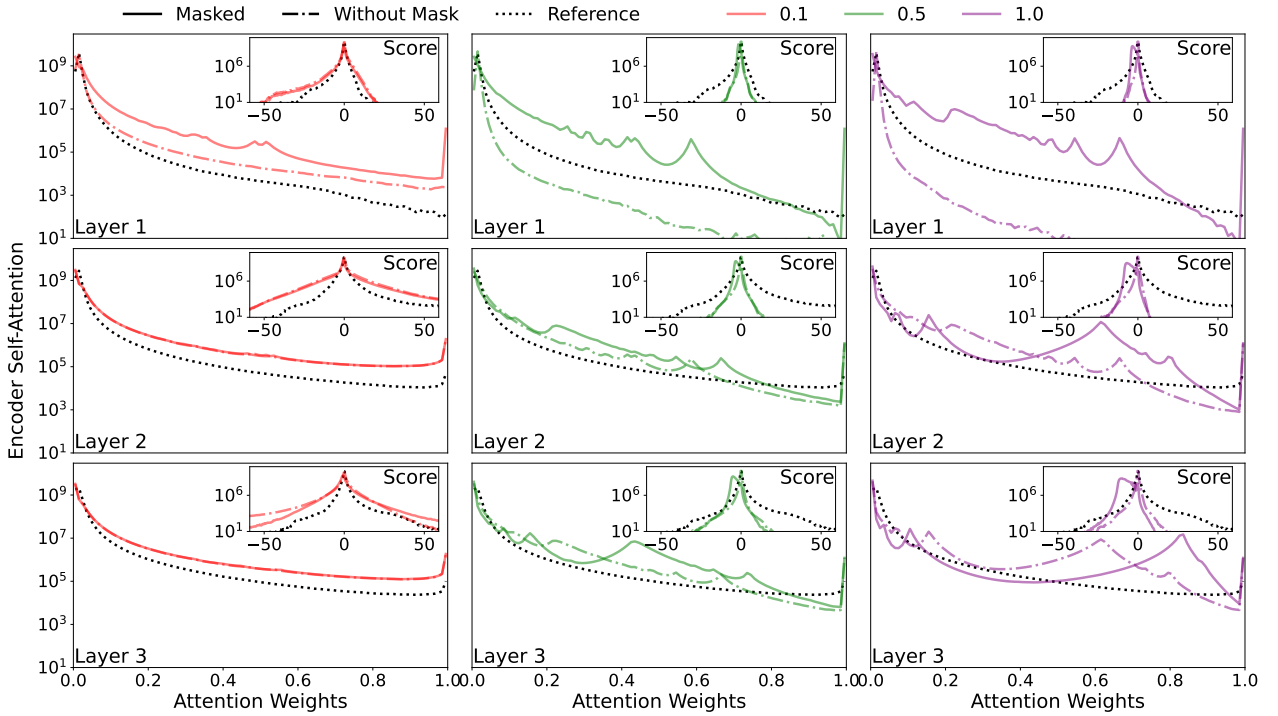
Figure A33: We present the *Powerformer* attention score (inset) and weight distributions on the Electricity dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
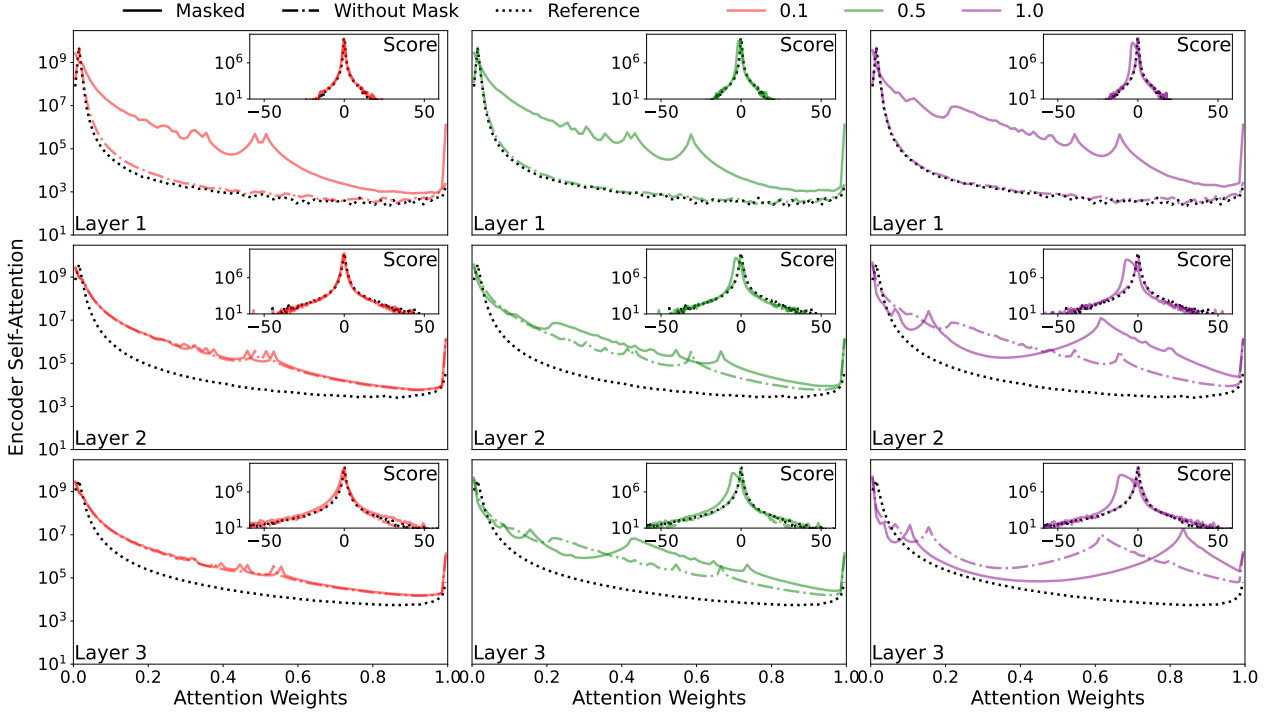


Figure A34: We present the *Powerformer* attention score (inset) and weight distributions on the Electricity dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
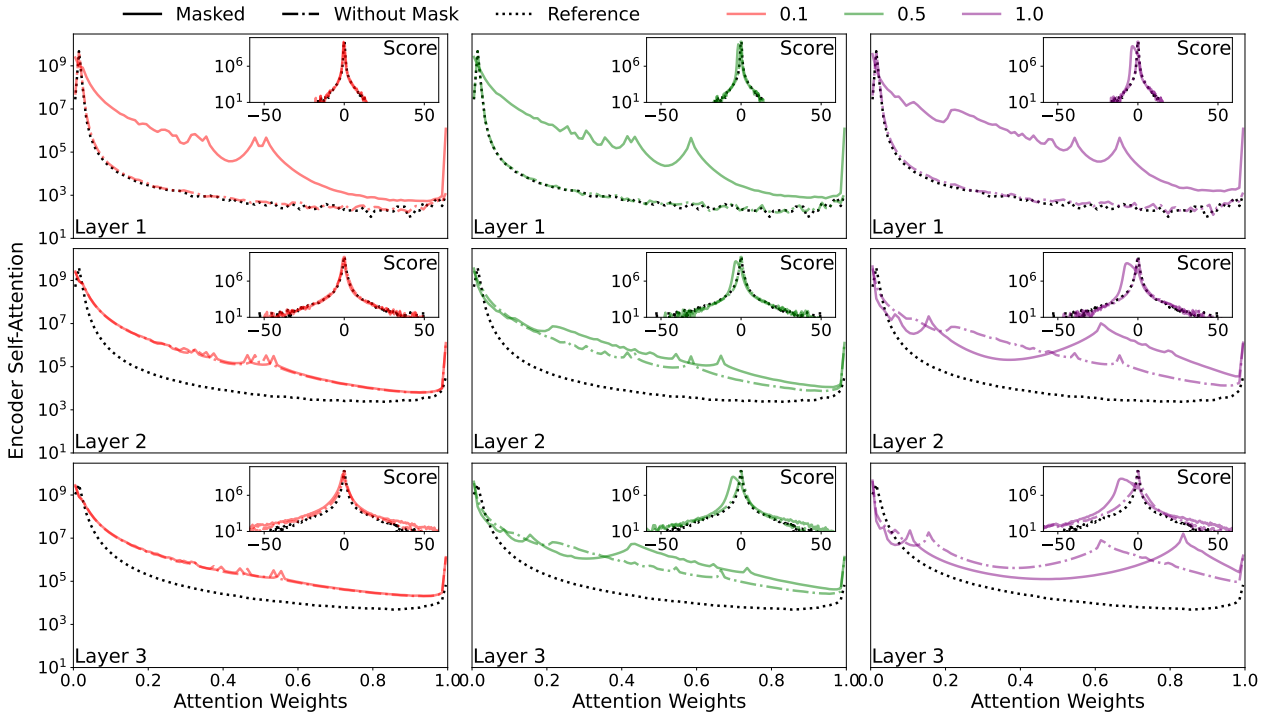
Figure A35: We present the *Powerformer* attention score (inset) and weight distributions on the Traffic dataset with a forecast and input length of 96 and 336, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

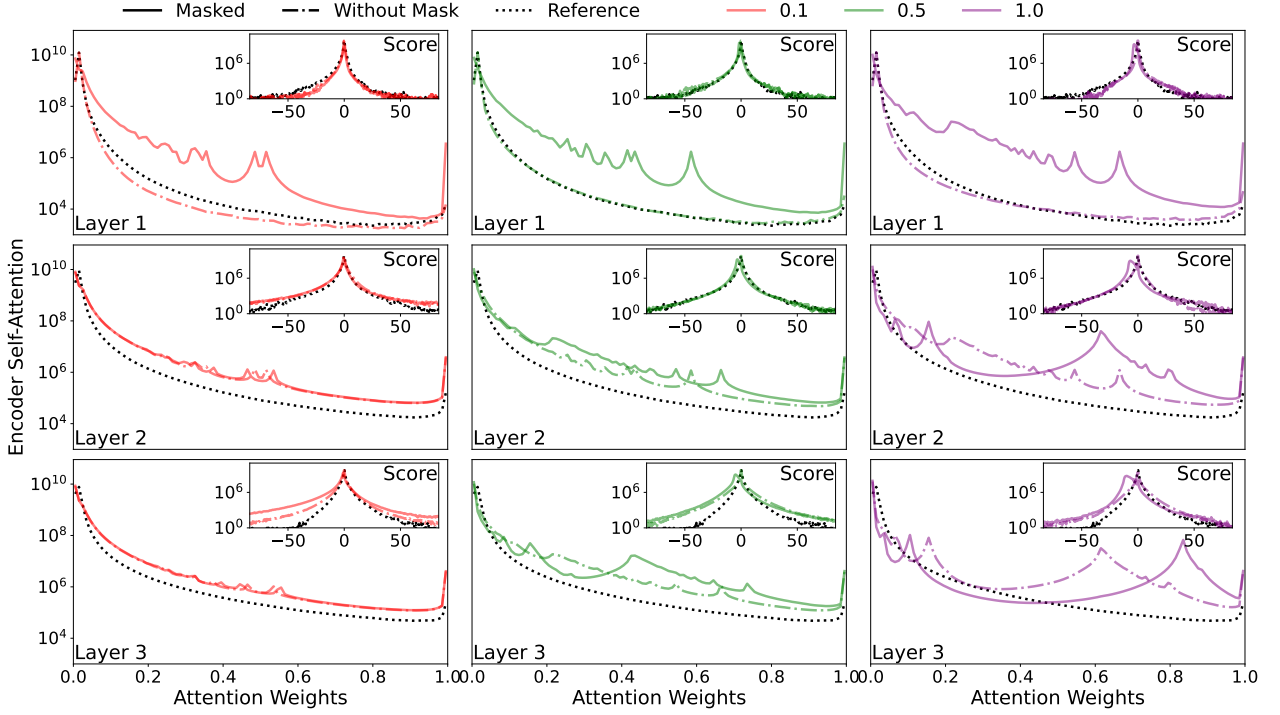

Figure A36: We present the *Powerformer* attention score (inset) and weight distributions on the Traffic dataset with a forecast and input length of 720 and 336, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{PL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{PL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Table A6: We compare *Powerformer's* performance with the similarity power-law mask $f^{(\mathrm{SPL})}(t)$ on standard time-series datasets for varying decay lengths. The best results are bolded and the second best are underlined.

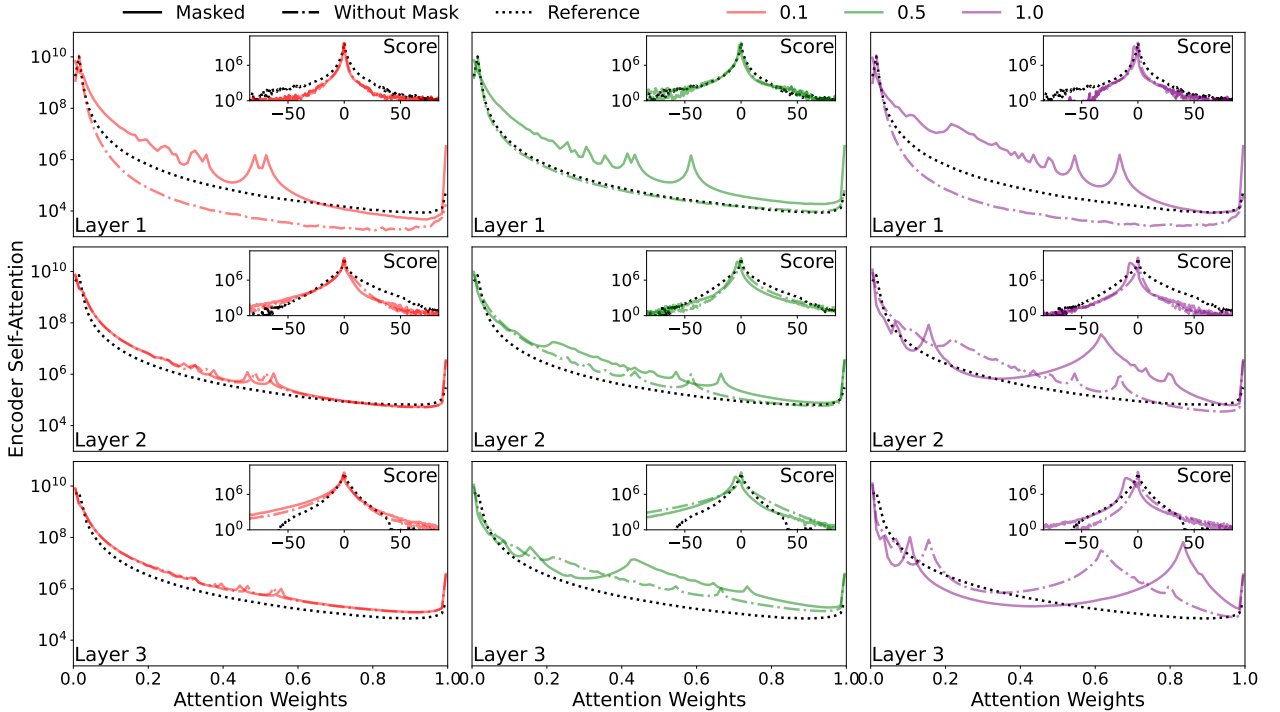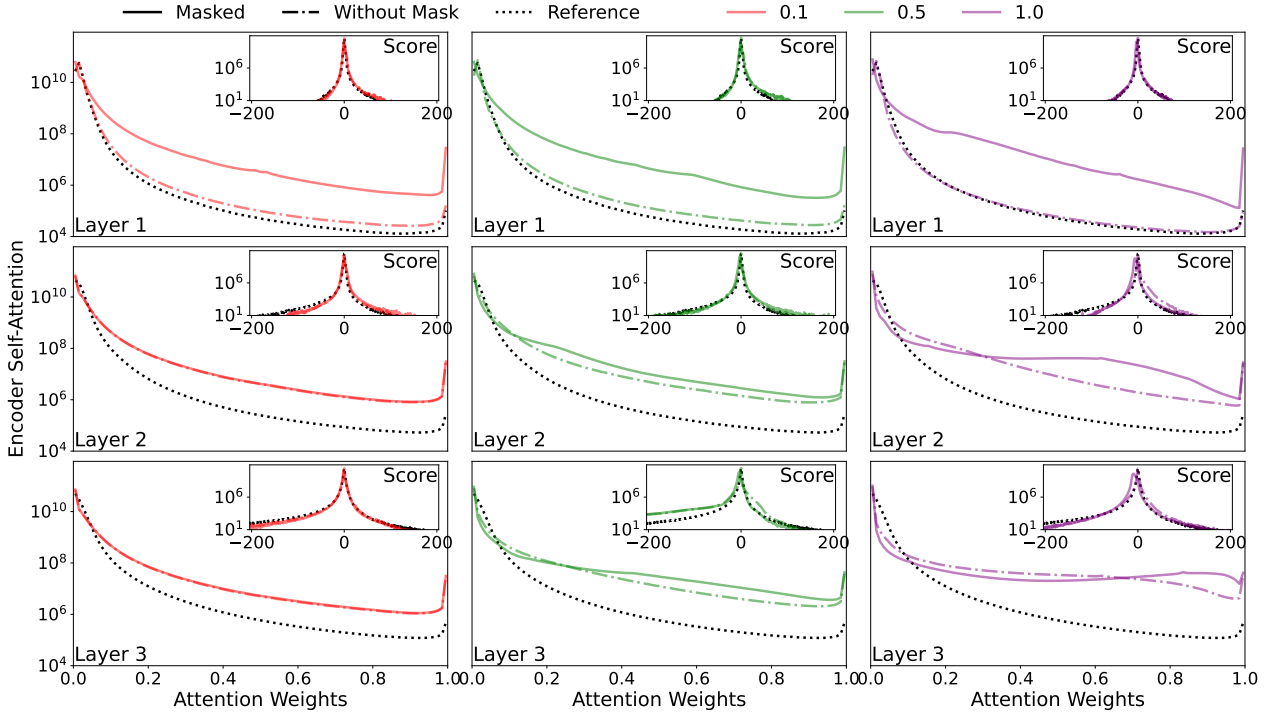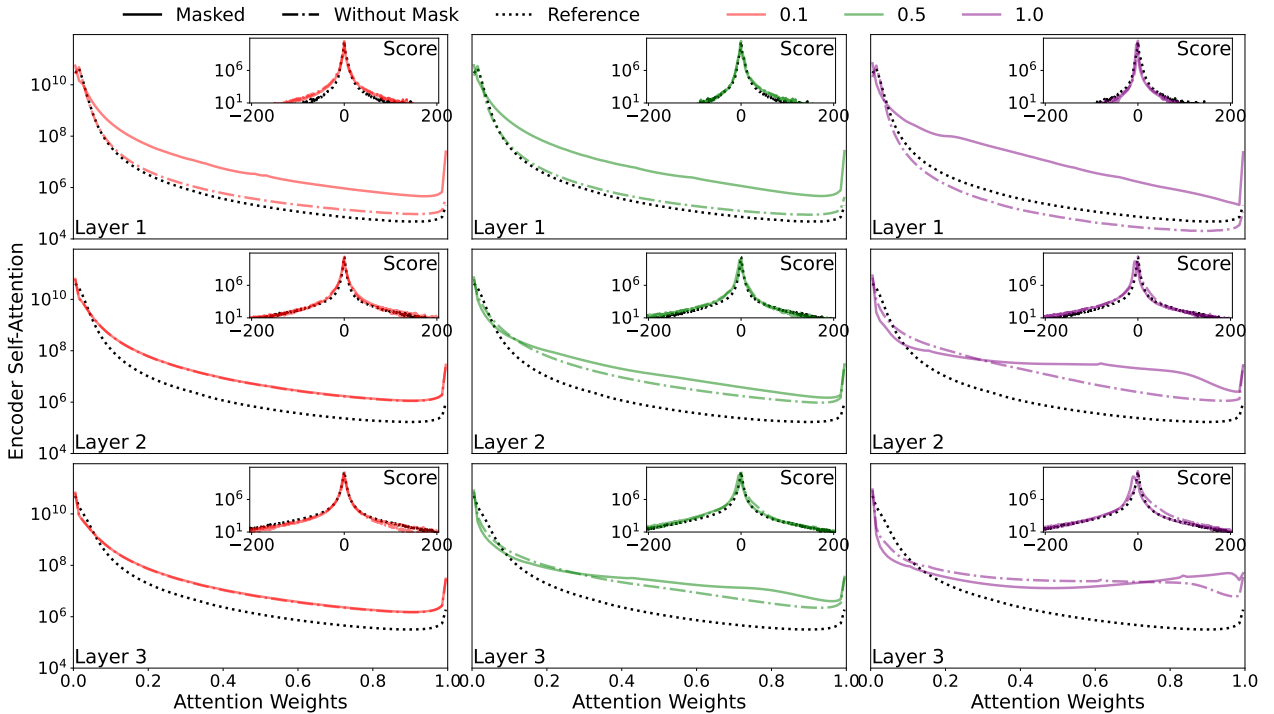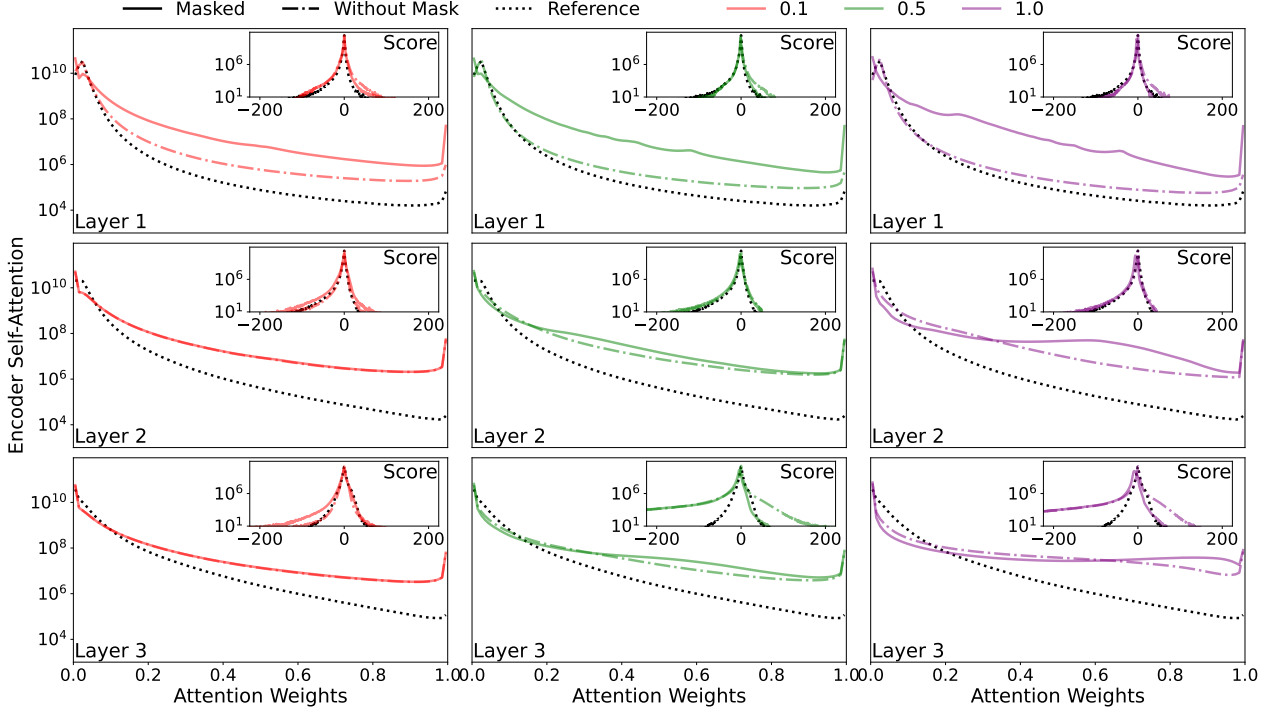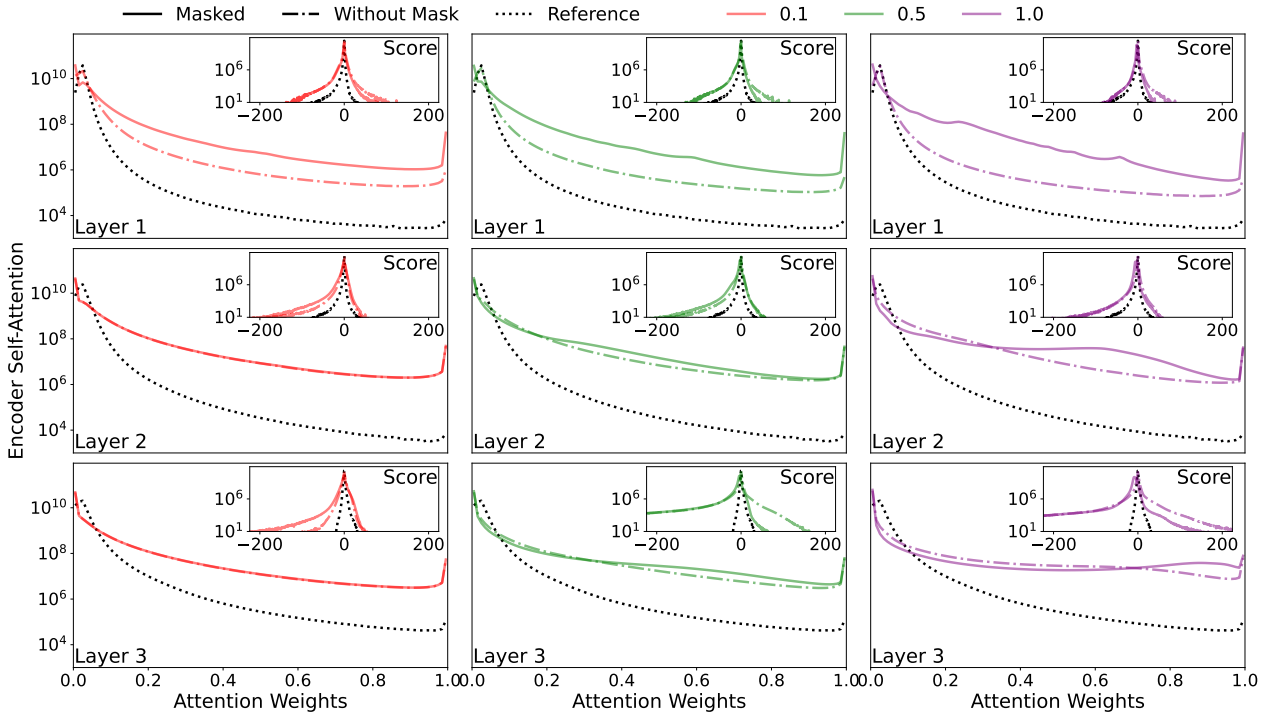| | | Delay | 0.1 | | 0.5 | | 1 | | 2 | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Metric | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 336 | 96 | **0.377** | **0.401** | <u>0.378</u> | **0.402** | <u>0.378</u> | <u>0.402</u> | 0.379 | 0.403 |
| | | 192 | **0.413** | **0.421** | <u>0.414</u> | **0.421** | 0.415 | <u>0.422</u> | 0.415 | <u>0.422</u> |
| | | 336 | **0.424** | **0.429** | <u>0.425</u> | **0.429** | <u>0.425</u> | **0.429** | <u>0.425</u> | **0.429** |
| | | 720 | **0.439** | **0.457** | <u>0.440</u> | **0.457** | <u>0.443</u> | <u>0.458</u> | <u>0.443</u> | <u>0.458</u> |
| | 512 | 96 | **0.370** | **0.400** | <u>0.371</u> | **0.400** | <u>0.371</u> | **0.400** | <u>0.371</u> | **0.400** |
| | | 192 | **0.404** | **0.420** | <u>0.405</u> | **0.420** | <u>0.405</u> | 0.421 | <u>0.405</u> | **0.420** |
| | | 336 | 0.418 | 0.433 | <u>0.417</u> | 0.429 | **0.415** | **0.428** | **0.415** | **0.428** |
| | | 720 | **0.439** | **0.459** | <u>0.441</u> | <u>0.460</u> | 0.442 | 0.461 | 0.442 | 0.461 |
| ETTh2 | 336 | 96 | **0.274** | **0.336** | <u>0.275</u> | **0.336** | <u>0.275</u> | **0.336** | <u>0.275</u> | **0.336** |
| | | 192 | **0.338** | **0.379** | <u>0.340</u> | **0.379** | **0.338** | **0.379** | **0.338** | **0.379** |
| | | 336 | <u>0.330</u> | **0.384** | 0.333 | <u>0.385</u> | **0.329** | <u>0.385</u> | **0.329** | <u>0.385</u> |
| | | 720 | <u>0.380</u> | **0.422** | 0.381 | **0.422** | 0.379 | **0.422** | 0.379 | **0.422** |
| | 512 | 96 | **0.274** | **0.337** | 0.276 | 0.338 | <u>0.275</u> | 0.338 | <u>0.275</u> | 0.338 |
| | | 192 | **0.340** | **0.381** | 0.342 | <u>0.382</u> | <u>0.341</u> | **0.381** | <u>0.341</u> | **0.381** |
| | | 336 | 0.332 | 0.388 | 0.333 | <u>0.387</u> | <u>0.331</u> | 0.387 | **0.330** | **0.386** |
| | | 720 | 0.386 | <u>0.428</u> | 0.381 | **0.424** | <u>0.381</u> | **0.424** | **0.380** | **0.424** |
| ETTm1 | 336 | 96 | 0.295 | 0.345 | **0.292** | **0.343** | 0.293 | <u>0.344</u> | 0.294 | 0.345 |
| | | 192 | 0.336 | 0.372 | **0.329** | **0.369** | <u>0.333</u> | <u>0.371</u> | 0.337 | 0.373 |
| | | 336 | 0.363 | <u>0.392</u> | **0.360** | **0.390** | 0.363 | 0.393 | <u>0.362</u> | <u>0.392</u> |
| | | 720 | <u>0.413</u> | **0.423** | 0.419 | 0.425 | 0.425 | 0.427 | **0.412** | <u>0.424</u> |
| | 512 | 96 | <u>0.291</u> | 0.345 | **0.289** | 0.345 | 0.293 | 0.346 | <u>0.291</u> | **0.345** |
| | | 192 | **0.333** | 0.372 | 0.334 | 0.374 | 0.334 | <u>0.374</u> | **0.333** | 0.374 |
| | | 336 | **0.362** | **0.391** | 0.365 | **0.391** | **0.362** | 0.396 | **0.362** | 0.396 |
| | | 720 | **0.419** | 0.427 | 0.426 | <u>0.418</u> | 0.425 | <u>0.418</u> | 0.425 | **0.417** |
| ETTm2 | 336 | 96 | **0.163** | **0.253** | <u>0.164</u> | <u>0.254</u> | 0.165 | 0.256 | 0.165 | <u>0.254</u> |
| | | 192 | **0.222** | **0.293** | **0.222** | <u>0.294</u> | <u>0.223</u> | <u>0.294</u> | <u>0.223</u> | <u>0.294</u> |
| | | 336 | **0.277** | **0.329** | <u>0.278</u> | <u>0.330</u> | 0.279 | <u>0.330</u> | <u>0.278</u> | **0.329** |
| | | 720 | 0.363 | <u>0.381</u> | **0.358** | 0.383 | **0.358** | **0.380** | <u>0.361</u> | 0.384 |
| | 512 | 96 | **0.164** | **0.255** | **0.164** | **0.255** | <u>0.165</u> | <u>0.256</u> | <u>0.165</u> | <u>0.256</u> |
| | | 192 | **0.222** | **0.295** | <u>0.223</u> | <u>0.296</u> | <u>0.224</u> | <u>0.296</u> | <u>0.224</u> | <u>0.296</u> |
| | | 336 | **0.274** | **0.329** | <u>0.275</u> | **0.329** | <u>0.275</u> | **0.329** | <u>0.275</u> | **0.329** |
| | | 720 | <u>0.354</u> | **0.379** | <u>0.354</u> | 0.380 | **0.353** | 0.380 | **0.353** | **0.379** |
| Weather | 336 | 96 | **0.151** | **0.199** | <u>0.154</u> | 0.202 | <u>0.154</u> | 0.202 | 0.156 | 0.204 |
| | | 192 | **0.195** | <u>0.243</u> | <u>0.197</u> | **0.242** | 0.198 | <u>0.244</u> | 0.199 | <u>0.244</u> |
| | | 336 | **0.247** | <u>0.283</u> | <u>0.248</u> | **0.282** | 0.249 | **0.282** | <u>0.248</u> | <u>0.283</u> |
| | | 720 | <u>0.318</u> | <u>0.334</u> | **0.317** | **0.332** | 0.319 | 0.335 | <u>0.319</u> | <u>0.334</u> |
| | 512 | 96 | **0.149** | **0.199** | <u>0.150</u> | <u>0.200</u> | 0.152 | 0.202 | 0.154 | 0.204 |
| | | 192 | **0.192** | **0.240** | <u>0.194</u> | <u>0.242</u> | 0.196 | 0.244 | 0.198 | 0.244 |
| | | 336 | **0.243** | **0.280** | <u>0.244</u> | **0.280** | 0.246 | 0.283 | 0.246 | <u>0.281</u> |
| | | 720 | **0.310** | **0.329** | <u>0.311</u> | <u>0.330</u> | 0.312 | 0.331 | 0.313 | <u>0.331</u> |
| Electricity | 336 | 96 | **0.131** | **0.224** | **0.131** | <u>0.225</u> | **0.131** | <u>0.225</u> | **0.131** | <u>0.225</u> |
| | | 192 | <u>0.147</u> | 0.240 | 0.148 | <u>0.242</u> | **0.146** | **0.240** | **0.146** | **0.240** |
| | | 336 | <u>0.164</u> | 0.258 | 0.163 | **0.258** | **0.162** | 0.258 | <u>0.163</u> | 0.258 |
| | | 720 | **0.201** | 0.291 | <u>0.202</u> | 0.292 | <u>0.202</u> | 0.293 | <u>0.204</u> | 0.295 |
| | 512 | 96 | **0.129** | **0.224** | 0.130 | **0.224** | **0.129** | **0.224** | **0.129** | **0.224** |
| | | 192 | <u>0.146</u> | 0.240 | 0.147 | 0.241 | **0.145** | 0.240 | **0.145** | 0.240 |
| | | 336 | **0.162** | <u>0.257</u> | 0.163 | 0.259 | 0.163 | 0.259 | 0.163 | 0.259 |
| | | 720 | **0.198** | 0.290 | <u>0.199</u> | <u>0.291</u> | 0.205 | 0.295 | 0.207 | 0.296 |
| Traffic | 336 | 96 | **0.372** | <u>0.254</u> | <u>0.374</u> | 0.256 | 0.375 | <u>0.255</u> | 0.376 | <u>0.255</u> |
| | | 192 | **0.390** | <u>0.261</u> | <u>0.391</u> | <u>0.262</u> | <u>0.391</u> | <u>0.262</u> | 0.392 | <u>0.262</u> |
| | | 336 | **0.402** | <u>0.268</u> | <u>0.403</u> | <u>0.269</u> | <u>0.403</u> | **0.268** | 0.404 | <u>0.269</u> |
| | | 720 | **0.435** | <u>0.286</u> | <u>0.436</u> | <u>0.287</u> | **0.435** | <u>0.287</u> | 0.436 | <u>0.287</u> |
| | 512 | 96 | **0.368** | **0.254** | <u>0.394</u> | <u>0.281</u> | <u>0.394</u> | <u>0.281</u> | 0.395 | 0.282 |
| | | 192 | **0.383** | **0.260** | <u>0.407</u> | <u>0.286</u> | <u>0.407</u> | <u>0.286</u> | 0.407 | 0.287 |
| | | 336 | **0.393** | **0.266** | <u>0.395</u> | **0.266** | <u>0.395</u> | **0.266** | <u>0.395</u> | **0.266** |
| | | 720 | <u>0.432</u> | <u>0.288</u> | 0.433 | 0.288 | 0.434 | 0.288 | **0.431** | **0.287** |

Figure A37: We present *Powerformer's* attention score and weight distributions on the ETTh1 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.



Figure A38: We present *Powerformer's* attention score and weight distributions on the ETTh1 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Figure A39: We present *Powerformer's* attention score and weight distributions on the ETTh2 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.



Figure A40: We present *Powerformer's* attention score and weight distributions on the ETTh2 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.
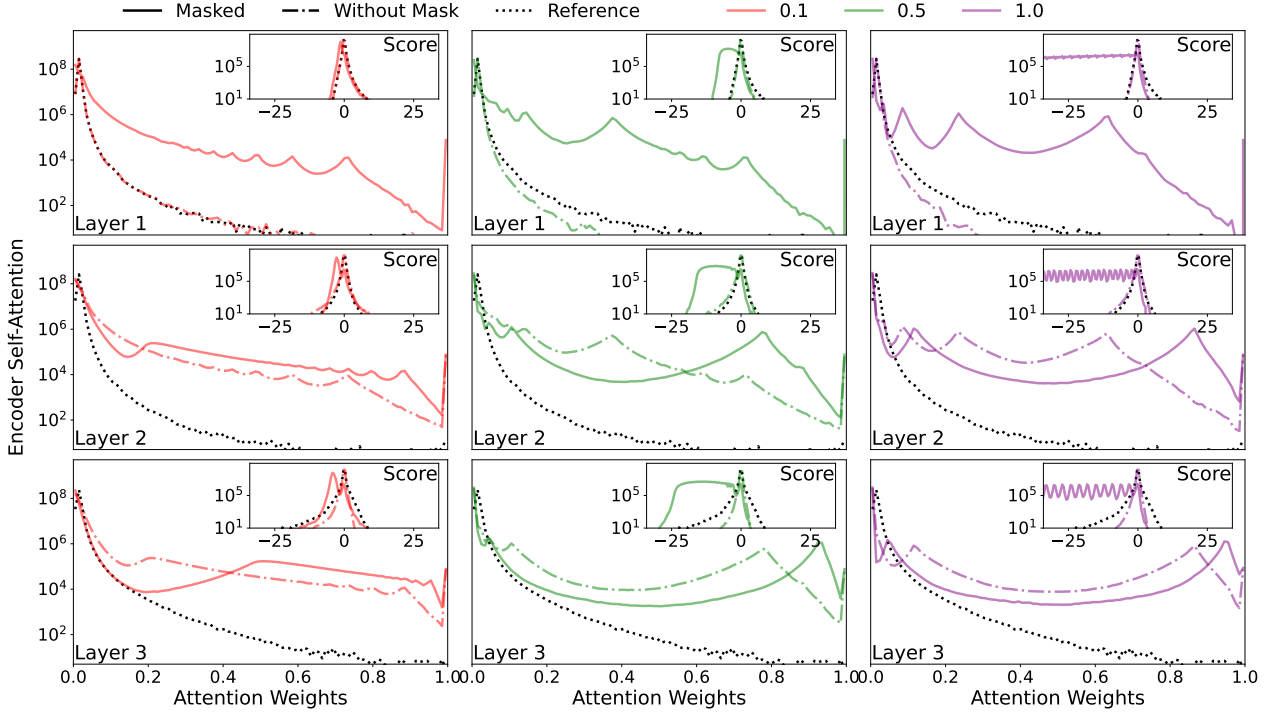
Figure A41: We present *Powerformer's* attention score and weight distributions on the ETTm1 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.



Figure A42: We present *Powerformer's* attention score and weight distributions on the ETTm1 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Figure A43: We present *Powerformer's* attention score and weight distributions on the ETTm2 dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{SPL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{SPL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.
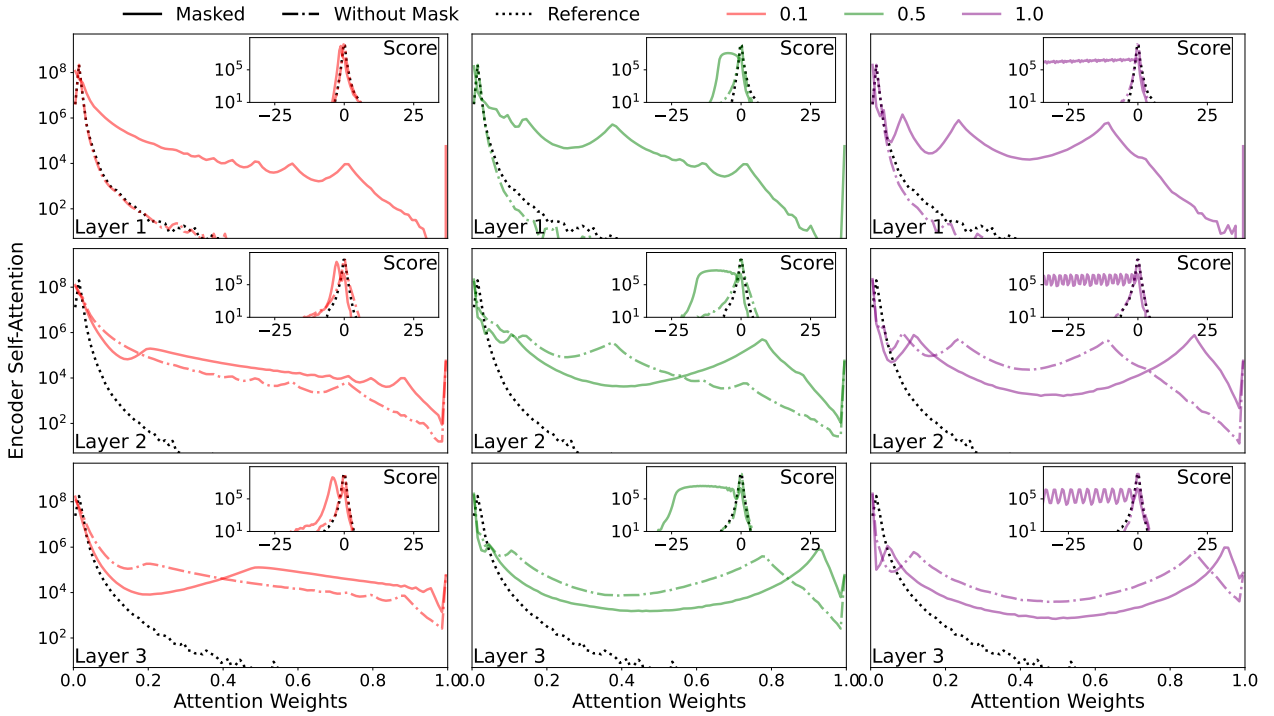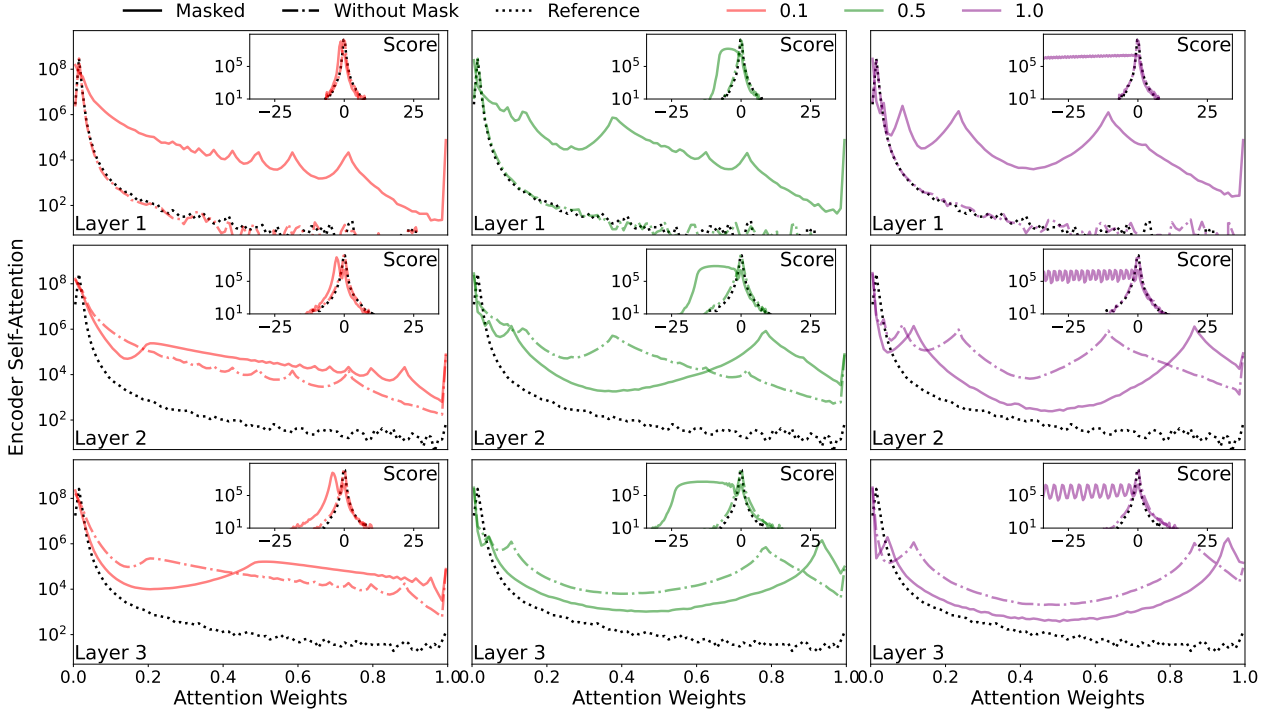


Figure A44: We present *Powerformer's* attention score and weight distributions on the ETTm2 dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{SPL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{SPL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.
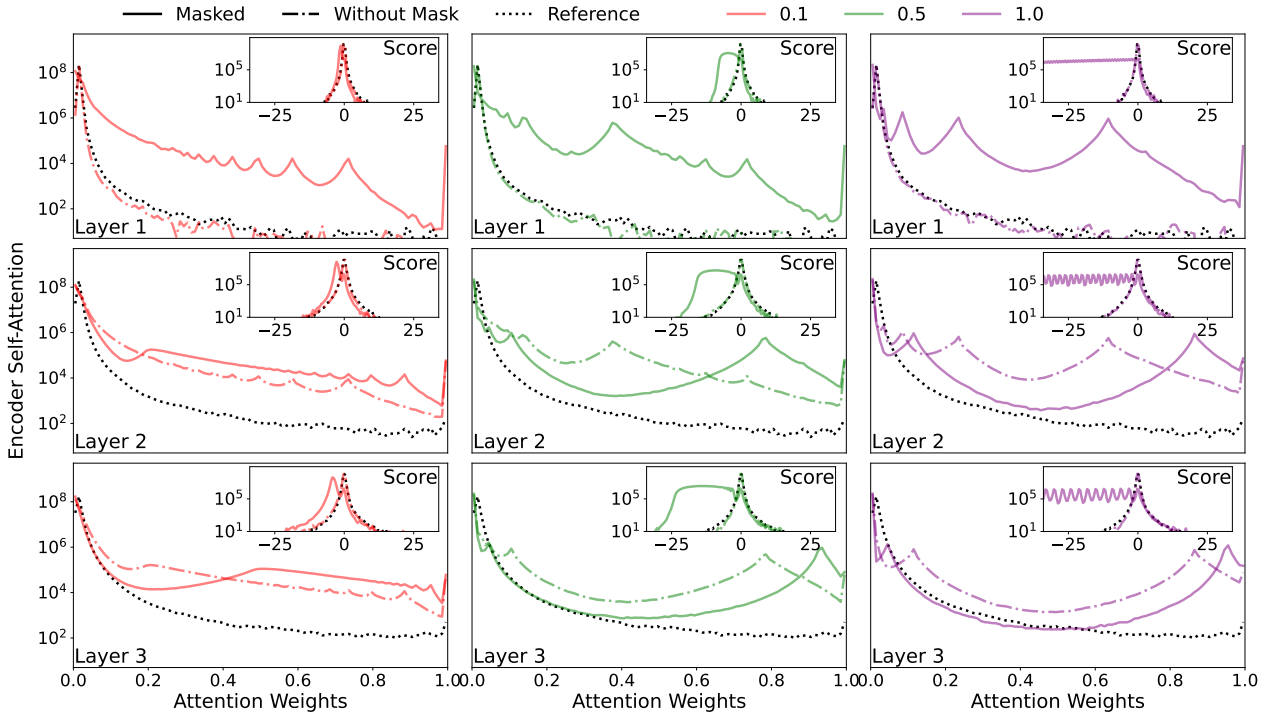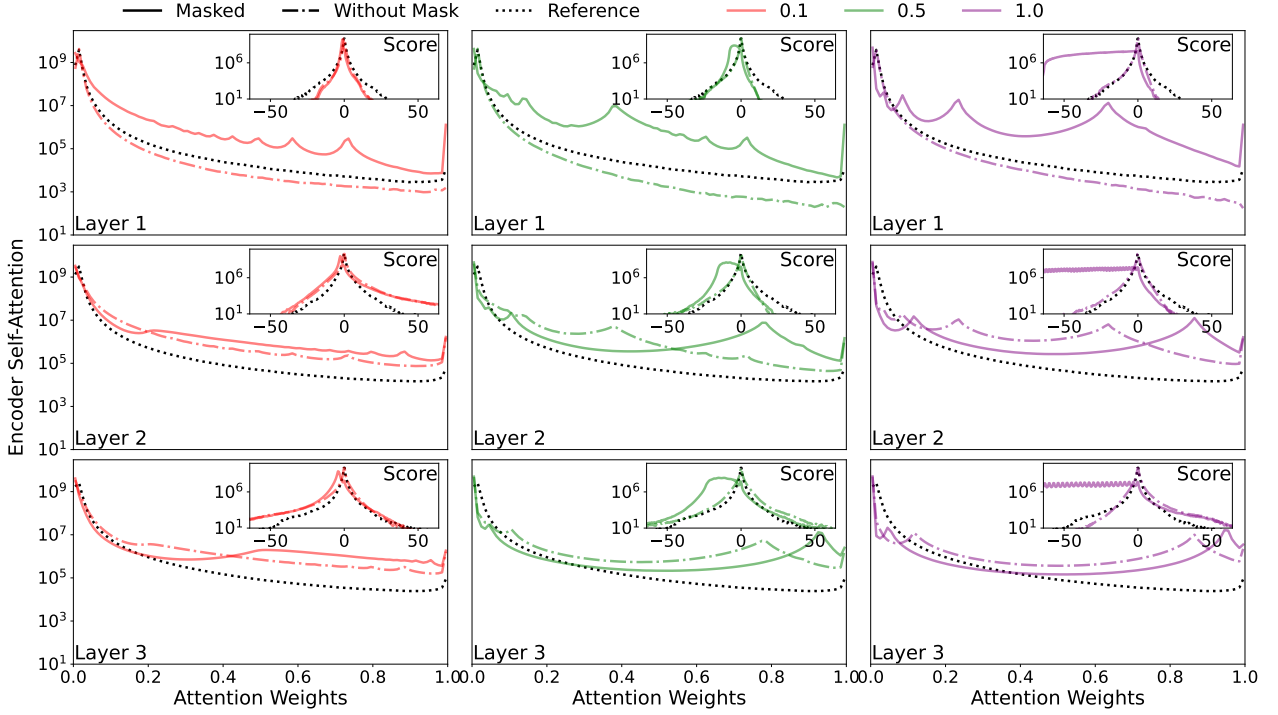
Figure A45: We present *Powerformer's* attention score and weight distributions on the Weather dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.



Figure A46: We present *Powerformer's* attention score and weight distributions on the Weather dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\mathrm{SPL})}(t)$ results before applying $\mathbf{M}^{(\mathrm{C,D})}$, and the solid lines represent WCMHA with $f^{(\mathrm{SPL})}(t)$ results after applying $\mathbf{M}^{(\mathrm{C,D})}$.

Figure A47: We present *Powerformer's* attention score and weight distributions on the Electricity dataset with a forecast and input length of 96 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{SPL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{SPL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.
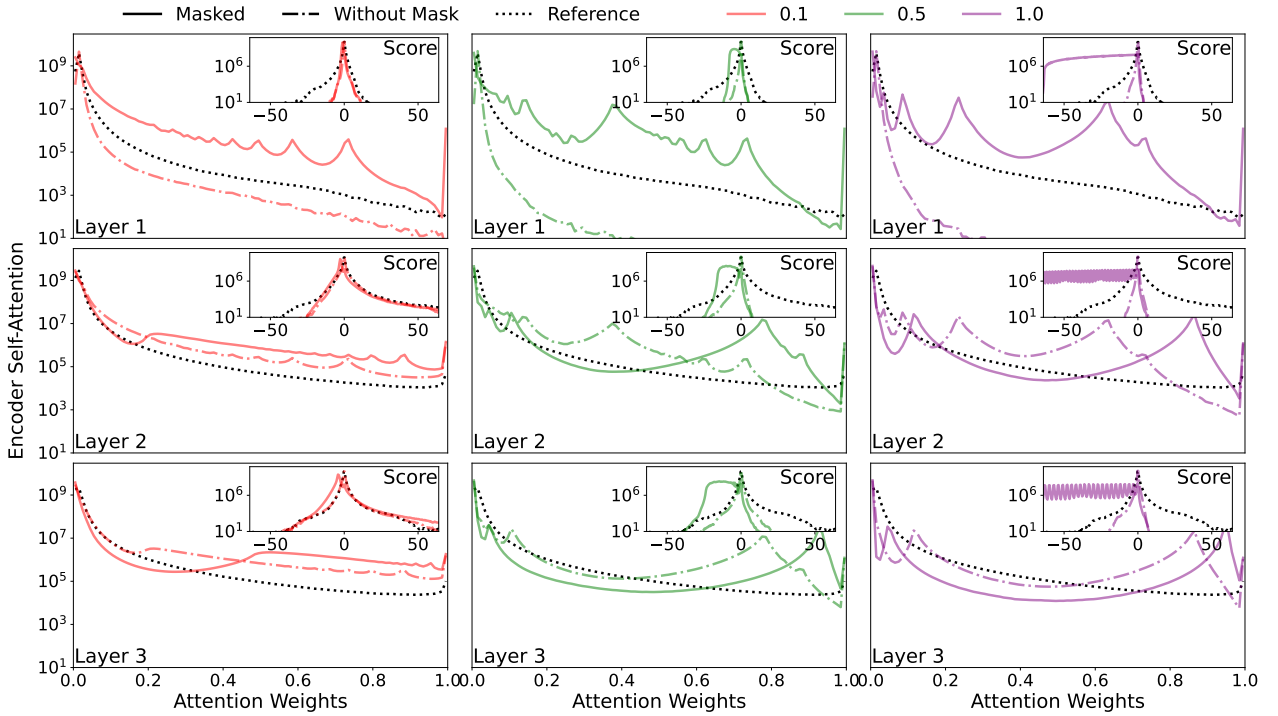


Figure A48: We present *Powerformer's* attention score and weight distributions on the Electricity dataset with a forecast and input length of 720 and 512, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{SPL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{SPL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.

Figure A49: We present *Powerformer's* attention score and weight distributions on the Traffic dataset with a forecast and input length of 96 and 336, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{SPL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{SPL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.
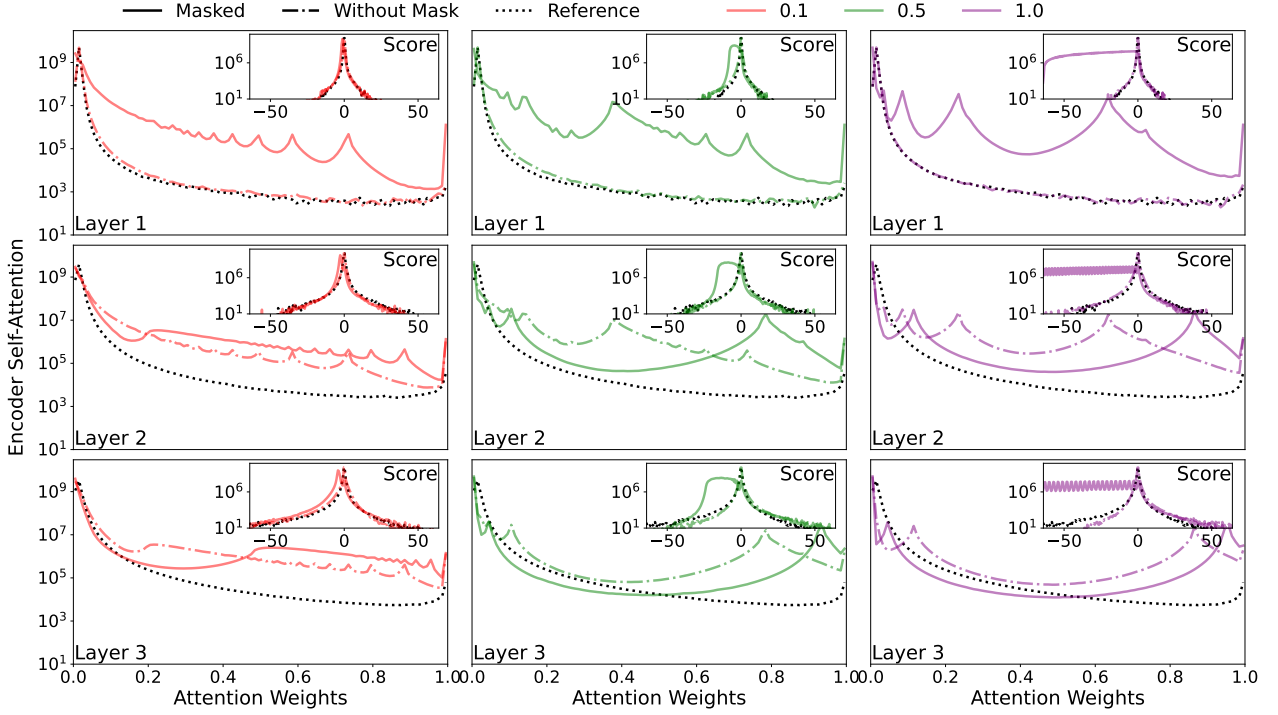


Figure A50: We present *Powerformer's* attention score and weight distributions on the Traffic dataset with a forecast and input length of 720 and 336, respectively. The dotted line represents the reference MHA results, the dashed-dotted line represents WCMHA with $f^{(\text{SPL})}(t)$ results before applying $\mathbf{M}^{(\text{C,D})}$, and the solid lines represent WCMHA with $f^{(\text{SPL})}(t)$ results after applying $\mathbf{M}^{(\text{C,D})}$.
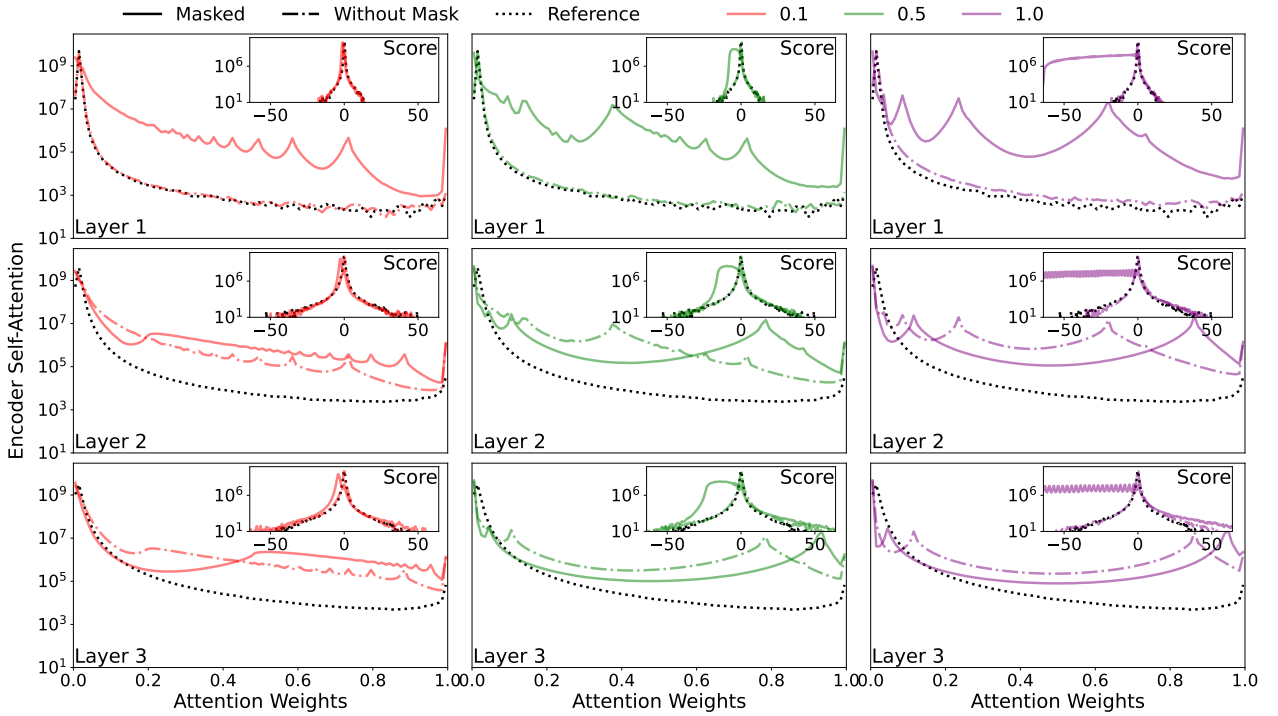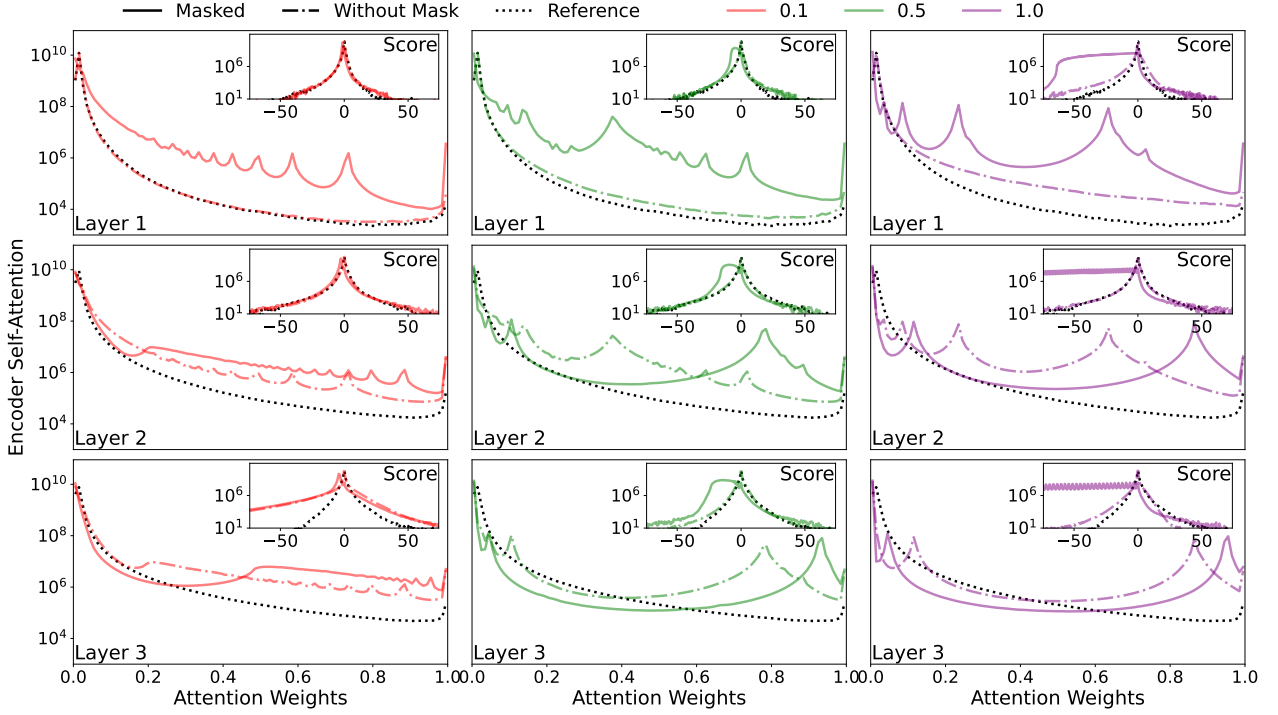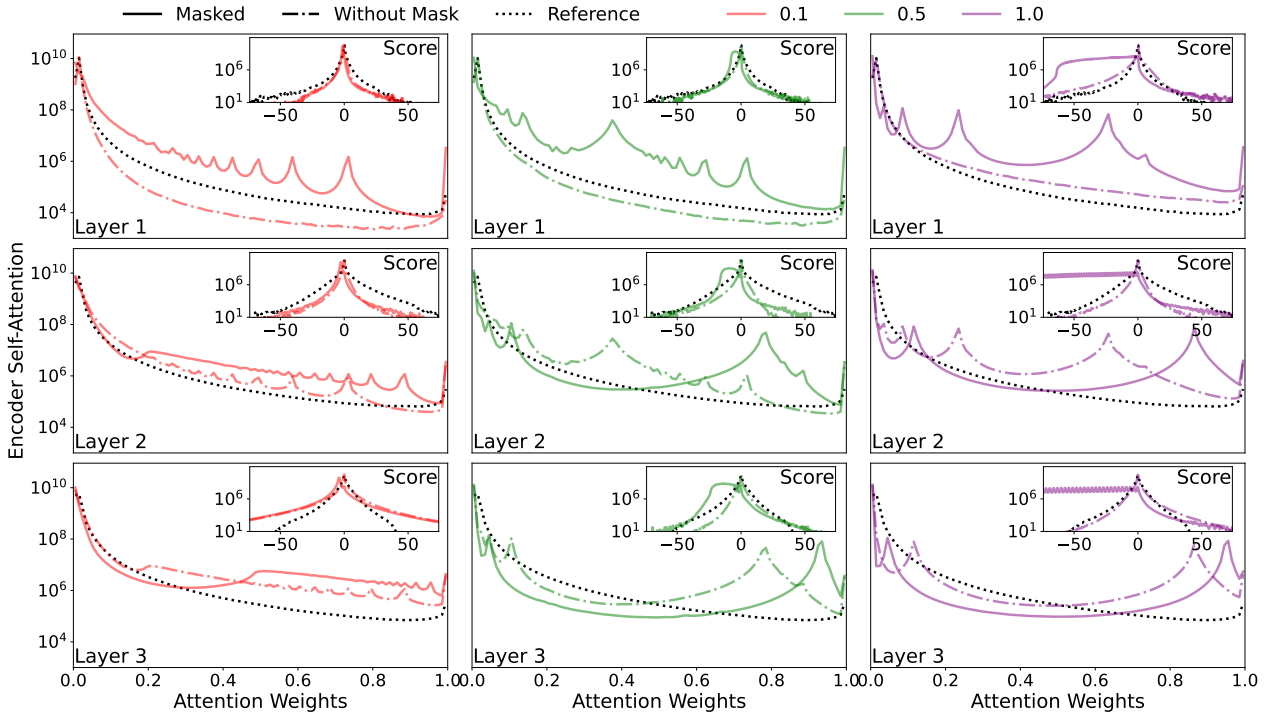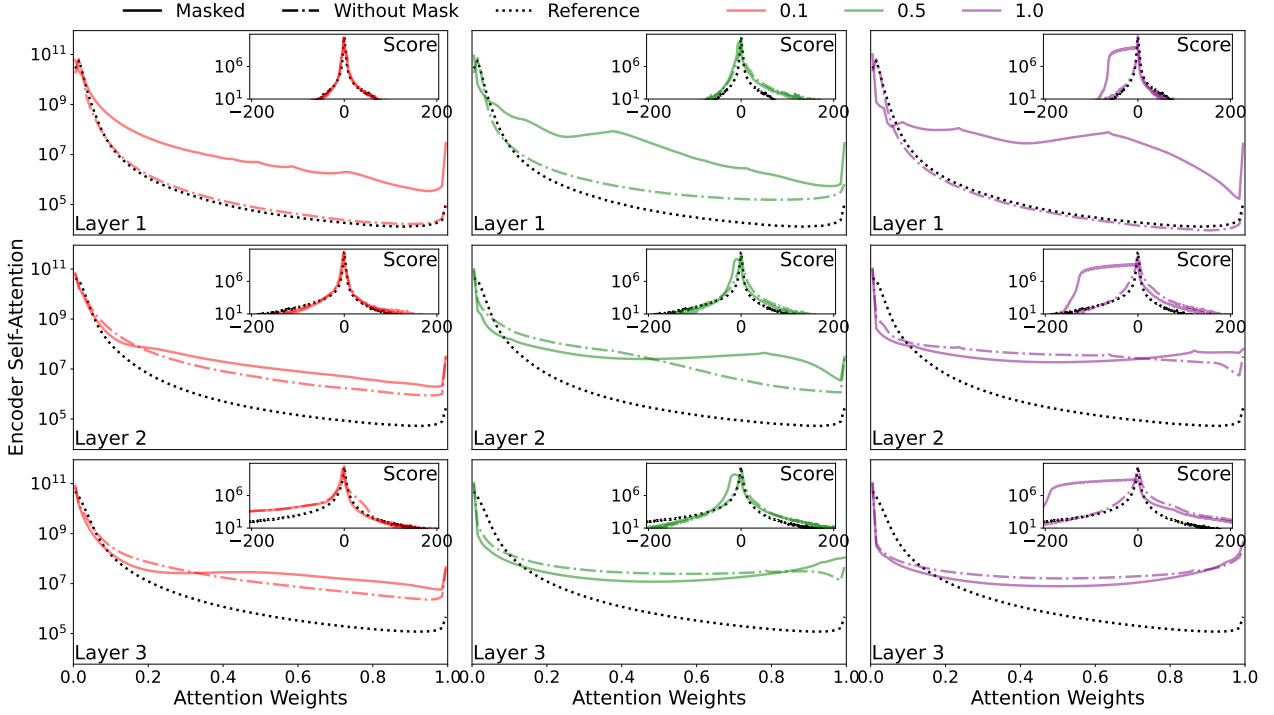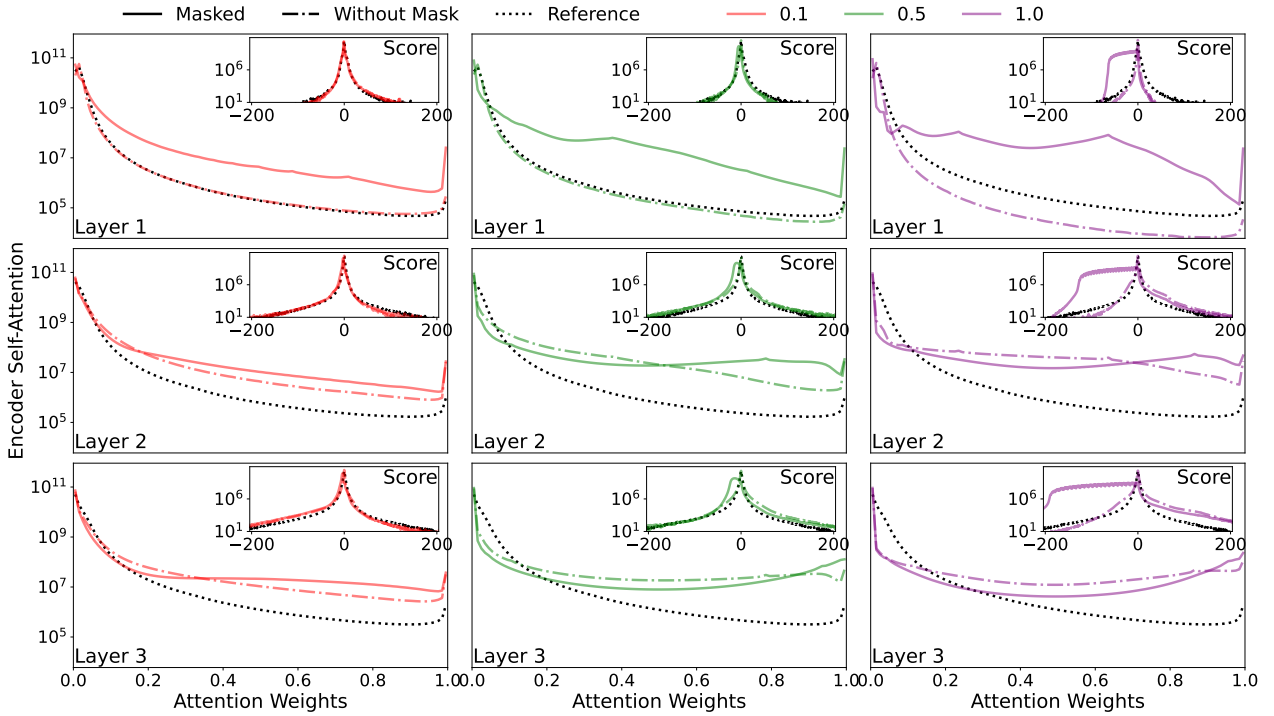
Table A7: We compare *Powerformer's* performance with WCMHA and the order 1 Butterworth Filter mask $f_1^{(\mathrm{BW})}(t)$ on standard time-series datasets for varying decay lengths. The best results are bolded and second best are underlined.

| | Delay | 2 | | 5 | | 10 | | 15 | | 20 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Metric | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 96 | <u>0.378</u> | <u>0.402</u> | **0.377** | <u>0.402</u> | **0.377** | <u>0.402</u> | **0.377** | **0.401** | **0.377** | **0.401** |
| | 192 | <u>0.415</u> | <u>0.422</u> | **0.414** | **0.421** | **0.414** | **0.421** | **0.414** | **0.421** | **0.414** | **0.421** |
| | 336 | <u>0.425</u> | **0.429** | <u>0.425</u> | **0.429** | <u>0.425</u> | **0.429** | <u>0.425</u> | **0.429** | **0.424** | **0.429** |
| | 720 | 0.442 | 0.459 | **0.439** | <u>0.458</u> | **0.439** | **0.457** | <u>0.440</u> | <u>0.458</u> | **0.439** | <u>0.458</u> |
| ETTh2 | 96 | <u>0.275</u> | **0.336** | <u>0.275</u> | **0.336** | **0.274** | **0.336** | **0.274** | **0.336** | **0.274** | **0.336** |
| | 192 | **0.338** | **0.379** | **0.338** | **0.379** | <u>0.339</u> | **0.379** | <u>0.339</u> | **0.379** | <u>0.339</u> | **0.379** |
| | 336 | <u>0.331</u> | <u>0.386</u> | <u>0.331</u> | <u>0.386</u> | <u>0.331</u> | **0.385** | <u>0.331</u> | **0.385** | **0.330** | **0.385** |
| | 720 | **0.380** | <u>0.422</u> | **0.380** | **0.421** | <u>0.382</u> | <u>0.422</u> | <u>0.382</u> | 0.424 | **0.380** | <u>0.422</u> |
| ETTm1 | 96 | **0.291** | **0.343** | <u>0.293</u> | **0.343** | 0.296 | **0.343** | 0.295 | <u>0.345</u> | 0.296 | 0.346 |
| | 192 | **0.329** | **0.369** | <u>0.333</u> | <u>0.371</u> | 0.337 | 0.372 | 0.338 | 0.373 | 0.336 | 0.373 |
| | 336 | <u>0.360</u> | **0.389** | **0.359** | <u>0.390</u> | 0.369 | 0.396 | 0.363 | 0.393 | 0.365 | 0.393 |
| | 720 | **0.414** | 0.426 | 0.425 | 0.429 | 0.418 | <u>0.425</u> | 0.417 | **0.424** | <u>0.416</u> | <u>0.425</u> |
| ETTm2 | 96 | <u>0.164</u> | **0.254** | **0.163** | **0.254** | 0.165 | <u>0.255</u> | 0.166 | 0.256 | 0.166 | <u>0.255</u> |
| | 192 | <u>0.223</u> | <u>0.294</u> | <u>0.223</u> | <u>0.294</u> | **0.222** | **0.293** | **0.222** | **0.293** | **0.222** | **0.293** |
| | 336 | 0.279 | 0.330 | 0.278 | 0.330 | **0.275** | **0.327** | <u>0.277</u> | 0.331 | <u>0.277</u> | <u>0.329</u> |
| | 720 | 0.358 | <u>0.381</u> | <u>0.353</u> | <u>0.381</u> | **0.351** | <u>0.381</u> | 0.359 | **0.380** | 0.359 | **0.380** |
| Weather | 96 | 0.155 | 0.204 | 0.154 | <u>0.203</u> | <u>0.153</u> | **0.202** | <u>0.153</u> | <u>0.203</u> | **0.152** | **0.202** |
| | 192 | 0.199 | **0.244** | <u>0.198</u> | **0.244** | <u>0.198</u> | **0.244** | <u>0.198</u> | <u>0.245</u> | **0.197** | <u>0.245</u> |
| | 336 | <u>0.250</u> | <u>0.283</u> | **0.248** | **0.282** | <u>0.250</u> | 0.284 | **0.248** | <u>0.283</u> | **0.248** | <u>0.283</u> |
| | 720 | <u>0.320</u> | <u>0.334</u> | **0.319** | **0.333** | 0.325 | 0.338 | 0.322 | 0.336 | **0.319** | 0.335 |

Table A8: We compare *Powerformer's* performance with WCMHA and the Butterworth filter mask $f_2^{(\mathrm{BW})}(t)$ on standard time-series datasets for varying decay lengths. The best results are bolded and second best are underlined.

| Delay | | 2 | | 5 | | 10 | | 15 | | 20 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Metric | | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 96 | <u>0.378</u> | <u>0.402</u> | **0.377** | <u>0.402</u> | **0.377** | **0.401** | **0.377** | **0.401** | **0.377** | **0.401** |
| | 192 | **0.414** | <u>0.422</u> | **0.414** | **0.421** | **0.414** | **0.421** | **0.414** | **0.421** | **0.414** | **0.421** |
| | 336 | <u>0.425</u> | **0.429** | <u>0.425</u> | **0.429** | <u>0.425</u> | **0.429** | **0.424** | **0.429** | **0.424** | **0.429** |
| | 720 | <u>0.442</u> | 0.459 | **0.439** | **0.457** | **0.439** | **0.457** | **0.439** | <u>0.458</u> | **0.439** | <u>0.458</u> |
| ETTh2 | 96 | <u>0.275</u> | **0.336** | <u>0.275</u> | **0.336** | **0.274** | **0.336** | **0.274** | **0.336** | **0.274** | **0.336** |
| | 192 | **0.338** | **0.379** | **0.338** | **0.379** | **0.338** | **0.379** | <u>0.339</u> | **0.379** | <u>0.339</u> | **0.379** |
| | 336 | <u>0.331</u> | 0.386 | <u>0.331</u> | 0.386 | <u>0.331</u> | <u>0.385</u> | **0.330** | <u>0.385</u> | **0.330** | **0.384** |
| | 720 | **0.380** | <u>0.422</u> | **0.380** | **0.421** | <u>0.381</u> | **0.421** | <u>0.381</u> | 0.423 | **0.380** | <u>0.422</u> |
| ETTm1 | 96 | **0.290** | **0.343** | <u>0.293</u> | <u>0.344</u> | 0.295 | 0.345 | 0.295 | 0.346 | 0.297 | <u>0.344</u> |
| | 192 | **0.329** | **0.370** | <u>0.333</u> | <u>0.371</u> | 0.339 | 0.373 | 0.337 | 0.373 | 0.336 | 0.373 |
| | 336 | **0.361** | <u>0.391</u> | **0.361** | <u>0.391</u> | <u>0.362</u> | 0.393 | 0.365 | 0.395 | 0.363 | **0.390** |
| | 720 | **0.413** | 0.425 | 0.422 | 0.428 | <u>0.415</u> | 0.423 | 0.417 | <u>0.422</u> | **0.413** | **0.421** |
| ETTm2 | 96 | <u>0.164</u> | <u>0.254</u> | <u>0.164</u> | <u>0.254</u> | 0.165 | **0.253** | **0.163** | **0.253** | 0.165 | 0.255 |
| | 192 | <u>0.223</u> | <u>0.294</u> | <u>0.223</u> | <u>0.294</u> | **0.222** | **0.293** | **0.222** | **0.293** | **0.222** | **0.293** |
| | 336 | 0.279 | <u>0.330</u> | 0.278 | <u>0.330</u> | **0.276** | 0.331 | **0.276** | **0.329** | <u>0.277</u> | **0.329** |
| | 720 | 0.358 | <u>0.381</u> | <u>0.354</u> | <u>0.381</u> | **0.353** | 0.382 | 0.359 | **0.380** | 0.360 | **0.380** |
| Weather | 96 | 0.156 | 0.204 | 0.154 | <u>0.203</u> | 0.154 | <u>0.203</u> | <u>0.153</u> | **0.202** | **0.152** | <u>0.203</u> |
| | 192 | 0.199 | **0.244** | 0.198 | **0.244** | 0.198 | **0.244** | <u>0.197</u> | **0.244** | **0.196** | **0.244** |
| | 336 | 0.250 | <u>0.283</u> | **0.248** | **0.282** | <u>0.249</u> | 0.284 | **0.248** | <u>0.283</u> | **0.248** | 0.284 |
| | 720 | 0.320 | **0.334** | **0.318** | **0.334** | 0.326 | 0.339 | 0.321 | 0.336 | <u>0.319</u> | <u>0.335</u> |

Table A9: We train *Powerformer* with $f^{(\mathrm{PL})}(t)$ and learnable time constant and present the results here. The base cases are the best results from *Powerformer* with constant decay weights: Table A5. The remaining columns designate the initialized time decay.

| Delay | | | Base Case | | 0.1 | | 0.25 | | 0.5 | | 0.75 | | 1 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Metric** | | | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE |
| ETTh1 | 336 | 96 | **0.376** | **0.401** | 0.377 | **0.401** | 0.377 | **0.401** | 0.377 | **0.401** | **0.376** | **0.401** | **0.376** | **0.401** |
| | | 192 | **0.413** | **0.421** | 0.413 | **0.421** | 0.413 | **0.421** | 0.413 | **0.421** | **0.413** | **0.421** | **0.413** | **0.421** |
| | | 336 | **0.424** | **0.430** | 0.425 | **0.430** | 0.425 | **0.430** | 0.425 | **0.430** | 0.425 | **0.430** | 0.425 | **0.430** |
| | | 720 | **0.437** | **0.455** | **0.437** | **0.455** | **0.437** | **0.455** | **0.437** | **0.455** | **0.437** | 0.456 | 0.438 | 0.456 |
| | 512 | 96 | **0.369** | **0.399** | **0.369** | **0.399** | **0.369** | **0.399** | 0.370 | **0.399** | 0.370 | **0.399** | 0.370 | **0.399** |
| | | 192 | **0.402** | **0.418** | 0.404 | 0.420 | 0.404 | 0.420 | **0.402** | **0.418** | 0.403 | **0.418** | 0.403 | **0.418** |
| | | 336 | **0.414** | **0.428** | **0.414** | **0.428** | **0.414** | **0.428** | 0.415 | 0.429 | 0.415 | 0.430 | 0.416 | 0.431 |
| | | 720 | **0.440** | **0.460** | **0.440** | **0.460** | **0.440** | **0.460** | **0.440** | **0.460** | 0.443 | 0.462 | 0.443 | **0.460** |
| ETTh2 | 336 | 96 | **0.274** | **0.335** | **0.274** | **0.335** | **0.274** | **0.335** | 0.275 | 0.336 | 0.275 | **0.335** | 0.275 | **0.335** |
| | | 192 | **0.338** | **0.379** | 0.341 | 0.381 | **0.338** | **0.379** | 0.340 | **0.379** | 0.340 | **0.379** | 0.341 | **0.379** |
| | | 336 | **0.325** | **0.379** | **0.325** | **0.379** | 0.326 | 0.380 | 0.331 | 0.384 | 0.333 | 0.383 | 0.334 | 0.383 |
| | | 720 | **0.376** | **0.419** | 0.377 | 0.420 | 0.379 | 0.421 | 0.381 | 0.423 | 0.381 | 0.422 | 0.381 | 0.422 |
| | 512 | 96 | **0.274** | **0.337** | **0.274** | **0.337** | **0.274** | **0.337** | **0.274** | **0.337** | 0.275 | **0.337** | 0.275 | 0.338 |
| | | 192 | **0.340** | **0.381** | **0.340** | **0.381** | **0.340** | **0.381** | 0.341 | 0.382 | 0.342 | 0.382 | 0.342 | 0.382 |
| | | 336 | **0.330** | **0.387** | **0.330** | **0.387** | 0.331 | **0.387** | 0.333 | 0.388 | 0.334 | 0.388 | 0.334 | 0.388 |
| | | 720 | **0.383** | **0.425** | **0.383** | 0.426 | 0.384 | 0.426 | **0.383** | 0.426 | 0.384 | 0.426 | **0.383** | **0.425** |
| ETTm1 | 336 | 96 | **0.290** | **0.342** | 0.292 | 0.343 | **0.290** | **0.342** | 0.292 | 0.343 | 0.293 | 0.344 | 0.292 | 0.345 |
| | | 192 | **0.330** | **0.369** | 0.331 | 0.371 | 0.334 | 0.372 | 0.332 | 0.371 | 0.332 | 0.370 | **0.330** | **0.369** |
| | | 336 | **0.359** | **0.389** | 0.362 | 0.392 | 0.361 | 0.391 | 0.362 | 0.390 | **0.359** | **0.389** | 0.361 | 0.390 |
| | | 720 | **0.412** | **0.421** | 0.413 | 0.423 | 0.413 | 0.422 | 0.413 | **0.421** | 0.413 | **0.421** | **0.412** | 0.423 |
| | 512 | 96 | **0.290** | **0.345** | 0.291 | **0.345** | **0.290** | **0.345** | **0.290** | **0.345** | 0.291 | 0.346 | 0.292 | 0.346 |
| | | 192 | **0.329** | **0.368** | 0.331 | 0.370 | 0.330 | 0.370 | **0.329** | **0.368** | 0.332 | 0.370 | 0.332 | 0.372 |
| | | 336 | **0.361** | **0.390** | 0.362 | 0.391 | 0.362 | **0.390** | 0.361 | 0.392 | 0.362 | **0.390** | 0.363 | **0.390** |
| | | 720 | **0.415** | **0.423** | 0.417 | 0.428 | **0.415** | 0.425 | 0.417 | **0.423** | 0.420 | 0.426 | 0.416 | 0.431 |
| ETTm2 | 336 | 96 | **0.162** | **0.252** | **0.162** | **0.252** | 0.165 | 0.255 | 0.166 | 0.256 | 0.164 | 0.254 | 0.164 | 0.254 |
| | | 192 | **0.222** | **0.292** | **0.222** | **0.292** | **0.222** | 0.293 | **0.222** | 0.293 | **0.222** | 0.293 | **0.222** | 0.293 |
| | | 336 | **0.277** | **0.329** | **0.277** | **0.329** | **0.277** | **0.329** | **0.277** | **0.329** | **0.277** | **0.329** | 0.278 | 0.330 |
| | | 720 | **0.359** | **0.380** | 0.363 | 0.381 | 0.363 | 0.381 | 0.362 | 0.382 | **0.359** | **0.380** | **0.359** | 0.384 |
| | 512 | 96 | **0.164** | **0.255** | 0.165 | **0.255** | 0.165 | **0.255** | 0.165 | 0.256 | 0.165 | 0.257 | **0.164** | **0.255** |
| | | 192 | **0.221** | **0.294** | 0.222 | 0.295 | 0.224 | 0.297 | 0.224 | 0.296 | 0.223 | 0.296 | **0.221** | **0.294** |
| | | 336 | **0.272** | **0.327** | 0.274 | 0.329 | 0.274 | 0.329 | 0.274 | 0.329 | **0.272** | **0.327** | 0.273 | **0.327** |
| | | 720 | **0.356** | **0.381** | 0.358 | 0.383 | 0.358 | 0.383 | 0.358 | 0.382 | 0.358 | 0.383 | **0.356** | **0.381** |
| Weather | 336 | 96 | **0.151** | **0.200** | **0.151** | 0.201 | **0.151** | **0.200** | **0.151** | **0.200** | 0.153 | 0.201 | 0.154 | 0.202 |
| | | 192 | 0.196 | **0.241** | 0.196 | 0.244 | 0.196 | 0.243 | 0.195 | 0.242 | 0.196 | **0.241** | 0.197 | 0.242 |
| | | 336 | 0.247 | **0.281** | 0.248 | 0.283 | 0.247 | 0.284 | 0.247 | 0.282 | **0.246** | **0.281** | 0.248 | 0.282 |
| | | 720 | **0.317** | 0.333 | 0.318 | 0.334 | 0.318 | 0.334 | 0.318 | 0.334 | 0.318 | 0.333 | **0.317** | **0.332** |
| | 512 | 96 | **0.147** | **0.197** | 0.148 | 0.198 | 0.148 | 0.199 | **0.147** | 0.198 | 0.148 | 0.198 | 0.149 | 0.199 |
| | | 192 | **0.191** | **0.239** | 0.193 | 0.241 | 0.194 | 0.242 | 0.193 | 0.240 | **0.191** | **0.239** | 0.193 | 0.241 |
| | | 336 | **0.243** | **0.279** | 0.244 | 0.281 | **0.243** | 0.281 | **0.243** | 0.280 | **0.243** | **0.279** | 0.244 | 0.280 |
| | | 720 | **0.311** | 0.329 | **0.311** | 0.330 | **0.311** | 0.330 | 0.314 | 0.331 | **0.311** | **0.328** | **0.311** | 0.329 |
| Electricity | 336 | 96 | **0.130** | **0.224** | **0.130** | **0.224** | 0.131 | **0.224** | 0.131 | **0.224** | 0.131 | **0.224** | 0.131 | **0.224** |
| | | 192 | **0.147** | **0.240** | 0.148 | 0.241 | **0.147** | **0.240** | **0.147** | **0.240** | 0.148 | **0.240** | **0.147** | 0.241 |
| | | 336 | **0.164** | **0.257** | **0.164** | 0.258 | **0.164** | 0.258 | **0.164** | 0.258 | **0.164** | 0.259 | **0.164** | 0.258 |
| | | 720 | 0.201 | **0.291** | 0.201 | 0.292 | **0.200** | **0.291** | 0.201 | 0.292 | 0.202 | 0.292 | 0.202 | 0.293 |
| | 512 | 96 | **0.129** | **0.223** | 0.130 | **0.223** | **0.129** | 0.224 | **0.129** | 0.224 | **0.129** | **0.223** | 0.130 | **0.223** |
| | | 192 | **0.146** | **0.240** | 0.147 | 0.241 | **0.146** | **0.240** | **0.146** | **0.240** | **0.146** | **0.240** | 0.147 | 0.241 |
| | | 336 | **0.162** | **0.257** | **0.162** | 0.258 | **0.162** | 0.258 | **0.162** | 0.258 | **0.162** | **0.257** | **0.162** | 0.258 |
| | | 720 | **0.198** | **0.289** | 0.199 | 0.290 | **0.198** | 0.290 | 0.199 | 0.290 | 0.200 | 0.291 | 0.200 | 0.291 |
| Traffic | 336 | 96 | **0.370** | **0.253** | **0.370** | **0.253** | **0.370** | **0.253** | 0.372 | 0.254 | 0.373 | 0.255 | 0.374 | 0.256 |
| | | 192 | **0.388** | **0.260** | **0.388** | 0.261 | **0.388** | **0.260** | 0.390 | 0.261 | 0.390 | 0.262 | 0.390 | 0.262 |
| | | 336 | **0.400** | **0.266** | **0.400** | 0.267 | **0.400** | **0.266** | 0.404 | 0.269 | 0.404 | 0.270 | 0.405 | 0.270 |
| | | 720 | 0.434 | **0.286** | 0.435 | **0.286** | **0.433** | **0.286** | 0.435 | 0.287 | 0.437 | 0.289 | 0.436 | 0.287 |
| | 512 | 96 | **0.365** | 0.252 | **0.365** | **0.251** | 0.366 | 0.252 | 0.366 | 0.253 | 0.367 | 0.254 | 0.394 | 0.282 |
| | | 192 | **0.382** | 0.258 | 0.390 | 0.269 | 0.383 | 0.261 | 0.383 | 0.261 | **0.382** | **0.260** | 0.407 | 0.287 |
| | | 336 | **0.391** | 0.265 | **0.391** | **0.264** | 0.392 | 0.265 | 0.392 | 0.265 | 0.393 | 0.265 | 0.394 | 0.265 |
| | | 720 | **0.430** | 0.286 | 0.433 | 0.288 | 0.431 | 0.288 | 0.432 | 0.287 | **0.430** | **0.285** | **0.430** | 0.286 |