

UNIFYING MULTI-SCALE DESIGN IN TIME-SERIES FORECASTING

Anonymous authors

Paper under double-blind review

ABSTRACT

Multi-scale modeling in time-series forecasting, which seeks to capture cross-scale relationships for modeling complex dependencies, is increasingly popular. While previous work lacks principled foundations, we unify existing scaling methods into a *scaling operator family*, providing a general theoretical basis for multi-scaling methods and revealing two key limitations of current models: **static** scaling and inflexible cross-scale modeling. To address these limitations, we propose SIGMA (Single Gaussian Multi-scale Architecture), a simple yet principled multi-scale framework. It enables **position-wise** scaling via the learnable discrete Gaussian (LDG) kernel grounded in scale-space theory, coupled with a lightweight MLP processor for efficient cross-scale interaction. We evaluate SIGMA comprehensively on long- and short-term forecasting benchmarks against state-of-the-art multi-scale baselines. SIGMA outperforms all competitors on both tasks, achieving the best performance in 55 out of 80 long-term evaluation settings. Beyond accuracy, SIGMA improves training speed by up to 3.8 times and reduces memory consumption by up to 5.3 times compared to the strongest competitors. Code is available at <https://anonymous.4open.science/r/SiGMA-ICLR2026>.

1 INTRODUCTION

Time-series forecasting is widely applied in diverse domains, including energy systems (Deb et al., 2017), climate and weather prediction (Dimri et al., 2020), and traffic management (Tedjopurnomo et al., 2020), to name a few. Despite its importance, forecasting remains fundamentally challenging due to the complex temporal structures of real-world time series (Kolambe, 2024). Inspired by the success of deep learning, recent work has introduced neural forecasting methods, including MLP-based methods (Ekambaram et al., 2023; Zeng et al., 2023; Yi et al., 2023; Challu et al., 2023) and Transformer-based methods (Nie et al., 2023; Liu et al., 2024; Shi et al., 2025).

Within this scope, multi-scale modeling has emerged as a powerful design principle. By constructing temporal dynamics at multiple resolutions, it explicitly models cross-scale interactions, thereby enhancing predictive performance. Specifically, coarse-grained scaling captures long-term dynamics, whereas fine-grained scaling preserves high-frequency fluctuations, yielding a more comprehensive representation (Chen et al., 2023). Building upon this principle, recent work has introduced diverse strategies from simple multi-rate sampling to hierarchical and frequency-aware mechanisms (Zhao et al., 2024; Wang et al., 2024, 2025; Naghashi et al., 2025; Hu et al., 2025a; Yang et al., 2025).

We visualize popular scaling operators in Figure 1. Each scaling operator takes a scale parameter s as an additional input to transform the given sequence, creating a representation with different length and resolution. Thus, the effectiveness of a multi-scale approach depends on not only the choice of the scaling function, but also its scale parameters and how to fuse multi-resolution representations.

However, existing methods exhibit two key limitations. (1) *Static scaling*. Most approaches employ a fixed scaling strategy (Wang et al., 2023, 2024; Hu et al., 2025a) for **different time steps**. This restricts the diversity of temporal views and limits the ability to capture time-varying characteristics inherent in real-world data. (2) *Inflexible cross-scale modeling*. Cross-scale interactions are modeled in a staged manner, processing each scale independently before combining them (Zhong et al., 2024; Wang et al., 2025). This procedure introduces sequential dependencies, reducing efficiency and limiting flexibility in modeling complex temporal relationships.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

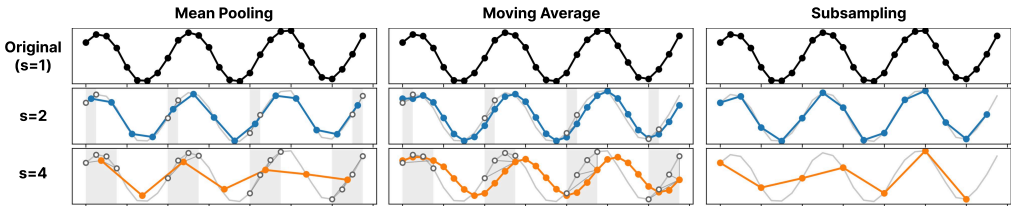


Figure 1: Examples of popular scaling operators used in existing methods. Each operator transforms the input into coarser representations, with a different scaling parameter s , capturing different levels of structure and abstraction. Refer to Table 1 for more information of existing scaling operators.

Before presenting our approach to alleviating these limitations, we introduce a novel and principled concept: a *family of scaling operators*, which unifies the scaling methods studied in the literature. We further demonstrate that any multi-scale forecasting model can be decomposed into two essential components: a *multi-scale generator* that creates multiple representations of the input series, and a *multi-scale processor* that fuses information across these representations. This abstraction provides a unified analytical framework for reinterpreting multi-scale time-series modeling.

Based on this unified view, we propose SIGMA (Single Gaussian Multi-scale Architecture), a simple yet principled framework for multi-scale modeling. It adopts *position-wise* scaling through the learnable discrete Gaussian (LDG) kernel, *where time-indexed scale parameters are learned from data* within a single operator grounded in scale-space theory. Then, a lightweight multi-layer perceptron (MLP) models cross-scale interactions in a single step, removing staged fusion. Together, these components provide a theoretically grounded and efficient approach that simplifies and strengthens multi-scale modeling.

We evaluate SIGMA extensively on long- and short-term forecasting benchmarks and observe consistent improvements over state-of-the-art multi-scale models. Especially on long-term tasks, SIGMA achieves the best performance in 55 out of 80 evaluation settings with up to 7.8% lower error than runner-up models. Beyond accuracy, SIGMA delivers clear efficiency gains, reducing memory usage by up to 3.8 \times and increasing training speed by 5.3 \times over the strongest baselines. We further support these findings with comprehensive ablation and case studies, confirming that SIGMA provides a principled, and efficient framework for multi-scale time-series modeling.

Our contributions are summarized as follows:

- **Unified framework:** We introduce the novel concept of a scaling operator family, which provides a principled foundation for multi-scale time-series modeling. It unifies existing methods under a generator-processor decomposition and clarifies their key limitations.
- **Simplified architecture:** We propose SIGMA, a simple yet principled multi-scale architecture. It enables *position-wise* scaling and efficient cross-scale fusion through the learnable discrete Gaussian (LDG) kernel with dynamic scaling parameters and a lightweight MLP processor.
- **Empirical validation:** We empirically show that SIGMA achieves state-of-the-art performance on extensive long- and short-term forecasting benchmarks while substantially reducing computational complexity compared to prior multi-scale models. We also conduct various ablation and case studies to provide deeper insights on the successful behavior of SIGMA.

2 PROBLEM DEFINITION AND RELATED WORK

Problem definition Given a time series, let $x \in \mathbb{R}^L$ denote an input sequence of length L from it, and $y \in \mathbb{R}^T$ denote the corresponding future trajectory of length T . We denote the input space as $\mathcal{X} \subset \mathbb{R}^L$ and the target space as $\mathcal{Y} \subset \mathbb{R}^T$. Our objective is to learn a time series forecasting model $h : \mathcal{X} \rightarrow \mathcal{Y}$ that maps a past window x to a T -step forecast of its future y .

Data assumption Throughout this work, we focus on non-trivial datasets whose forecastability ϕ is strictly less than 1. The forecastability of a time series is defined as one minus the spectral entropy

Table 1: Representative scaling operator families with associated scale parameters s . Each family simplifies the sequence in a structured manner as the scale parameter s increases, where $k \in [1, L_s]$ is the element index and L_s denotes the length of the output sequence. For wavelet decomposition, $f(\mathbf{x}|s)$ yields the approximation coefficients at level s using the scaling function ϕ .

Operator f	Equation	Scale parameter s	L_s	Usage
Average pooling	$[f(\mathbf{x} s)]_k = (1/s) \sum_{i=0}^{s-1} x_{ks+i}$	Window length	P/s	(Wang et al., 2024)
Max pooling	$[f(\mathbf{x} s)]_k = \max_{0 \leq i < s} x_{ks+i}$	Window length	P/s	(Challu et al., 2023)
Moving average	$[f(\mathbf{x} s)]_k = (1/s) \sum_{i=0}^{s-1} x_{k-i}$	Window length	P	(Wang et al., 2023)
Subsampling	$[f(\mathbf{x} s)]_k = x_{ks}$	Stride	P/s	(Liu et al., 2022a)
Segmentation	$[f(\mathbf{x} s)]_k = x_k$	Length divisor	P/s	(Wu et al., 2023)
Wavelet decomposition	$[f(\mathbf{x} s)]_k = \langle \mathbf{x}, \phi_{s,k} \rangle$	Decomposition level	$P/2^s$	(Murad et al., 2025)

of its Fourier decomposition (Goerg, 2013), formally given as

$$\phi(\mathbf{x}) = 1 - \frac{H(\mathbf{x})}{\log(2\pi)} \in [0, 1], \quad H(\mathbf{x}) = - \int_{-\pi}^{\pi} p_{\mathbf{x}}(\omega) \log p_{\mathbf{x}}(\omega) d\omega, \quad (1)$$

where $p_{\mathbf{x}}(\omega)$ denotes the normalized spectral density of \mathbf{x} on $[-\pi, \pi]$, characterizing how concentrated the spectrum of \mathbf{x} is. Intuitively, larger forecastability ϕ indicates a more concentrated spectrum and thus a stronger predictable structure (Wang et al., 2024; 2025). Most real-world datasets including every benchmark used in our experiments satisfy $\phi < 1$ (see Appendix B). This serves as a mild theoretical assumption and does not restrict the practical applicability of our framework.

Time series forecasting The field of time series forecasting has seen a significant evolution from traditional statistical methods to modern deep learning architectures. While statistical models like ARIMA (Box et al., 2015) and state-space models (Durbin & Koopman, 2012) are still used, they are often limited in their ability to capture complex nonlinear and long-range dependencies (Zhou et al., 2021). Among deep forecasting models that learn representations directly from raw sequences, multi-layer perceptron architectures have recently been recognized for their effectiveness and efficiency (Ekambaram et al., 2023; Zeng et al., 2023; Yi et al., 2023; Challu et al., 2023). Concurrently, Transformer-based approaches have demonstrated strong empirical results, particularly for long-horizon forecasting tasks (Nie et al., 2023; Liu et al., 2024; Shi et al., 2025).

Multi-scale design in time series Time series often exhibit complex patterns across diverse temporal scales, from rapid fluctuations to long-term trends. To capture such heterogeneity, recent work incorporates explicit multi-scale mechanisms to disentangle and integrate temporal dependencies. These methods include hierarchical decompositions that construct multi-resolution representations via downsampling (Liu et al., 2022b; Challu et al., 2023), frequency-domain or wavelet-based analyses that emphasize periodic structure through spectral transforms (Wu et al., 2023; Cai et al., 2024; Murad et al., 2025), and aggregation schemes that fuse representations across scales (Wang et al., 2023; 2024; 2025). However, existing multi-scale methods rely on fixed scaling and staged cross-scale fusion, limiting adaptability to diverse temporal patterns and flexibility in modeling cross-scale interactions. This highlights the need for more flexible multi-scale modeling.

3 FOUNDATIONS OF SCALING IN TIME SERIES

In time series analysis, *scaling* denotes transformations that change the temporal or frequency resolution of data. As illustrated in Figure 1, various operations such as downsampling and moving averages are commonly employed in previous work. However, these techniques have largely been treated as ad-hoc design choices, lacking a unifying theoretical principle. We address this gap by introducing a rigorous definition of scaling operator families. To the best of our knowledge, this is the first work that establishes a unified mathematical foundation for analyzing scaling operations in time series. Detailed proofs of the theorems are provided in Appendix A.

Definition 3.1 (Scaling operator family). *A set of transformations $\mathcal{F} = \{f(\mathbf{x}|s) \mid s \in \mathbb{Z}_+\}$ with $f(\cdot|s) : \mathcal{X} \rightarrow \mathcal{X}_s \subset \mathbb{R}^{L_s}$ conditioned on $s \in \mathbb{Z}_+$ is a scaling operator family if it satisfies*

- *Non-expansiveness:* For $s_i \in \mathbb{Z}_+$ and $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, $\|f(\mathbf{x}|s_i) - f(\mathbf{x}'|s_i)\|_2 \leq \|\mathbf{x} - \mathbf{x}'\|_2$.
- *Energy reduction:* For any $s_i, s_j \in \mathbb{Z}_+$ such that $s_i = m \cdot s_j$ for some $m > 1$, $\|f(\mathbf{x}|s_i)\|_2 \leq \|f(\mathbf{x}|s_j)\|_2$ for all $\mathbf{x} \in \mathcal{X}$ and $\|f(\mathbf{x}'|s_i)\|_2 < \|f(\mathbf{x}'|s_j)\|_2$ for some $\mathbf{x}' \in \mathcal{X}$.

The scale parameter s determines the level of scaling, and also the output dimension L_s .

Theorem 3.2. Common sequence operations such as max/mean/min-pooling (Zheng et al., 2014), moving averages (Box et al., 2015), subsampling (Hannan, 2009), segmentation (Keogh et al., 2004) and wavelet decompositions (Percival & Walden, 2000) are scaling operator families.

Definition 3.1 formalizes the expected behavior of scaling. Scaling operators reduce or at least preserve the distances between the resulting sequences, as they are designed to capture diverse properties from a sequence but without adding external information. At the same time, the *energy reduction* property formalizes the intuition that time series with coarser scales should be simpler than finer ones along multiplicative chains of scale parameters: as the scale parameter s increases, the result exhibits less overall energy, with small fluctuations being smoothed out.

Table 1 summarizes the operator families given in Theorem 3.2 with their associated scale parameters. These operators can be grouped into two categories based on how they transform the given sequence: (i) *smoothing operators*, which attenuate fine-grained variability by summarizing adjacent observations (e.g., pooling and moving average), and (ii) *structural operators*, which progressively decompose or resample the signals by enforcing a specific structure (e.g., subsampling and wavelet decomposition). Previous works have typically used the term *scaling* only for one group, but our unified Definition 3.1 suggests that both categories of operators can be understood in the same lens in terms of the non-expansiveness and *energy reduction* properties.

To empirically verify that the operator families in Table 1 satisfy the requirements of Definition 3.1, we evaluate their average pairwise contraction and induced average energy on the Traffic dataset. As shown in Figure 2, all operator families exhibit strictly smaller output differences than input differences, confirming non-expansiveness, and their energies decrease monotonically over multiplicative scales, confirming energy reduction. These observations demonstrate that the operator families in Table 1 satisfy both conditions of Definition 3.1 on *real-world datasets*. Additional experiments on other datasets are provided in Appendix C.

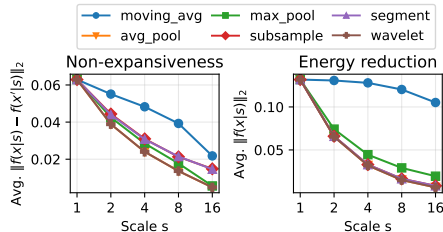


Figure 2: The non-expansiveness and energy reduction of the six scaling operator families on the Traffic dataset.

Theorem 3.3. Trivial operations such as constant mappings, permutations, additive shifts, scalar multiplications, and general linear transformations do not form scaling operator families.

Theorem 3.3 highlights that not all parameterized transformations qualify as scaling operator families. For example, while a scalar multiplication can change the scale of each observation, it does not guarantee *energy* reduction and thus falls outside our definition. Putting Theorem 3.2 and Theorem 3.3 together, we claim that our definition of scaling operator families is designed carefully to include only the operations that have been considered as “scaling” in the literature, to provide more interpretable and theoretically grounded approaches to multi-scale time series modeling.

4 REVISITING MULTI-SCALE METHODS IN TIME-SERIES FORECASTING

Based on the definition of scaling operator families, we unify existing multi-scale time series forecasting methods and analyze their limitations from the perspective of unified modules.

4.1 GENERALIZATION OF MULTI-SCALE METHODS

We formally define multi-scale time series modeling as the composition of two fundamental components: *multi-scale generation* and *multi-scale processing*. For an input sequence $\mathbf{x} \in \mathbb{R}^L$, scaling parameters $\mathbf{s} = (s_1, \dots, s_M)$ where M is the number of scaling parameters, and a scaling operator family \mathcal{F} , a *multi-scale generator* g is defined as

$$g(\mathbf{x}|\mathbf{s}) = (f(\mathbf{x}|s_1), \dots, f(\mathbf{x}|s_M)), \quad (2)$$

which produces a collection of representations of \mathbf{x} across the scaling parameters in \mathbf{s} .

A *multi-scale processor* p_θ is a function that integrates the representations $g(\mathbf{x}|s)$ to capture meaningful cross-scale interactions with its learnable parameters θ . For example, some approaches employ weighted summation and aggregation modules to integrate information (Wu et al., 2023; Cai et al., 2024), while others adopt cross-scale mixing mechanisms that dynamically enable interaction across temporal resolutions (Hu et al., 2025a; Wang et al., 2024; 2025).

These components form the building blocks of a general architecture for multi-scale forecasting.

Definition 4.1 (Multi-scale forecasting model). *A multi-scale forecasting model for time series is a function that maps an input sequence \mathbf{x} to a predicted output sequence $\hat{\mathbf{y}} \in \mathbb{R}^T$:*

$$\hat{\mathbf{y}} = \left(\bigcirc_{i=1}^K (p_{\theta_i} \circ g_i) \right) (\mathbf{x}|s), \quad (3)$$

where the operator \bigcirc represents a composition of K stacked blocks. The i -th block consists of a multi-scale generator g_i and a multi-scale processor p with parameters θ_i .

Theorem 4.2. *For each of the following multi-scale model architectures: (Liu et al., 2022b;a; Wu et al., 2023; Wang et al., 2023; Cai et al., 2024; Wang et al., 2024; 2025; Hu et al., 2025a; Murad et al., 2025; Naghashi et al., 2025), there exists a set of multi-scale generators $\{g_i\}_{i=1}^K$ and multi-scale processors $\{p_{\theta_i}\}_{i=1}^K$ such that the model is an instance of Definition 4.1.*

Definition 4.1 provides a unifying perspective that encompasses a wide range of existing methods in time series analysis. The depth K specifies the number of generation-processing blocks, which determines the complexity of interactions across different scale parameters. A comprehensive categorization of existing models and the proof of Theorem 4.2 are provided in Appendix A.

4.2 LIMITATIONS OF EXISTING MULTI-SCALE METHODS

Based on our unifying framework of multi-scale methods, we first present two notable observations of these models fundamentally characterized by the design of their generators and processors.

Observation 1. Static multi-scale generator Existing models select a finite set of scale parameters $s = \{s_1, \dots, s_M\}$ in advance, e.g., window length, stride, and dilation, and consider them as fixed hyperparameters. A dyadic progression $s_i = 2^{i-1}$ is the most common choice due to its simplicity and computational efficiency (Wang et al., 2023; 2024; Hu et al., 2025a). Under fixed s , the multi-scale generator applies an identical set of scale operators uniformly over the entire time series.

Observation 2. Staged multi-scale processor Multi-scale processors commonly adopt a two-stage approach. First, temporal dependencies are processed within each scaled representation in isolation. Second, a separate cross-scale fusion mechanism integrates the outputs to capture interactions across different levels of scaling. The fusion process follows a fixed traversal order, such as coarse-to-fine, fine-to-coarse, or stacked multi-pass schedules (Zhong et al., 2024; Wang et al., 2025).

From these observations, we identify two key limitations.

Limitation 1. Rigid scaling strategy Fixed sets of scales cannot adjust to each target dataset. Since time series exhibit diverse periodicities and trends, a pre-defined s is rarely optimal. Some scaling operator families reinforce this limitation: moving averages with fixed windows fail to capture long-range dependencies, while fixed-rate subsampling imposes grid patterns that overlook informative resolutions. Such designs limit the extensibility of multi-scale modeling.

Limitation 2. Lack of flexibility Staged multi-scale processors constrain both flexibility and efficiency. In the first stage, temporal dependencies and scale information are modeled in isolation, preventing intertwined relationships. In the second stage, cross-scale fusion depends on the completion of all per-scale computations, creating sequential bottlenecks. Explicitly modeling cross-scale dependencies further enforces sequential processing, hindering efficiency and parallelization.

5 SIGMA: SIMPLIFIED MULTI-SCALE MODELING

We introduce SIGMA, a simple yet principled framework for multi-scale time series modeling. Motivated by the limitations of existing approaches, SIGMA ensures **learnable multi-scale** generation and flexible processing for any given time series; it instantiates the multi-scale generator with the learnable discrete Gaussian (LDG) kernel grounded in scale-space theory, and adopts a lightweight MLP as the multi-scale processor. Detailed proofs for all theorems are provided in Appendix A.

270 5.1 SIMPLIFIED MULTI-SCALE GENERATOR

271 We first extend scaling operator families to have multiple, continuous scale parameters $\mathbf{s} \in \mathbb{R}_+^M$ for
272 each transformation f , unlike the single, discrete parameter in Definition 3.1.

273 **Definition 5.1** (Extended scaling operator family). *A set of transformations $\mathcal{F} = \{f(\mathbf{x}|\mathbf{s}) \mid \mathbf{s} \in \mathbb{R}_+^M\}$ with $f(\cdot|\mathbf{s}) : \mathcal{X} \rightarrow \mathcal{X}_{\mathbf{s}} \subset \mathbb{R}^{L_{\mathbf{s}}}$ is an extended scaling operator family if it satisfies*

- 274 • *Consistency: For $s \in \mathbb{Z}_+$, $\{f(\mathbf{x}|s\mathbf{1})\}$ is a scaling operator family.*
- 275 • *Differentiability: For any $\mathbf{x} \in \mathcal{X}$, $f(\mathbf{x}|\mathbf{s})$ is continuously differentiable with respect to \mathbf{s} .*

276 The output dimension $L_{\mathbf{s}}$ depends on the set of scale parameters $\mathbf{s} = (s_1, \dots, s_M)$.

277 This generalized definition presents two key advantages. First, scale parameters \mathbf{s} can be optimized
278 directly from data via gradient-based methods, since $f(\mathbf{x}|\mathbf{s})$ is differentiable over \mathbb{R}_+^M . This allows
279 the model to [learn scale parameters from data](#) rather than relying on fixed, manual choices. Second,
280 a single extended operator f can express a rich range of scales simultaneously, removing the need
281 for stacked or hierarchical modules with discrete, pre-specified scale parameters. Together, these
282 properties directly address the [rigid scaling strategy](#) identified in Section 4.2.

283 Grounded in these principles, a variety of constructions could realize an extended scaling operator
284 family. In this work, we propose the *learnable discrete Gaussian* (LDG) kernel k , defined as

$$285 k(\mathbf{x}|\mathbf{s}) = \mathbf{K}(\mathbf{s})\mathbf{x}, \quad [\mathbf{K}(\mathbf{s})]_{i,j} = e^{-s_i} I_{|i-j|}(s_i), \quad (4)$$

286 where $[\cdot]_{i,j}$ denotes the (i, j) -th element of a matrix, $\mathbf{K}(\mathbf{s}) \in \mathbb{R}^{L \times L}$ with L being the length of the
287 sequence \mathbf{x} , and $I_n(\cdot)$ is the modified Bessel function of the first kind of integer order n . We model
288 the scale parameters \mathbf{s} to be learnable. As a result, k applies a dynamic scaling to each element of
289 \mathbf{x} , allowing a rich transformation without manually deciding various scale parameters.

290 Figure 3 illustrates how the LDG kernel actually works; it
291 yields smooth, scale-controlled transformations at each time
292 step with learnable scale parameters. It also provides a sym-
293 metric and effectively unbounded receptive field for capturing
294 long-range dependencies. Theorem 5.2 formally estab-
295 lishes that this kernel constitutes a valid extended scaling op-
296 erator family, confirming its suitability as a principled gen-
297 erator. Moreover, this Gaussian instantiation is the uniquely
298 determined operator by the discrete scale-space axioms as
299 shown in Theorem 5.3. This result aligns our method with scale-space theory, a principled frame-
300 work for multi-resolution signal analysis originating in computer vision (Witkin, 1987; Lindeberg,
301 2013).

302 **Theorem 5.2.** *The LDG kernel family, $\{k(\cdot|\mathbf{s}) \mid \mathbf{s} \in \mathbb{R}_+^M\}$, is an extended scaling operator family.*

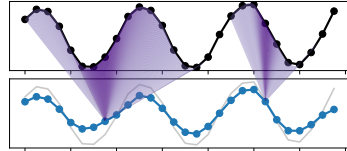
303 **Theorem 5.3.** *The LDG kernel family is the unique extended scaling operator family of symmetric
304 kernels that satisfies the discrete scale-space axioms (Lindeberg, 2002).*

305 5.2 SIMPLIFIED MULTI-SCALE PROCESSOR

306 To complement our principled generator, we propose a *simplified multi-scale processor* that over-
307 comes the inflexibility of existing methods. In previous work, heuristic scaling produced disjoint
308 representations, forcing complex processors to reconstruct cross-scale relationships. In contrast, our
309 LDG kernel yields all scales from a single, continuous, theoretically grounded operator, preserving
310 smooth inter-scale structure. As a result, instead of modeling temporal dependencies and cross-scale
311 interactions in separate stages, we are able to process both simultaneously in a single step.

312 We implement this processor with a lightweight, two-layer multi-layer perceptron (MLP) operating
313 on the concatenated multi-scale representations, followed by a linear projection to the forecasting
314 horizon. By a single forward pass, it jointly refines temporal features and fuses information across
315 scales, enabling efficient and flexible learning of nonlinear interactions without pre-wired pathways.
316 This one-shot composition directly addresses the lack of flexibility identified in Section 4.2.

317 Given an input sequence $\mathbf{x} \in \mathbb{R}^L$, we first produce d -dimensional embeddings $\mathbf{X} = \text{Embed}(\mathbf{x}) \in \mathbb{R}^{L \times d}$
318 as done in (Wu et al., 2021; Wang et al., 2024; 2025). Based on these embeddings, SIGMA,
319



320 Figure 3: The LDG kernel performs smoothing with [position-wise scales](#).

Table 2: Long-term forecasting results across eight datasets with horizons $T \in \{96, 192, 336, 720\}$ and input length fixed at 96. SIGMA achieves the smallest forecasting errors in 55 out of 80 evaluation settings and the second-best in 19 cases. These results confirm the effectiveness of principled multi-scale modeling for diverse long-horizon forecasting tasks.

Method	SIGMA (Ours)		AMD (2025a)		MultiPatch. (2025)		WPMixer (2025)		TimeMixer (2024)		MSGNet (2024)		MICN (2023)		TimesNet (2023)		Pyra. (2022b)		
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
Weather	96	0.160	0.204	0.182	0.227	0.172	0.211	0.164	0.210	0.166	0.214	0.161	0.209	0.192	0.250	0.172	0.221	0.195	0.281
	192	0.209	0.248	0.231	0.266	0.218	0.254	0.209	0.250	0.208	0.251	0.217	0.257	0.233	0.289	0.225	0.265	0.245	0.322
	336	0.270	0.293	0.283	0.302	0.275	0.296	0.264	0.290	0.265	0.293	0.280	0.303	0.283	0.332	0.289	0.309	0.307	0.365
	720	0.348	0.345	0.357	0.350	0.355	0.348	0.344	0.343	0.345	0.345	0.373	0.362	0.354	0.388	0.361	0.355	0.394	0.420
	Avg	0.247	0.273	0.263	0.286	0.255	0.278	0.245	0.273	0.246	0.276	0.258	0.283	0.266	0.315	0.262	0.288	0.285	0.347
Electricity	96	0.146	0.241	0.187	0.269	0.173	0.259	0.167	0.259	0.157	0.248	0.169	0.281	0.171	0.284	0.165	0.268	0.283	0.377
	192	0.163	0.257	0.191	0.274	0.181	0.267	0.179	0.268	0.169	0.260	0.189	0.298	0.178	0.290	0.182	0.284	0.296	0.391
	336	0.179	0.273	0.206	0.290	0.199	0.285	0.197	0.288	0.187	0.277	0.201	0.310	0.189	0.301	0.198	0.299	0.306	0.401
	720	0.213	0.304	0.248	0.323	0.239	0.318	0.232	0.315	0.227	0.312	0.238	0.340	0.208	0.318	0.231	0.325	0.305	0.393
	Avg	0.175	0.269	0.208	0.289	0.198	0.282	0.194	0.282	0.185	0.274	0.199	0.307	0.187	0.298	0.194	0.294	0.298	0.390
Traffic	96	0.431	0.288	0.544	0.345	0.471	0.318	0.528	0.347	0.477	0.309	0.599	0.353	0.516	0.309	0.590	0.318	0.678	0.384
	192	0.444	0.296	0.527	0.335	0.480	0.319	0.511	0.337	0.488	0.312	0.634	0.372	0.535	0.317	0.614	0.327	0.672	0.377
	336	0.461	0.303	0.538	0.339	0.499	0.329	0.519	0.337	0.506	0.319	0.663	0.391	0.548	0.322	0.640	0.342	0.681	0.381
	720	0.494	0.320	0.573	0.358	0.537	0.351	0.548	0.350	0.535	0.332	0.721	0.420	0.574	0.332	0.662	0.350	0.709	0.395
	Avg	0.458	0.302	0.546	0.344	0.497	0.329	0.527	0.343	0.501	0.318	0.654	0.384	0.543	0.320	0.627	0.334	0.685	0.384
Exchange	96	0.084	0.204	0.083	0.201	0.089	0.208	0.086	0.202	0.090	0.210	0.104	0.230	0.093	0.226	0.112	0.242	0.630	0.645
	192	0.174	0.297	0.175	0.297	0.187	0.308	0.176	0.296	0.185	0.305	0.200	0.322	0.184	0.331	0.214	0.334	0.935	0.782
	336	0.322	0.411	0.326	0.412	0.357	0.436	0.343	0.421	0.351	0.428	0.394	0.460	0.322	0.443	0.379	0.452	1.204	0.873
	720	0.833	0.687	0.847	0.693	0.913	0.724	0.900	0.712	0.911	0.713	1.027	0.772	0.780	0.694	0.961	0.746	1.956	1.117
	Avg	0.353	0.400	0.358	0.401	0.387	0.419	0.376	0.408	0.384	0.414	0.431	0.446	0.345	0.424	0.416	0.443	1.181	0.854
ETTh1	96	0.379	0.393	0.385	0.396	0.377	0.397	0.382	0.404	0.381	0.398	0.398	0.418	0.425	0.435	0.419	0.432	0.701	0.630
	192	0.430	0.425	0.437	0.425	0.427	0.428	0.438	0.427	0.441	0.434	0.444	0.445	0.505	0.484	0.474	0.464	0.850	0.713
	336	0.481	0.446	0.480	0.445	0.469	0.449	0.499	0.464	0.475	0.449	0.484	0.471	0.606	0.556	0.499	0.475	1.060	0.777
	720	0.480	0.468	0.485	0.469	0.499	0.485	0.487	0.471	0.522	0.494	0.509	0.498	0.752	0.648	0.529	0.500	1.905	0.803
	Avg	0.443	0.433	0.447	0.434	0.443	0.440	0.451	0.442	0.455	0.444	0.459	0.458	0.572	0.531	0.480	0.468	0.879	0.731
ETTh2	96	0.289	0.339	0.291	0.340	0.293	0.347	0.291	0.342	0.296	0.348	0.327	0.369	0.358	0.405	0.327	0.369	1.439	0.922
	192	0.369	0.394	0.373	0.391	0.372	0.396	0.376	0.397	0.375	0.395	0.408	0.417	0.497	0.484	0.404	0.412	5.640	1.894
	336	0.415	0.427	0.416	0.427	0.420	0.431	0.435	0.439	0.430	0.438	0.429	0.439	0.618	0.552	0.459	0.456	4.800	1.844
	720	0.431	0.446	0.424	0.441	0.431	0.450	0.459	0.462	0.456	0.460	0.446	0.459	0.856	0.666	0.451	0.459	4.466	1.824
	Avg	0.376	0.402	0.376	0.400	0.379	0.406	0.390	0.410	0.389	0.410	0.402	0.421	0.582	0.527	0.410	0.424	4.086	1.621
ETTm1	96	0.323	0.359	0.330	0.365	0.319	0.358	0.322	0.358	0.322	0.359	0.328	0.370	0.322	0.373	0.334	0.374	0.592	0.514
	192	0.360	0.381	0.372	0.384	0.363	0.385	0.361	0.382	0.364	0.384	0.371	0.395	0.361	0.402	0.404	0.409	0.645	0.568
	336	0.392	0.404	0.406	0.405	0.398	0.410	0.389	0.402	0.397	0.407	0.410	0.419	0.409	0.437	0.418	0.421	0.776	0.643
	720	0.455	0.442	0.471	0.440	0.460	0.448	0.469	0.447	0.456	0.444	0.494	0.465	0.506	0.498	0.484	0.455	0.936	0.730
	Avg	0.383	0.397	0.395	0.399	0.385	0.400	0.385	0.397	0.385	0.398	0.401	0.412	0.399	0.427	0.410	0.415	0.737	0.614
ETTm2	96	0.174	0.257	0.183	0.267	0.177	0.259	0.175	0.257	0.176	0.258	0.179	0.263	0.186	0.283	0.187	0.266	0.387	0.464
	192	0.239	0.299	0.246	0.306	0.243	0.304	0.240	0.299	0.241	0.302	0.250	0.308	0.278	0.352	0.257	0.309	0.676	0.623
	336	0.296	0.337	0.305	0.342	0.305	0.346	0.304	0.342	0.304	0.345	0.311	0.345	0.405	0.437	0.322	0.349	1.196	0.836
	720	0.394	0.394	0.404	0.397	0.407	0.404	0.398	0.397	0.405	0.402	0.416	0.406	0.546	0.515	0.427	0.409	3.588	1.460
	Avg	0.276	0.322	0.285	0.328	0.283	0.328	0.279	0.324	0.281	0.327	0.289	0.330	0.354	0.397	0.298	0.333	1.462	0.846

our multi-scale forecasting model, is succinctly represented by the following equation:

$$\hat{y} = W_1(\text{MLP}(H) + H)W_2, \quad H = K(s)X \parallel (I - K(s))X \in \mathbb{R}^{2L \times d}, \quad (5)$$

where $s \in \mathbb{R}_+^L$ is the set of learnable scale parameters, and $W_1 \in \mathbb{R}^{T \times 2L}$ and $W_2 \in \mathbb{R}^{d \times 1}$ denote the projection heads. These scale parameters are optimized as dataset-level parameters during training and kept fixed at inference. We decompose the input into a smoothed component $K(s)X$ and a residual component $(I - K(s))X$ analogous to the classical trend-seasonal decomposition (Cleveland et al., 1990). To further stabilize optimization and preserve scale-specific information, we incorporate a skip connection that adds each scale’s input back to its transformed output (He et al., 2016).

We adopt channel independence, processing each variable separately, which makes the architecture naturally applicable to multivariate settings (Zeng et al., 2023). We further apply reversible instance normalization to each time series variable to remove distribution shifts (Kim et al., 2022).

6 EXPERIMENTS

We conduct a comprehensive empirical evaluation of SIGMA on standard long-term and short-term forecasting datasets against eight state-of-the-art baseline models. We also provide deeper analyses to investigate its efficiency, robustness, and underlying design principles.

Table 3: Short-term forecasting results in the M4 benchmark dataset with various temporal granularities. SIGMA achieves the best performance in 11 of the 15 cases. This highlights its effectiveness in capturing multi-scale patterns across a diverse range of time series types.

	Method	SIGMA (Ours)	AMD (2025a)	MultiPatch. (2025)	WPMixer (2025)	TimeMixer (2024)	MSGNet (2024)	MICN (2023)	TimesNet (2023)	Pyra. (2022b)
Yearly	SMAPE	<u>13.314</u>	13.447	13.296	13.632	13.326	13.354	14.580	13.482	14.987
	MASE	2.989	3.022	3.009	3.075	<u>3.002</u>	2.989	3.382	3.056	3.361
	OWA	0.783	0.792	0.785	0.804	0.785	<u>0.784</u>	0.871	0.797	0.881
Quarterly	SMAPE	10.060	10.259	10.166	10.299	10.281	10.446	11.389	<u>10.116</u>	11.706
	MASE	1.177	1.211	<u>1.178</u>	1.217	1.206	1.248	1.380	1.186	1.397
	OWA	0.886	0.907	<u>0.892</u>	0.911	0.907	0.929	1.020	<u>0.892</u>	1.041
Monthly	SMAPE	12.750	12.898	12.810	12.945	12.984	12.970	13.797	<u>12.775</u>	14.444
	MASE	0.936	0.952	<u>0.942</u>	0.959	0.964	0.976	1.077	0.945	1.142
	OWA	0.882	0.895	<u>0.887</u>	0.900	0.903	0.908	0.985	<u>0.887</u>	1.038
Others	SMAPE	4.867	<u>4.822</u>	4.849	4.925	4.739	5.521	6.123	4.978	6.115
	MASE	3.316	<u>3.245</u>	3.271	3.259	3.234	3.829	4.196	3.265	4.156
	OWA	1.037	<u>1.021</u>	1.028	1.028	1.014	1.174	1.300	1.048	1.282
Weighted Average	SMAPE	11.840	11.987	<u>11.889</u>	12.067	12.002	12.080	13.015	11.910	13.495
	MASE	1.585	1.605	<u>1.591</u>	1.623	1.604	1.647	1.836	1.604	1.864
	OWA	0.868	0.881	<u>0.872</u>	0.887	0.882	0.898	0.983	0.876	1.015

Datasets We evaluate SIGMA on both long-term and short-term forecasting benchmarks, including (long-term) Weather, ETT (ETTh1, ETTh2, ETTm1, ETTm2), Electricity, Exchange-Rate, Traffic, and (short-term) M4. We follow the standard protocol and partition all datasets into training, validation, and test sets in chronological order, using a 6:2:2 ratio for the ETT datasets and a 7:1:2 ratio for the others (Zhou et al., 2021). Consistent with previous work, the observation window is fixed to $L = 96$, and forecasting horizons are set to $T \in \{96, 192, 336, 720\}$ for long-term forecasting tasks (Wu et al., 2023). The M4 benchmark consists of 100,000 univariate time series collected at multiple temporal frequencies ranging from yearly to hourly. We follow the official benchmark protocol, where the prediction lengths are fixed according to the frequency of each series (Makridakis et al., 2018). Comprehensive details for each dataset are provided in Appendix B.

Baselines We compare SIGMA against eight state-of-the-art baselines that adopt multi-scale modeling strategies: AMD (Hu et al., 2025a), MultiPatchFormer (Naghashi et al., 2025), WPMixer (Murad et al., 2025), TimeMixer (Wang et al., 2024), MSGNet (Cai et al., 2024), MICN (Wang et al., 2023), TimesNet (Wu et al., 2023) and Pyraformer (Liu et al., 2022b).

Experimental settings For each benchmark, we adopt evaluation metrics used in previous work. For long-term forecasting, we use mean squared error (MSE) and mean absolute error (MAE) (Wu et al., 2023). For the short-term benchmark, we report symmetric mean absolute percentage error (SMAPE), mean absolute scaled error (MASE), and the overall weighted average (OWA) (Oreshkin et al., 2020). We apply the same batch size, number of training epochs, and early-stopping strategy across all methods to maintain a unified setting (Wu et al., 2023). Each experiment is repeated three times, and the averaged results are reported. Additional details are provided in Appendix B.

6.1 COMPARISON WITH BASELINE MODELS

Long-term forecasting Table 2 shows that SIGMA achieves the best performance in 55 out of 80 evaluation settings for long-term forecasting across all datasets and horizons. The improvements are particularly seen in the high-dimensional Electricity and Traffic datasets. SIGMA achieves 5.4% and 7.8% reduction in MSE on average relative to the second-best model, respectively. It highlights the importance of [learning position-wise scale parameters from data to capture](#) fine-grained temporal dynamics through principled multi-scale modeling when applied to real-world multivariate time series. In comparison to the competitors that employ complex multi-scale architectures, the performance improvement of SIGMA with a simple, parameter-efficient architecture further verifies the significance of our design.

Short-term forecasting Table 3 shows that SIGMA also establishes clear superiority over the competitors on the M4 benchmark for short-term forecasting, achieving the best results in 11 out of 15 cases. It achieves the best average results, particularly in the Quarterly and Monthly datasets, which contain the largest number of series. On the other hand, SIGMA exhibits limited performance on the “Others” category, which lacks sufficient data as it accounts for less than 5% of the benchmark. This

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

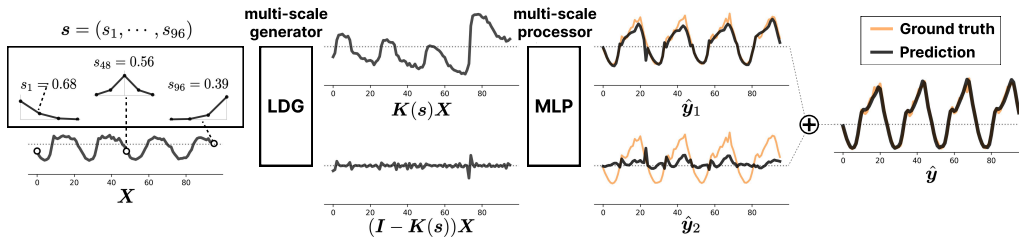


Figure 5: Case study on Traffic for predict-96 setting. The multi-scale generator employs the LDG kernel to extract trend-seasonal representations of the time series, while the multi-scale processor uses an MLP to produce complementary coarse- and fine-grained components. By integrating these complementary signals, SIGMA achieves more accurate and effective predictions.

suggests that learning position-wise scale parameters is most effective when trained on sufficiently large datasets, while less-structured series favor stronger inductive biases to extract meaningful patterns. Nevertheless, SIGMA outperforms all other baselines by a large margin overall.

6.2 DEEPER ANALYSIS ON SIGMA

Ablation study To validate the design principles of SIGMA, we conduct an ablation study on the ETTh1 dataset under the predict-720 setting, comparing SIGMA against five of its variants: ① integrates the trend-seasonal mixing multi-scale processor with TimeMixer with LDG kernel multi-scale generator; ② changes the LDG kernel to have a single scale parameter applied to all elements; ③ removes the scaling mechanism entirely and relies only on the raw input; ④ employs moving average, a non-learnable scale operator family; ⑤ uses unnormalized convolution, a learnable kernel that does not belong to the scale operator family.

As shown in Table 4, most modifications result in performance degradation. This indicates that the expressivity of the LDG kernel in ① and ②, the effect of scaling in ③, and the position-wise learnable scaling in ④ all play a critical role. While ① achieves performance comparable to SIGMA through a different multi-scale processor design, its multi-stage mixing architecture introduces additional computational overhead, whereas similar accuracy can be achieved with a simpler MLP-based processor. The significant performance drop in ⑤ demonstrates that invalid scaling operations which do not belong to the scale operator family can be sub-optimal, as they can transform the inherent characteristic of time series as well.

Efficiency analysis We assess efficiency in terms of per-iteration (i.e., batch) training time and memory footprint on the ETTh1 dataset under the predict-720 setting, with all methods executed on a single GPU using the same batch size. As shown in Figure 4, SIGMA attains the fastest training time and the most compact memory footprint while achieving the lowest MSE. Specifically, SIGMA reduces memory consumption by 3.8 times and improves training speed by 5.3 times compared to AMD, the previous state-of-the-art baseline.

This significant efficiency gap comes from fundamental architectural differences. AMD stacks three downsampling stages and applies per-variable multi-scale mixing in a sequential manner. While this approach may improve accuracy over other baselines, it imposes substantial memory and runtime overhead. In contrast, SIGMA employs a single-kernel formulation that preserves full parallelism in its representation updates. Our principled simplification of multi-scale modeling yields a more favorable speed-memory trade-off without compromising modeling capacity, leading to superior overall performance.

Table 4: Ablation study on ETTh1 with predict-720 setting.

Method	MSE	MAE
SIGMA	0.480	0.468
①	0.486	0.467
②	0.489	0.473
③	0.492	0.475
④	0.493	0.475
⑤	0.524	0.492

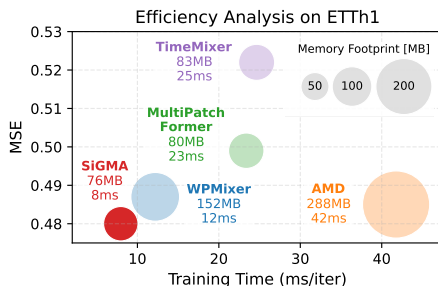


Figure 4: Efficiency analysis on ETTh1 with predict-720 setting. SIGMA achieves the best trade-off between accuracy and efficiency.

486 **Case study** We conduct a case study on the Traffic dataset under the predict-96 setting to illustrate
487 the behavior of SIGMA. As shown in Figure 5, SIGMA learns a scale parameter at each timestep. The
488 LDG kernel acts as a principled smoothing operator, allowing the model to selectively emphasize
489 broad trends or localized variations. The resulting representations are then processed by the MLP,
490 which captures richer temporal dependencies and adjusts the balance between coarse- and fine-
491 grained components. By integrating these components, SIGMA yields predictions that closely align
492 with the ground truth, validating the effectiveness of its multi-scale design.

493 494 7 CONCLUSION

495
496 In this work, we introduce a concept of *scaling operator families*, which allows us to unify existing
497 multi-scale models within a principled framework and decompose them into a pair of a multi-scale
498 generator and a multi-scale processor. Building on this foundation, we propose SIGMA, a simple
499 yet principled architecture that employs the learnable discrete Gaussian (LDG) kernel to enable
500 expressive scaling, where a lightweight MLP processor is then applied for efficient cross-scale fu-
501 sion. Extensive experiments on both long- and short-term forecasting benchmarks demonstrate that
502 SIGMA consistently outperforms state-of-the-art baselines, while substantially reducing memory
503 consumption and computational time. These findings underscore the effectiveness of SIGMA as a
504 reliable and generalizable solution for time-series forecasting. Future work can include extending
505 our framework to multivariate scaling and sample-specific [dynamic](#) modeling to further enhance its
506 expressiveness and applicability to cover a more diverse range of time series data.

507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

540 REPRODUCIBILITY STATEMENT

541

542 Code and scripts to reproduce all experiments are available at <https://anonymous.4open.science/r/SiGMA-ICLR2026> and also included in the supplementary material. We ensure
 543 reproducibility by using publicly available datasets with standard splits and protocols. All exper-
 544 iments are conducted in PyTorch under unified training and evaluation settings across baselines.
 545 We report results averaged over three runs, and provide ablation studies and efficiency analyses for
 546 verification.
 547

548

549 ETHICS STATEMENT

550

551 This work develops a general framework for time-series forecasting and does not involve sensitive
 552 personal data. The potential impacts include positive applications such as improving forecasting in
 553 energy, climate, and transportation. We encourage responsible use of our method in domains that
 554 align with ethical and societal values.
 555

556

557 REFERENCES

558

559 George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.

560

561 Wanlin Cai, Yuxuan Liang, Xianggen Liu, Jianshuai Feng, and Yuankai Wu. Msgnet: Learning
 562 multi-scale inter-series correlations for multivariate time series forecasting. In *Proceedings of the
 563 AAAI conference on artificial intelligence*, volume 38, pp. 11141–11149, 2024.

564 Cristian Challu, Kin G Olivares, Boris N Oreshkin, Federico Garza Ramirez, Max Mergenthaler
 565 Canseco, and Artur Dubrawski. Nhits: Neural hierarchical interpolation for time series forecast-
 566 ing. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 6989–6997,
 567 2023.

568

569 Ling Chen, Donghui Chen, Zongjiang Shang, Binqing Wu, Cen Zheng, Bo Wen, and Wei Zhang.
 570 Multi-scale adaptive graph neural network for multivariate time series forecasting. *IEEE Trans-
 571 actions on Knowledge and Data Engineering*, 35(10):10748–10761, 2023.

572 Robert B Cleveland, William S Cleveland, Jean E McRae, Irma Terpenning, et al. Stl: A seasonal-
 573 trend decomposition. *J. off. Stat*, 6(1):3–73, 1990.

574

575 Chirag Deb, Fan Zhang, Junjing Yang, Siew Eang Lee, and Kwok Wei Shah. A review on time series
 576 forecasting techniques for building energy consumption. *Renewable and Sustainable Energy
 577 Reviews*, 74:902–924, 2017.

578 Tripti Dimri, Shamshad Ahmad, and Mohammad Sharif. Time series analysis of climate variables
 579 using seasonal arima approach. *Journal of Earth System Science*, 129(1):149, 2020.

580

581 James Durbin and Siem Jan Koopman. *Time series analysis by state space methods*. Oxford univer-
 582 sity press, 2012.

583 Vijay Ekambaram, Arindam Jati, Nam Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam.
 584 Tsmixer: Lightweight mlp-mixer model for multivariate time series forecasting. In *Proceedings
 585 of the 29th ACM SIGKDD conference on knowledge discovery and data mining*, pp. 459–469,
 586 2023.

587

588 Georg Goerg. Forecastable component analysis. In *International conference on machine learning*,
 589 pp. 64–72. PMLR, 2013.

590 Edward James Hannan. *Multiple time series*. John Wiley & Sons, 2009.

591

592 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
 593 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
 770–778, 2016.

- 594 Yifan Hu, Peiyuan Liu, Peng Zhu, Dawei Cheng, and Tao Dai. Adaptive multi-scale decomposi-
595 tion framework for time series forecasting. In *Proceedings of the AAAI Conference on Artificial*
596 *Intelligence*, volume 39, pp. 17359–17367, 2025a.
- 597
- 598 Yifan Hu, Guibin Zhang, Peiyuan Liu, Disen Lan, Naiqi Li, Dawei Cheng, Tao Dai, Shu-Tao Xia,
599 and Shirui Pan. Timefilter: Patch-specific spatial-temporal graph filtration for time series fore-
600 casting. In *Forty-second International Conference on Machine Learning*, 2025b.
- 601
- 602 Eamonn Keogh, Selina Chu, David Hart, and Michael Pazzani. Segmenting time series: A survey
603 and novel approach. *Data mining in time series databases*, 57(2004):1–22, 2004.
- 604
- 605 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-
606 versible instance normalization for accurate time-series forecasting against distribution shift. In
607 *International Conference on Learning Representations*, 2022.
- 608
- 609 Diederik Kinga, Jimmy Ba Adam, et al. A method for stochastic optimization. In *International*
610 *conference on learning representations (ICLR)*, volume 5. California, 2015.
- 611
- 612 Milind Kolambe. Forecasting the future: A comprehensive review of time series prediction tech-
613 niques. *Journal of Electrical Systems*, 20:575–586, 04 2024. doi: 10.52783/jes.1478.
- 614
- 615 Tony Lindeberg. Scale-space for discrete signals. *IEEE transactions on pattern analysis and*
616 *machine intelligence*, 12(3):234–254, 2002.
- 617
- 618 Tony Lindeberg. *Scale-space theory in computer vision*, volume 256. Springer Science & Business
619 Media, 2013.
- 620
- 621 Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia LAI, Lingna Ma, and Qiang Xu. SCINet:
622 Time series modeling and forecasting with sample convolution and interaction. In Alice H. Oh,
623 Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information*
624 *Processing Systems*, 2022a.
- 625
- 626 Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar.
627 Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and fore-
628 casting. In *International Conference on Learning Representations*, 2022b.
- 629
- 630 Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.
631 itransformer: Inverted transformers are effective for time series forecasting. In *The Twelfth Inter-*
632 *national Conference on Learning Representations*, 2024.
- 633
- 634 Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. The m4 competition: Re-
635 sults, findings, conclusion and way forward. *International Journal of forecasting*, 34(4):802–808,
636 2018.
- 637
- 638 Md Mahmuddun Nabi Murad, Mehmet Aktukmak, and Yasin Yilmaz. Wpmixer: Efficient multi-
639 resolution mixing for long-term time series forecasting. In *Proceedings of the AAAI Conference*
640 *on Artificial Intelligence*, volume 39, pp. 19581–19588, 2025.
- 641
- 642
- 643
- 644
- 645
- 646
- 647
- Vahid Naghshahi, Mounir Boukadoum, and Abdoulaye Banire Diallo. A multiscale model for multi-
variate time series forecasting. *Scientific Reports*, 15(1):1565, 2025.
- 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.
- Boris N. Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural ba-
sis expansion analysis for interpretable time series forecasting. In *International Conference on*
Learning Representations, 2020.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-
performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

- 648 Donald B Percival and Andrew T Walden. *Wavelet methods for time series analysis*, volume 4.
649 Cambridge university press, 2000.
- 650
- 651 Xiaoming Shi, Shiyu Wang, Yuqi Nie, Dianqi Li, Zhou Ye, Qingsong Wen, and Ming Jin. Time-
652 mome: Billion-scale time series foundation models with mixture of experts. In *The Thirteenth*
653 *International Conference on Learning Representations*, 2025.
- 654 David Alexander Tedjopurnomo, Zhifeng Bao, Baihua Zheng, Farhana Murtaza Choudhury, and
655 Alex Kai Qin. A survey on modern deep neural network for traffic prediction: Trends, meth-
656 ods and challenges. *IEEE Transactions on Knowledge and Data Engineering*, 34(4):1544–1561,
657 2020.
- 658 Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. MICN: Multi-
659 scale local and global context modeling for long-term series forecasting. In *The Eleventh Inter-*
660 *national Conference on Learning Representations*, 2023.
- 661
- 662 Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y. Zhang, and
663 JUN ZHOU. Timemixer: Decomposable multiscale mixing for time series forecasting. In *The*
664 *Twelfth International Conference on Learning Representations*, 2024.
- 665
- 666 Shiyu Wang, Jiawei LI, Xiaoming Shi, Zhou Ye, Baichuan Mo, Wenze Lin, Ju Shengtong, Zhixuan
667 Chu, and Ming Jin. Timemixer++: A general time series pattern machine for universal predictive
668 analysis. In *The Thirteenth International Conference on Learning Representations*, 2025.
- 669
- 670 Andrew P Witkin. Scale-space filtering. In *Readings in computer vision*, pp. 329–332. Elsevier,
1987.
- 671
- 672 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans-
673 formers with auto-correlation for long-term series forecasting. *Advances in neural information*
674 *processing systems*, 34:22419–22430, 2021.
- 675
- 676 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet:
677 Temporal 2d-variation modeling for general time series analysis. In *The Eleventh International*
Conference on Learning Representations, 2023.
- 678
- 679 Zhifei Yang, Jia Zhang, and Zeyang Li. Multi-scale time series prediction model based on deep
680 learning and its application. *PLoS One*, 20(7):e0325474, 2025.
- 681
- 682 Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian, Long-
683 bing Cao, and Zhendong Niu. Frequency-domain mlps are more effective learners in time series
forecasting. *Advances in Neural Information Processing Systems*, 36:76656–76679, 2023.
- 684
- 685 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series
686 forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp.
11121–11128, 2023.
- 687
- 688 Shubao Zhao, Ming Jin, Zhaoxiang Hou, Chengyi Yang, Zengxiang Li, Qingsong Wen, and
689 Yi Wang. Himtm: Hierarchical multi-scale masked time series modeling with self-distillation
690 for long-term forecasting. In *Proceedings of the 33rd ACM International Conference on Informa-*
691 *tion and Knowledge Management*, pp. 3352–3362, 2024.
- 692
- 693 Yi Zheng, Qi Liu, Enhong Chen, Yong Ge, and J Leon Zhao. Time series classification using multi-
694 channels deep convolutional neural networks. In *International conference on web-age information*
management, pp. 298–310. Springer, 2014.
- 695
- 696 Shuhan Zhong, Sizhe Song, Weipeng Zhuo, Guanyao Li, Yang Liu, and S.-H. Gary Chan. A multi-
697 scale decomposition mlp-mixer for time series analysis. *Proc. VLDB Endow.*, 17(7):1723–1736,
698 March 2024.
- 699
- 700 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
701 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings*
of the AAAI conference on artificial intelligence, volume 35, pp. 11106–11115, 2021.

A PROOFS

A.1 PROOF OF THEOREM 3.2

Lemma A.1. *If $T : \mathbb{R}^m \rightarrow \mathbb{R}^n$ satisfies $\|T(\mathbf{x}) - T(\mathbf{y})\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2$ for all \mathbf{x}, \mathbf{y} , then $\|T(\mathbf{x})\|_2 \leq \|\mathbf{x}\|_2$ for all \mathbf{x} .*

Proof. $\|T(\mathbf{x})\|_2 = \|T(\mathbf{x}) - T(\mathbf{0})\|_2 \leq \|\mathbf{x} - \mathbf{0}\|_2 = \|\mathbf{x}\|_2$. \square

Lemma A.2. *Suppose for scales s_i, s_j we have $f(\cdot|s_i) = G_{i,j} \circ f(\cdot|s_j)$ for some non-expansive $G_{i,j}$. Then, $\|f(\mathbf{x}|s_i)\|_2 \leq \|f(\mathbf{x}|s_j)\|_2$ for all \mathbf{x} . Moreover, if there exists some \mathbf{y} in the range of $f(\cdot|s_j)$ with $\|G_{i,j}(\mathbf{y})\|_2 < \|\mathbf{y}\|_2$, then there exists $\mathbf{x}' \in \mathcal{X}$ such that $\|f(\mathbf{x}'|s_i)\|_2 < \|f(\mathbf{x}'|s_j)\|_2$.*

Proof. By Lemma A.1, we have $\|G_{i,j}(\mathbf{u})\|_2 \leq \|\mathbf{u}\|_2$ for all \mathbf{u} . Thus for any \mathbf{x} ,

$$\|f(\mathbf{x}|s_i)\|_2 = \|G_{i,j}(f(\mathbf{x}|s_j))\|_2 \leq \|f(\mathbf{x}|s_j)\|_2.$$

If some \mathbf{y} in the range of $f(\cdot|s_j)$ satisfies $\|G_{i,j}(\mathbf{y})\|_2 < \|\mathbf{y}\|_2$, pick \mathbf{x}' with $f(\mathbf{x}'|s_j) = \mathbf{y}$. Then

$$\|f(\mathbf{x}'|s_i)\|_2 = \|G_{i,j}(\mathbf{y})\|_2 < \|\mathbf{y}\|_2 = \|f(\mathbf{x}'|s_j)\|_2.$$

\square

Average pooling. Define

$$(f_{\text{avg}}(\mathbf{x}|s))_m := \frac{1}{s} \sum_{i \in B_m} x_i,$$

where $B_m = \{(m-1)s + 1, \dots, ms\}$ for $m = 1, \dots, L_s := P/s$.

Non-expansiveness. For any \mathbf{x}, \mathbf{y} and any block B_m ,

$$\begin{aligned} |f_{\text{avg}}(\mathbf{x}|s)_m - f_{\text{avg}}(\mathbf{y}|s)_m| &= \frac{1}{s} \left| \sum_{i \in B_m} (x_i - y_i) \right| \\ &\leq \frac{1}{s} \sqrt{s} \|\mathbf{x}|_{B_m} - \mathbf{y}|_{B_m}\|_2 \leq \|\mathbf{x}|_{B_m} - \mathbf{y}|_{B_m}\|_2. \end{aligned}$$

Summing squares over m yields

$$\|f_{\text{avg}}(\mathbf{x}|s) - f_{\text{avg}}(\mathbf{y}|s)\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2,$$

so non-expansiveness holds.

Energy reduction. Let $s_i = m \cdot s_j$ with $m > 1$. Each s_i -block is a union of m consecutive s_j -blocks.

Let $\mathbf{y} = f_{\text{avg}}(\mathbf{x}|s_j) \in \mathbb{R}^{L_{s_j}}$, and define

$$(G_{i,j}(\mathbf{y}))_k := \frac{1}{m} \sum_{r=0}^{m-1} y_{km+r},$$

i.e., average pooling on \mathbf{y} with block size m . Then one checks directly that

$$f_{\text{avg}}(\mathbf{x}|s_i) = G_{i,j}(f_{\text{avg}}(\mathbf{x}|s_j)),$$

and $G_{i,j}$ is the same blockwise averaging operator as above, hence non-expansive. By Lemma A.2, for all \mathbf{x} ,

$$\|f_{\text{avg}}(\mathbf{x}|s_i)\|_2 \leq \|f_{\text{avg}}(\mathbf{x}|s_j)\|_2.$$

To see strictness for some input, choose \mathbf{y} with two distinct values in each m -block, so that the block average has strictly smaller energy than \mathbf{y} . Since $f_{\text{avg}}(\cdot|s_j)$ is surjective onto $\mathbb{R}^{L_{s_j}}$, there exists \mathbf{x}' with $f_{\text{avg}}(\mathbf{x}'|s_j) = \mathbf{y}$, and Lemma A.2 yields

$$\|f_{\text{avg}}(\mathbf{x}'|s_i)\|_2 < \|f_{\text{avg}}(\mathbf{x}'|s_j)\|_2.$$

756 **Max- and min-pooling.** Define

$$757 (f_{\max}(\mathbf{x}|s))_m := \max_{i \in B_m} x_i, (f_{\min}(\mathbf{x}|s))_m := \min_{i \in B_m} x_i.$$

759 *Non-expansiveness.* For any $\mathbf{u}, \mathbf{v} \in \mathbb{R}^s$,

$$760 \left| \max_i u_i - \max_i v_i \right| \leq \max_i |u_i - v_i| \leq \|\mathbf{u} - \mathbf{v}\|_2,$$

761 and similarly for the minimum. Thus each block map is 1-Lipschitz; by summing over blocks
762 we get

$$763 \|f_{\max}(\mathbf{x}|s) - f_{\max}(\mathbf{y}|s)\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2$$

764 and likewise for f_{\min} .

765 *Energy reduction.* Let $s_i = m \cdot s_j$ with $m > 1$. Define $y = f_{\max}(\mathbf{x}|s_j)$ and

$$766 (G_{i,j}(y))_k := \max_{r \in \{0, \dots, m-1\}} y_{km+r},$$

767 i.e., max-pooling over m consecutive entries of y . Then

$$768 f_{\max}(\mathbf{x}|s_i) = G_{i,j}(f_{\max}(\mathbf{x}|s_j)),$$

769 and the same max-Lipschitz argument as above shows $G_{i,j}$ is non-expansive. Thus by Lemma A.2,
770 for all \mathbf{x} ,

$$771 \|f_{\max}(\mathbf{x}|s_i)\|_2 \leq \|f_{\max}(\mathbf{x}|s_j)\|_2.$$

772 To get strictness, choose \mathbf{y} such that in each group of m entries there is at least one strictly smaller
773 than the maximum; then $\|G_{i,j}(\mathbf{y})\|_2 < \|\mathbf{y}\|_2$. As before, we can realize such a \mathbf{y} as $f_{\max}(\mathbf{x}'|s_j)$
774 by setting one large value and smaller ones within each block. The same reasoning applies to min-
775 pooling by symmetry.

776 **Moving average.** Define

$$777 (f_{\text{ma}}(\mathbf{x}|s))_t := \frac{1}{s} \sum_{i=0}^{s-1} x_{t-i}, \quad t = 1, \dots, P.$$

778 This is convolution with $h_s = \frac{1}{s} \mathbf{1}_s$.

779 *Non-expansiveness.* By Young's inequality,

$$780 \|f_{\text{ma}}(\mathbf{x}|s) - f_{\text{ma}}(\mathbf{y}|s)\|_2 = \|h_s * (\mathbf{x} - \mathbf{y})\|_2 \leq \|h_s\|_1 \|\mathbf{x} - \mathbf{y}\|_2 = \|\mathbf{x} - \mathbf{y}\|_2.$$

781 Thus $f_{\text{ma}}(\cdot|s)$ is non-expansive.

782 *Energy reduction.* Let $X(\omega)$ and $Y_s(\omega)$ denote the Fourier transforms of \mathbf{x} and $f_{\text{ma}}(\mathbf{x}|s)$, respec-
783 tively, and $H_s(\omega)$ the frequency response of h_s . We have

$$784 Y_s(\omega) = H_s(\omega)X(\omega), H_s(\omega) = \frac{1}{s} \sum_{k=0}^{s-1} e^{-i\omega k} = \frac{1}{s} \frac{\sin(s\omega/2)}{\sin(\omega/2)} e^{-i\omega(s-1)/2}.$$

785 Thus

$$786 \|f_{\text{ma}}(\mathbf{x}|s)\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |H_s(\omega)|^2 |X(\omega)|^2 d\omega.$$

787 For $s_i = m \cdot s_j$, write $\alpha = s_j \omega / 2$. Then

$$788 \frac{|H_{s_i}(\omega)|}{|H_{s_j}(\omega)|} = \frac{1}{m} \frac{|\sin(m\alpha)|}{|\sin(\alpha)|} \leq 1,$$

789 since $|\sin(m\alpha)| \leq m |\sin(\alpha)|$ for all $m \in \mathbb{N}$ and all α . Hence $|H_{s_i}(\omega)|^2 \leq |H_{s_j}(\omega)|^2$ for all ω , and
790 therefore, for any \mathbf{x} ,

$$791 \|f_{\text{ma}}(\mathbf{x}|s_i)\|_2^2 \leq \|f_{\text{ma}}(\mathbf{x}|s_j)\|_2^2.$$

To see when the inequality is strict, define

$$D(\omega) := |H_{s_j}(\omega)|^2 - |H_{s_i}(\omega)|^2 \geq 0.$$

The above argument shows $D(\omega) \geq 0$ for all ω , and $H_{s_i} \not\equiv H_{s_j}$ implies that $D(\omega) > 0$ on a set of nonzero measure. For any \mathbf{x} ,

$$\|f_{\text{ma}}(\mathbf{x}|s_j)\|_2^2 - \|f_{\text{ma}}(\mathbf{x}|s_i)\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} D(\omega) |X(\omega)|^2 d\omega.$$

Therefore,

$$\|f_{\text{ma}}(\mathbf{x}|s_j)\|_2^2 > \|f_{\text{ma}}(\mathbf{x}|s_i)\|_2^2$$

for every input \mathbf{x} whose spectrum $|X(\omega)|^2$ places nonzero mass on the region where $D(\omega) > 0$. Equivalently, the inequality is strict for any non-degenerate \mathbf{x} whose energy is not concentrated entirely on the set where $|H_{s_i}(\omega)| = |H_{s_j}(\omega)|$.

Subsampling. Define

$$(f_{\text{sub}}(\mathbf{x}|s))_m := x_{(m-1)s+1}, \quad L_s = P/s.$$

Non-expansiveness. For each block B_m , define $g_s(\mathbf{u}) = u_1$. Then

$$|g_s(\mathbf{u}) - g_s(\mathbf{v})| = |u_1 - v_1| \leq \|\mathbf{u} - \mathbf{v}\|_2.$$

Applying the blockwise arguments shows that $f_{\text{sub}}(\cdot|s)$ is non-expansive.

Energy reduction. If $s_i = m \cdot s_j$, pure decimation satisfies

$$f_{\text{sub}}(\mathbf{x}|s_i) = f_{\text{sub}}(f_{\text{sub}}(\mathbf{x}|s_j)|m).$$

So $f_{\text{sub}}(\cdot|s_i) = G_{i,j} \circ f_{\text{sub}}(\cdot|s_j)$ with $G_{i,j}$ another subsampling operator, which is non-expansive. Hence Lemma A.2 yields

$$\|f_{\text{sub}}(\mathbf{x}|s_i)\|_2 \leq \|f_{\text{sub}}(\mathbf{x}|s_j)\|_2$$

for all \mathbf{x} . The strict inequality holds by setting \mathbf{x} as nonzero values concentrated on indices that are dropped at the coarser stride.

Segmentation. Define

$$f_{\text{seg}}(\mathbf{x}|s) := (x_1, \dots, x_{L_s}) \in \mathbb{R}^{L_s}, \quad L_s = P/s,$$

Non-expansiveness. This is an orthogonal projection onto the first L_s coordinates:

$$f_{\text{seg}}(\mathbf{x}|s) = R_s \mathbf{x}, \quad R_s R_s^\top = I_{L_s}.$$

Thus

$$\|f_{\text{seg}}(\mathbf{x}|s) - f_{\text{seg}}(\mathbf{y}|s)\|_2 = \|R_s(\mathbf{x} - \mathbf{y})\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2,$$

so f_{seg} is non-expansive and $\|f_{\text{seg}}(\mathbf{x}|s)\|_2 \leq \|\mathbf{x}\|_2$.

Energy reduction. If $s_i = m \cdot s_j$, then

$$f_{\text{seg}}(\mathbf{x}|s_i) = R_{s_i} \mathbf{x} = R_{s_i} R_{s_j}^\top f_{\text{seg}}(\mathbf{x}|s_j).$$

Here $G_{i,j} := R_{s_i} R_{s_j}^\top$ is an orthogonal projection from $\mathbb{R}^{L_{s_j}}$ to $\mathbb{R}^{L_{s_i}}$, hence non-expansive. Thus by Lemma A.2, for all \mathbf{x} ,

$$\|f_{\text{seg}}(\mathbf{x}|s_i)\|_2 \leq \|f_{\text{seg}}(\mathbf{x}|s_j)\|_2.$$

The inequality is strict exactly when at least one of the discarded coordinates is nonzero, i.e., when

$$x_k \neq 0 \quad \text{for some } k \in \{L_{s_i} + 1, \dots, L_{s_j}\}.$$

For such \mathbf{x} we have

$$\sum_{k=L_{s_i}+1}^{L_{s_j}} x_k^2 > 0 \Rightarrow \|f_{\text{seg}}(\mathbf{x}|s_i)\|_2 < \|f_{\text{seg}}(\mathbf{x}|s_j)\|_2.$$

Wavelet decomposition. Let $\{\phi_{s,k}\}_k$ be an orthonormal basis of V_s , and define the wavelet approximation coefficients by

$$f_{\text{wav}}(\mathbf{x}|s)_k := \langle \mathbf{x}, \phi_{s,k} \rangle.$$

For standard wavelets, the approximation spaces $\{V_s\}$ are typically constructed as nested subspaces indexed by dyadic scales $s = 2^j$. Here, we consider a more general multiplicative scale set $\mathcal{S} \subset \mathbb{Z}_+$ and assume that to each $s \in \mathcal{S}$ we associate an approximation subspace $V_s \subset \mathbb{R}^{L_s}$ with orthogonal projector $P_s : \mathbb{R}^L \rightarrow V_s$, such that

$$s_i = m s_j \text{ with } m > 1 \implies V_{s_i} \subset V_{s_j} \text{ and } P_{s_i} = P_{s_i} P_{s_j}.$$

Equivalently, if $C_s : V_s \rightarrow \mathbb{R}^{L_s}$ denotes the isometry mapping $P_s \mathbf{x}$ to its coordinates in the basis $\{\phi_{s,k}\}_k$, then $f_{\text{wav}}(\mathbf{x}|s) = C_s P_s \mathbf{x} \in \mathbb{R}^{L_s}$.

Non-expansiveness. For any $s \in \mathcal{S}$, P_s is an orthogonal projector, hence

$$\|P_s \mathbf{x} - P_s \mathbf{y}\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2.$$

Since C_s is an isometry on V_s ,

$$\|f_{\text{wav}}(\mathbf{x}|s) - f_{\text{wav}}(\mathbf{y}|s)\|_2 = \|C_s(P_s \mathbf{x} - P_s \mathbf{y})\|_2 = \|P_s \mathbf{x} - P_s \mathbf{y}\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2,$$

so $f_{\text{wav}}(\cdot|s)$ is non-expansive.

Energy reduction. Fix $s_i, s_j \in \mathcal{S}$ with $s_i = m s_j, m > 1$. For any \mathbf{x} , we can write

$$f_{\text{wav}}(\mathbf{x}|s_j) = C_{s_j} P_{s_j} \mathbf{x}, \quad f_{\text{wav}}(\mathbf{x}|s_i) = C_{s_i} P_{s_i} \mathbf{x} = C_{s_i} P_{s_i} P_{s_j} \mathbf{x}.$$

Define a linear map $G_{i,j} : \mathbb{R}^{L_{s_j}} \rightarrow \mathbb{R}^{L_{s_i}}$ by

$$G_{i,j} := C_{s_i} P_{s_i} P_{s_j}^\top C_{s_j}^\top, \quad \text{so that } f_{\text{wav}}(\mathbf{x}|s_i) = G_{i,j}(f_{\text{wav}}(\mathbf{x}|s_j)).$$

Here C_{s_j}, C_{s_i} are isometries and P_{s_i} is an orthogonal projector, so

$$\|G_{i,j}\|_2 \leq 1,$$

i.e., $G_{i,j}$ is non-expansive. By Lemma A.2, for all \mathbf{x} ,

$$\|f_{\text{wav}}(\mathbf{x}|s_i)\|_2 = \|G_{i,j}(f_{\text{wav}}(\mathbf{x}|s_j))\|_2 \leq \|f_{\text{wav}}(\mathbf{x}|s_j)\|_2.$$

For strictness, pick any $\mathbf{x} \in V_{s_j} \setminus V_{s_i}$. Then $P_{s_j} \mathbf{x} = \mathbf{x}$, while $P_{s_i} \mathbf{x}$ is the orthogonal projection of \mathbf{x} onto the strictly smaller subspace V_{s_i} , so by Pythagoras theorem, $\|P_{s_i} \mathbf{x}\|_2 < \|\mathbf{x}\|_2$. Then,

$$\|f_{\text{wav}}(\mathbf{x}|s_i)\|_2 = \|C_{s_i} P_{s_i} \mathbf{x}\|_2 < \|C_{s_j} \mathbf{x}\|_2 = \|f_{\text{wav}}(\mathbf{x}|s_j)\|_2.$$

A.2 PROOF OF THEOREM 3.3

Proof. We will show that each of operations is scale-degenerate in the sense that, for some pair $s_i > s_j$, we have

$$\|f(\mathbf{x}|s_i)\|_2 = \|f(\mathbf{x}|s_j)\|_2 \quad \text{for all } \mathbf{x} \in \mathcal{X},$$

so there exists no input \mathbf{x} that yields strict inequality. In other words, all such families are degenerate in the scale parameter and therefore violate the energy reduction condition in Definition 3.1.

Constant mappings. If $f(\mathbf{x}|s) = \mathbf{c}$ for some fixed $\mathbf{c} \in \mathbb{R}^{L_s}$, then

$$\|f(\mathbf{x}|s_i)\|_2 = \|f(\mathbf{x}|s_j)\|_2 = \|\mathbf{c}\|_2.$$

Permutations. If $f(\mathbf{x}|s) = \Pi_s \mathbf{x}$ for a permutation matrix Π_s , then

$$\|f(\mathbf{x}|s_i)\|_2 = \|f(\mathbf{x}|s_j)\|_2 = \|\mathbf{x}\|_2.$$

Additive shifts. If $f(\mathbf{x}|s) = \mathbf{x} + \mathbf{b}$, then

$$\|f(\mathbf{x}|s_i)\|_2 = \|f(\mathbf{x}|s_j)\|_2 = \|\mathbf{x} + \mathbf{b}\|_2.$$

Scalar multiplications. If $f(\mathbf{x}|s) = c\mathbf{x}$ for a constant $c \in \mathbb{R}$, then

$$\|f(\mathbf{x}|s_i)\|_2 = \|f(\mathbf{x}|s_j)\|_2 = |c|\|\mathbf{x}\|_2.$$

General linear maps. If $f(\mathbf{x}|s) = \mathbf{W}\mathbf{x}$, then

$$\|f(\mathbf{x}|s_i)\|_2 = \|f(\mathbf{x}|s_j)\|_2 = \|\mathbf{W}\mathbf{x}\|_2.$$

□

A.3 PROOF OF THEOREM 4.2

Proof. Each listed architecture admits a decomposition into (i) multi-scale generators g_i and (ii) multi-scale processors p_{θ_i} . We specify $(K, \{g_i\}, \{p_{\theta_i}\})$ for each.

Here is a cleaner, camera-ready rewrite that keeps all your formulas intact, standardizes notation to \mathbf{x} , and tightens phrasing/formatting.

Pyraformer (Liu et al., 2022b).

- **Choice of K .** $K = 1$.
- **Generator g_1 (CSCM: strided-conv subsampling and pyramid).** Let $\mathbf{x}^{(0)} = \mathbf{x} \in \mathbb{R}^{L \times d}$ and fix $C \geq 2$ as a stride and size of kernel. For $s = 1, \dots, M$ define

$$\mathbf{x}^{(s)} = \text{Conv1D}(\mathbf{x}^{(s-1)}; \text{kernel} = C, \text{stride} = C) \in \mathbb{R}^{L_s \times d_s}, \quad L_s = \left\lfloor \frac{L}{C^s} \right\rfloor.$$

- **Processor p_{θ_1} (PAM).** Index pyramid nodes by $u = (s, t)$ and let $\mathcal{N}(u)$ be the intra-, inter-scale neighborhood. With multi-head projections

$$Q_u = \mathbf{x}_u W_Q^{(s)}, \quad K_v = \mathbf{x}_v W_K^{(s(v))}, \quad V_v = \mathbf{x}_v W_V^{(s(v))},$$

define

$$\alpha_{uv} = \frac{\exp(\langle Q_u, K_v \rangle / \sqrt{d_k})}{\sum_{w \in \mathcal{N}(u)} \exp(\langle Q_u, K_w \rangle / \sqrt{d_k})}, \quad Z_u = \sum_{v \in \mathcal{N}(u)} \alpha_{uv} V_v,$$

and residual connection, feed-forward network update

$$Y_u = \text{LN}(\text{LN}(\mathbf{x}_u + Z_u) + \text{FFN}(\text{LN}(\mathbf{x}_u + Z_u))).$$

Upsample each $Y^{(s)}$ to length L via U_s and fuse:

$$H = \text{Fuse}(U_0 Y^{(0)}, \dots, U_M Y^{(M)}) \in \mathbb{R}^{L \times d}, \quad \hat{\mathbf{y}} = \text{TP}(H) \in \mathbb{R}^T.$$

SCINet (Liu et al., 2022a).

- **Choice of K .** $K = L$, which is the number of SCI levels.
- **Generator g_i (even/odd subsampling).** Let $\mathbf{x} \in \mathbb{R}^{L \times d}$ be the current 1D sequence at level i . Define

$$S_{\text{even}}(\mathbf{x}) = (\mathbf{x}_0, \mathbf{x}_2, \dots) \in \mathbb{R}^{\lceil L/2 \rceil \times d}, \quad S_{\text{odd}}(\mathbf{x}) = (\mathbf{x}_1, \mathbf{x}_3, \dots) \in \mathbb{R}^{\lfloor L/2 \rfloor \times d},$$

and set

$$g_i(\mathbf{x}) = (\mathbf{x}_{\text{even}}, \mathbf{x}_{\text{odd}}), \quad \mathbf{x}_{\text{even}} = S_{\text{even}}(\mathbf{x}), \quad \mathbf{x}_{\text{odd}} = S_{\text{odd}}(\mathbf{x}).$$

- **Processor p_{θ_i} (SCI-Block: cross-gating, interaction, re-interleaving).** Let $\phi, \psi, \rho, \eta : \mathbb{R}^{* \times d} \rightarrow \mathbb{R}^{* \times d}$ be 1D conv modules. Cross-gating:

$$\begin{aligned} \mathbf{x}'_{\text{odd}} &= \mathbf{x}_{\text{odd}} \odot \exp(\phi(\mathbf{x}_{\text{even}})), & \mathbf{x}'_{\text{even}} &= \mathbf{x}_{\text{even}} \odot \exp(\psi(\mathbf{x}_{\text{odd}})). \\ \tilde{\mathbf{x}}_{\text{odd}} &= \mathbf{x}'_{\text{odd}} \pm \rho(\mathbf{x}'_{\text{even}}), & \tilde{\mathbf{x}}_{\text{even}} &= \mathbf{x}'_{\text{even}} \pm \eta(\mathbf{x}'_{\text{odd}}). \end{aligned}$$

Re-interleave to a single stream:

$$p_{\theta_i}(\mathbf{x}_{\text{even}}, \mathbf{x}_{\text{odd}}) = \text{Interleave}(\tilde{\mathbf{x}}_{\text{even}}, \tilde{\mathbf{x}}_{\text{odd}}) \in \mathbb{R}^{L \times d}.$$

Apply L levels recursively:

$$\hat{\mathbf{y}} = (p_{\theta_L} \circ g_L \circ \dots \circ p_{\theta_1} \circ g_1)(\mathbf{x}).$$

TimesNet (Wu et al., 2023).

- **Choice of K .** $K = B$, the number of TimesBlocks.
- **Generator g_i (frequency-to-time imaging).** Given $\mathbf{x} \in \mathbb{R}^{L \times d}$, compute $\hat{\mathbf{x}} = \text{FFT}(\mathbf{x})$ and select k dominant periods $\{p_m\}_{m=1}^k$. For each p , set $s = \lfloor L/p \rfloor$, $f = p$, and fold to a time-period image

$$\tilde{\mathbf{x}}^{(p)} = \text{Reshape}(\mathbf{x}, s = \lfloor L/p \rfloor, f = p) \in \mathbb{R}^{d \times s \times f}, \quad g_i(\mathbf{x}) = \{\tilde{\mathbf{x}}^{(p_m)}\}_{m=1}^k.$$

- **Processor p_{θ_i} (2D conv mixing, adaptive aggregation).** Process each $\tilde{\mathbf{x}}^{(p_m)}$ with a 2D block to obtain $\hat{\mathbf{x}}_{2D}^{(p_m)}$, fold back to $\hat{\mathbf{x}}_{1D}^{(p_m)} \in \mathbb{R}^{L \times d}$, and aggregate

$$Z_i = \sum_{m=1}^k \text{Conv2D}_{\theta_i}(\tilde{\mathbf{x}}^{(p_m)}), \quad \hat{\mathbf{x}}_{1D}^{(p_m)} = \text{FoldBack}(\text{Proc}_{\theta_i}(\tilde{\mathbf{x}}^{(p_m)})).$$

Let $A_{\ell-1}^{f_1}, \dots, A_{\ell-1}^{f_k}$ be amplitudes from block $(\ell - 1)$ and define

$$(\tilde{A}_{\ell-1}^{f_1}, \dots, \tilde{A}_{\ell-1}^{f_k}) = \text{Softmax}(A_{\ell-1}^{f_1}, \dots, A_{\ell-1}^{f_k}),$$

then

$$p_{\theta_i}(\{\tilde{\mathbf{x}}^{(p_m)}\}_{m=1}^k) = \mathbf{x}_{\ell}^{1D} = \sum_{m=1}^k \tilde{A}_{\ell-1}^{f_m} \odot \hat{\mathbf{x}}_{1D}^{(p_m)} \in \mathbb{R}^{L \times d}.$$

MICN (Wang et al., 2023).

- **Choice of K .** $K = H$, the number of multi-branch inception stages.
- **Generator g_i (downsampling, isometric convolution).** Given branch scale K_i and current sequence $\mathbf{x}^{(i-1)} \in \mathbb{R}^{L_{i-1} \times d}$,

$$\mathbf{x}_{\text{gen}}^{(i)} = \text{Conv1D}(\text{AvgPool}_{K_i}(\mathbf{x}^{(i-1)}), \text{stride} = K_i) \in \mathbb{R}^{L_i \times d_i}, \quad L_i = \left\lfloor \frac{L_{i-1}}{K_i} \right\rfloor,$$

and set $g_i(\mathbf{x}^{(i-1)}) = \mathbf{x}_{\text{gen}}^{(i)}$.

- **Processor p_{θ_i} (inception mixer: local and global fusion).** With local kernel set \mathcal{R} and global branch G_i ,

$$\text{LocalMix}_i(x) = \sum_{r \in \mathcal{R}} U_{i,r} \text{Conv1D}_{k=r}(x),$$

$$\text{GlobalMix}_i(x) = G_i(\text{Conv1D}_{\text{global}}(x)),$$

and the stage output

$$p_{\theta_i}(\mathbf{x}_{\text{gen}}^{(i)}) = \text{Merge}(\text{LocalMix}_i(\mathbf{x}_{\text{gen}}^{(i)}), \text{GlobalMix}_i(\mathbf{x}_{\text{gen}}^{(i)})) \in \mathbb{R}^{L_i \times d_i}.$$

MSGNet (Cai et al., 2024).

- **Choice of K .** $K = 1$.
- **Generator g_1 (Time imaging).** Let $\mathbf{x} \in \mathbb{R}^{N \times L \times d}$ (nodes \times time \times channels). Choose $\{p_m\}_{m=1}^M$ and set $s_m = \lfloor L/p_m \rfloor$, $f_m = p_m$. Reshape per node:

$$\tilde{\mathbf{x}}^{(m)} = \text{Reshape}(\mathbf{x}, s_m, f_m) \in \mathbb{R}^{N \times d \times s_m \times f_m}, \quad g_1(\mathbf{x}) = \{\tilde{\mathbf{x}}^{(m)}\}_{m=1}^M.$$

- **Processor p_{θ_1} (multi-scale adaptive graph convolution and multi-head attention).** (i) *Graph filtering per scale.* Build a learnable adjacency

$$A^{(m)} = \text{Softmax}_{\text{row}}(\Phi_m(\text{Pool}_{t,f}(\tilde{\mathbf{x}}^{(m)}))) \in \mathbb{R}^{N \times N}, \quad \tilde{L}^{(m)} = \mathbf{I} - D^{(m)-\frac{1}{2}} A^{(m)} D^{(m)-\frac{1}{2}},$$

and apply a K_g -order Chebyshev filter

$$\mathcal{G}^{(m)} = \sum_{k=0}^{K_g} \theta_k^{(m)} T_k(\tilde{L}^{(m)}) \tilde{\mathbf{x}}^{(m)} \in \mathbb{R}^{N \times d \times s_m \times f_m}.$$

(ii) *Multi-head attention over time-period tokens per scale.* Flatten (s_m, f_m) to $S_m = s_m f_m$ tokens and let $Z^{(m)} \in \mathbb{R}^{(N \cdot d) \times S_m}$. For head h ,

$$Q_h^{(m)} = Z^{(m)} W_{Q,h}^{(m)}, \quad K_h^{(m)} = Z^{(m)} W_{K,h}^{(m)}, \quad V_h^{(m)} = Z^{(m)} W_{V,h}^{(m)},$$

$$\text{MHA}^{(m)}(Z^{(m)}) = [\text{Concat}_h \text{Softmax}(\frac{Q_h^{(m)} K_h^{(m)\top}}{\sqrt{d_k}}) V_h^{(m)}] W_O^{(m)}.$$

Unflatten and sum across scales:

$$Y^{(m)} = \text{Unflatten}(\text{MHA}^{(m)}(Z^{(m)}), s_m, f_m), \quad \hat{\mathbf{x}} = \sum_{m=1}^M \text{FoldBack}(Y^{(m)}) \in \mathbb{R}^{N \times L \times d},$$

and predict

$$p_{\theta_1}(\{\tilde{\mathbf{x}}^{(m)}\}_{m=1}^M) = \text{TP}(\hat{\mathbf{x}}) \in \mathbb{R}^{N \times T}.$$

TimeMixer (Wang et al., 2024).

- **Choice of K .** $K = 1$.
- **Generator g_1 (downsampling by average pooling).**

$$g_1(\mathbf{x}) = \{\mathbf{x}_m\}_{m=0}^M, \quad \mathbf{x}_0 = \mathbf{x}, \quad \mathbf{x}_m = \text{AvgPool}_{2^m}(\mathbf{x}) \in \mathbb{R}^{\lfloor L/2^m \rfloor \times d} \quad (m \geq 1).$$

- **Processor p_{θ_1} (PDM^L, FMM).** Let $X^0 = \{\mathbf{x}_m\}_{m=0}^M$. For $\ell = 1, \dots, L$ and each m ,

$$(s_m^\ell, t_m^\ell) = \text{Decompose}(X_m^{\ell-1}), \quad X_m^{\ell-1} = s_m^\ell + t_m^\ell.$$

Bottom-up seasonal mixing ($s_{-1}^\ell = 0$):

$$s_m^\ell \leftarrow s_m^\ell + \text{BottomUpMixing}(s_{m-1}^\ell), \quad m = 0, \dots, M.$$

Top-down trend mixing ($t_{M+1}^\ell = 0$):

$$t_m^\ell \leftarrow t_m^\ell + \text{TopDownMixing}(t_{m+1}^\ell), \quad m = M, \dots, 0.$$

Scale-wise residual with cross-scale aggregation:

$$X^\ell = X^{\ell-1} + \text{FFN}(\text{SMix}(\{s_m^\ell\}_{m=0}^M) + \text{TMix}(\{t_m^\ell\}_{m=0}^M)).$$

Per-scale predictors and future multi-predictor mixing:

$$\hat{\mathbf{x}}_m = \text{Pred}_m(X_m^L) \in \mathbb{R}^T, \quad p_{\theta_1}(X^0) = \text{FMM}(\{\hat{\mathbf{x}}_m\}_{m=0}^M) \in \mathbb{R}^T.$$

TimeMixer++ (Wang et al., 2025).

- **Choice of K .** $K = L$, which is the one generator–processor pair per MixerBlock.
- **Generator g_i (recursive downsampling, multi-resolution time imaging).**

$$\mathbf{x}_m = \text{Conv}(\mathbf{x}_{m-1}, \text{stride} = 2) \in \mathbb{R}^{\lfloor T/2^m \rfloor \times C}, \quad m = 1, \dots, M, \quad X_{\text{init}} = \{\mathbf{x}_0, \dots, \mathbf{x}_M\}.$$

Input projection (channel mixing & embedding).

$$\mathbf{x}_M \leftarrow \text{Channel-Attn}(Q_M, K_M, V_M),$$

where

$$X^0 = \{\mathbf{x}_m^0\}_{m=0}^M = \text{Embed}(X_{\text{init}}), \quad \mathbf{x}_m^0 \in \mathbb{R}^{\lfloor T/2^m \rfloor \times d_{\text{model}}}.$$

Per-block imaging. From $X^{i-1} = \{\mathbf{x}_m^{i-1}\}$, estimate top- K frequencies at the coarsest scale

$$A, \{f_k\}_{k=1}^K, \{p_k\}_{k=1}^K = \text{FFT}(\mathbf{x}_M^{i-1}), \quad p_k = \frac{T}{f_k},$$

and for each m, k ,

$$z_m^{(i,k)} = \text{Reshape}_{1\text{D} \rightarrow 2\text{D}}(\text{Padding}_{m,k}(\mathbf{x}_m^{i-1})) \in \mathbb{R}^{p_k \times f_{m,k} \times d_{\text{model}}}, \quad f_{m,k} = \left\lfloor \frac{\lfloor T/2^m \rfloor}{p_k} \right\rfloor,$$

so

$$g_i(X^{i-1}) = \{z_m^{(i,k)} : m = 0, \dots, M; k = 1, \dots, K\}.$$

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

- **Processor** p_{θ_i} (**MixerBlock**).

$$X^i = \text{LayerNorm}(X^{i-1} + \text{MixerBlock}_{\theta_i}(g_i(X^{i-1}))),$$

with

$$\text{(TID)} \quad s_m^{(i,k)} = \text{Attention}_{\text{col}}(\cdot), \quad t_m^{(i,k)} = \text{Attention}_{\text{row}}(\cdot),$$

$$\text{(MCM; seasonal bottom-up)} \quad s_m^{(i,k)} \leftarrow s_m^{(i,k)} + \text{Conv2D}_{\downarrow 2}(s_{m-1}^{(i,k)}),$$

$$\text{(MCM; trend top-down)} \quad t_m^{(i,k)} \leftarrow t_m^{(i,k)} + \text{TransConv2D}_{\uparrow 2}(t_{m+1}^{(i,k)}),$$

$$\text{(2D} \rightarrow \text{1D)} \quad z_m^{(i,k)} \leftarrow \text{Reshape}_{2\text{D} \rightarrow 1\text{D}}(s_m^{(i,k)} + t_m^{(i,k)}),$$

$$\text{(MRM; adaptive period fusion)} \quad \widehat{A}^{f_k} = \text{Softmax}(A^{f_1}, \dots, A^{f_K})_k, \quad \mathbf{x}_m^i = \sum_{k=1}^K \widehat{A}^{f_k} \odot z_m^{(i,k)},$$

- *Output projection.* With $X^L = \{\mathbf{x}_m^L\}_{m=0}^M$,

$$\widehat{\mathbf{y}}_m = \text{Head}_m(\mathbf{x}_m^L) \in \mathbb{R}^T, \quad \widehat{\mathbf{y}} = \text{Ensemble}(\{\widehat{\mathbf{y}}_m\}_{m=0}^M) \in \mathbb{R}^T.$$

AMD (Hu et al., 2025a).

- **Choice of K .** $K = 1$.
- **Generator g_1 (multi-scale moving-average, downsampling).** For $\mathbf{x} \in \mathbb{R}^T$ and rate $d \geq 2$,

$$\tau^{(1)} = \mathbf{x}, \quad \tau^{(i)} = \text{AvgPool}_{k_i}(\tau^{(i-1)}) \in \mathbb{R}^{\lfloor T/d^{i-1} \rfloor}, \quad i = 2, \dots, h,$$

and define $g_1(\mathbf{x}) = \{\tau^{(i)}\}_{i=1}^h$.

- **Processor** p_{θ_1} (**MDM, DDI, AMS**).

Coarse-to-fine residual mixing:

$$\xi^{(h)} = \tau^{(h)}, \quad \xi^{(i)} = \tau^{(i)} + \text{MLP}_i(\xi^{(i+1)}), \quad i = h-1, \dots, 1,$$

and set $u = \xi^{(1)} \in \mathbb{R}^T$.

DDI: with patches $\widehat{U} = \mathcal{P}_P(u) \in \mathbb{R}^{n_P \times P}$,

$$Z = \widehat{U} + \text{MLP}_t(\widehat{U}), \quad \widehat{V} = Z + \beta [\text{MLP}_c(Z^\top)]^\top,$$

then $v = \text{Unpatch}(\widehat{V}) \in \mathbb{R}^T$.

AMS (dense MoE): selector

$$S = \text{Softmax}(\text{TopK}(\text{Softmax}(Q(u)), k)),$$

$$Q(u) = \text{Decomp}(u) + \psi \text{Softplus}(\text{Decomp}(u)W_{\text{noise}}),$$

with $\psi \sim \mathcal{N}(0, 1)$ and

$$\text{TopK}(z, k) = \begin{cases} \alpha \log(z + 1), & z < v_k, \\ \alpha (e^z - 1), & z \geq v_k, \end{cases}$$

then

$$\widehat{\mathbf{y}} = \sum_{j=1}^m S_j \text{Predictor}_j(v) \in \mathbb{R}^T.$$

WPMixer (Murad et al., 2025).

- **Choice of K .** $K = 1$.
- **Generator g_1 (wavelet decomposition; time-frequency multiresolution).** Given $\mathbf{x} \in \mathbb{R}^{C \times L}$, perform an m -level DWT with mother wavelet ψ :

$$[\mathbf{x}_m^A, \mathbf{x}_m^D, \mathbf{x}_{m-1}^D, \dots, \mathbf{x}_1^D] = \text{Decomp}(\mathbf{x}, \psi, m),$$

where $\mathbf{x}_i^A, \mathbf{x}_i^D \in \mathbb{R}^{C \times L_i}$ are approximation/detail series at level i . Retain

$$\mathcal{W} = \{\mathbf{x}_i^W\}_{i=1}^m, \quad \mathbf{x}_m^W = \mathbf{x}_m^A, \quad \mathbf{x}_i^W = \mathbf{x}_i^D \quad (i < m), \quad g_1(\mathbf{x}) = \{\mathbf{x}_1^W, \dots, \mathbf{x}_m^W\}.$$

1134 • **Processor** p_{θ_1} **per-resolution branch.** For each $i \in \{1, \dots, m\}$:

$$\begin{aligned} 1135 & \tilde{\mathbf{x}}_i^W = \text{RevIN-Norm}(\mathbf{x}_i^W), \\ 1136 & \mathbf{x}_i^P = \text{Patch}(\tilde{\mathbf{x}}_i^W) \in \mathbb{R}^{C \times N_i \times P}, \\ 1137 & \mathbf{x}_i^d = \text{Embedding}(\mathbf{x}_i^P) \in \mathbb{R}^{C \times N_i \times d}. \end{aligned}$$

1140 *Patch Mixer:*

$$1141 \quad \mathbf{x}'_i = P(\text{BN}(\mathbf{x}_i^d)) \in \mathbb{R}^{d \times C \times N_i}, \quad \mathbf{x}''_i = \text{L}_2(\text{GELU}(\text{L}_1(\mathbf{x}'_i))).$$

1142 *Embedding Mixer:*

$$1143 \quad \hat{\mathbf{x}}_i = \text{BN}(P(\mathbf{x}''_i)) \in \mathbb{R}^{C \times N_i \times d}, \quad Y_i^d = \hat{\mathbf{x}}_i + \text{L}'_2(\text{GELU}(\text{L}'_1(\hat{\mathbf{x}}_i))).$$

1144 *Head:*

$$1145 \quad Y_i^f = \text{Flatten}(Y_i^d) \in \mathbb{R}^{C \times (N_i d)}, \quad Y_i^h = \text{Linear}(Y_i^f) \in \mathbb{R}^{C \times T_i}.$$

1146 *Reconstruction and denormalization:*

$$1147 \quad Y = \text{Reconstruction}_{\psi}(Y_m^A, Y_m^D, \dots, Y_1^D) \in \mathbb{R}^{C \times T}, \quad \hat{\mathbf{x}} = \text{RevIN-Denorm}(Y) \in \mathbb{R}^{T \times C}.$$

1148 **MultiPatchFormer (Naghashi et al., 2025).**

- 1151 • **Choice of K .** $K = S$, which is the number of stages.
 1152 • **Generator g_i (multi-scale patchifying).** Let $\mathbf{x}_{\text{in}} \in \mathbb{R}^L$. For scale $i \in \{1, \dots, S\}$ choose (P_i, S_i)
 1153 and apply

$$1154 \quad E^{(i)} = \text{Conv1D}(\text{Pad}(\mathbf{x}_{\text{in}}), \text{kernel} = P_i, \text{stride} = S_i, \text{out_ch} = d_i) \in \mathbb{R}^{P_N \times d_i},$$

1155 with scales chosen so P_N is identical across i . Define

$$1156 \quad g_i(\mathbf{x}_{\text{in}}) = E^{(i)}, \quad Z^{(i)} = \text{Concat}_d(Z^{(i-1)}, E^{(i)}) \in \mathbb{R}^{P_N \times d^{(i)}},$$

1157 where $d^{(i)} = \sum_{j=1}^i d_j$.

- 1158 • **Processor p_{θ_i} .** Temporal self-attention over P_N tokens:

$$1159 \quad Q_i = Z^{(i)} W_Q^{(i)}, \quad K_i = Z^{(i)} W_K^{(i)}, \quad V_i = Z^{(i)} W_V^{(i)} \in \mathbb{R}^{P_N \times d_h},$$

$$1160 \quad \text{MHA}_t(Z^{(i)}) = \left[\text{Concat}_h \text{Softmax}\left(\frac{Q_i K_i^\top}{\sqrt{d_h}}\right) V_i \right] W_O^{(i)} \in \mathbb{R}^{P_N \times d^{(i)}}.$$

1161 Fuse to width d_* and apply LN+FFN:

$$1162 \quad U^{(i)} = \Phi_i(\text{Concat}_d(\text{MHA}_t(Z^{(i)}), H^{(i-1)})) \in \mathbb{R}^{P_N \times d_*}, \quad H^{(i)} = \text{FFN}(\text{LN}(U^{(i)} + U^{(i)})).$$

1163 Decoder at $i = S$: reduce to $z = \frac{1}{P_N} \sum_{p=1}^{P_N} H_p^{(S)}$ and, with $T = M\Delta$,

$$1164 \quad \hat{y}^{(1)} = \mathcal{D}_1(z) \in \mathbb{R}^\Delta, \quad \hat{y}^{(m)} = \mathcal{D}_m(z, \text{Concat}(\hat{y}^{(1)}, \dots, \hat{y}^{(m-1)})) \in \mathbb{R}^\Delta,$$

$$1165 \quad p_{\theta_i} = \begin{cases} (\text{MHA}_t \& \text{Fuse} \& \text{LN+FFN}), & i < S, \\ (\text{MHA}_t \& \text{Fuse} \& \text{LN+FFN} \& \text{MultiStepDecode}), & i = S. \end{cases}$$

1166 In every case, g_i yields a scale-specific view and p_{θ_i} performs learned processing within/across
 1167 scales, establishing Definition 4.1.

1168 □

1181 A.4 PROOF OF THEOREM 5.2

1182 *Proof.* We verify (i) *differentiability* in s and (ii) *consistency* to a (discrete) scaling operator family.

1183 **Differentiability in s .** For each fixed pair (i, j) , the mapping $s_i \mapsto e^{-s_i} I_{|i-j|}(s_i)$ is real-analytic
 1184 on \mathbb{R}_+ because e^{-s_i} and the modified Bessel function $I_\nu(s_i)$ are analytic for $s_i > 0$ and integer
 1185 ν . Hence $s \mapsto \mathbf{K}(s)$ is C^∞ on \mathbb{R}_+^M , and so is $s \mapsto k(\mathbf{x}|s) = \mathbf{K}(s)\mathbf{x}$ for any \mathbf{x} . This proves the
 1186 differentiability requirement in Definition 5.1.

1188 **Consistency to a discrete scaling operator family.** Fix $s \in \mathbb{Z}_+$ and set $\mathbf{s} = s\mathbf{1}$. Then $\mathbf{K}(s\mathbf{1})$
 1189 becomes space-invariant:

$$1190 [\mathbf{K}(s\mathbf{1})]_{i,j} = e^{-s} I_{|i-j|}(s) =: g_s(i-j),$$

1191 so $k(\mathbf{x}|s\mathbf{1})$ is the convolution of \mathbf{x} with the kernel $g_s \in \ell^1(\mathbb{Z})$. Two classical properties hold:
 1192

1193 (a) *Semigroup and normalization.* For all $s, t \geq 0$, $g_{s+t} = g_s * g_t$, and $\sum_{n \in \mathbb{Z}} g_s(n) = 1$. Thus
 1194 $\mathbf{K}(s\mathbf{1})$ is a symmetric, doubly-stochastic Markov operator.

1195 (b) *Frequency response and contraction.* Let $\widehat{g}_s(\omega)$ denote the discrete-time Fourier transform
 1196 (DTFT) of g_s . A standard computation gives

$$1197 \widehat{g}_s(\omega) = \exp(s(\cos \omega - 1)) = \exp(-2s \sin^2(\omega/2)), \quad \omega \in [-\pi, \pi],$$

1198 hence $|\widehat{g}_s(\omega)| \leq 1$ with equality only at $\omega = 0$ for $s > 0$.

1200 We now verify the two axioms in Definition 3.1 for the discrete family $f(\mathbf{x}|s) := k(\mathbf{x}|s\mathbf{1})$:
 1201

1202 *Non-expansiveness.* By Plancherel,

$$1203 \|f(\mathbf{x}|s) - f(\mathbf{x}'|s)\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |\widehat{g}_s(\omega)|^2 |\widehat{\mathbf{x}} - \widehat{\mathbf{x}}'|^2 d\omega \leq \|\mathbf{x} - \mathbf{x}'\|_2^2,$$

1204 since $|\widehat{g}_s(\omega)| \leq 1$. Thus $\|f(\mathbf{x}|s) - f(\mathbf{x}'|s)\|_2 \leq \|\mathbf{x} - \mathbf{x}'\|_2$.
 1205

1206 *Energy reduction.* Let $s_i, s_j \in \mathbb{Z}_+$ with $s_i = ms_j$ for some integer $m > 1$. Let $X(\omega)$ denote the
 1207 DTFT of \mathbf{x} ; then by Plancherel,

$$1208 \|f(\mathbf{x}|s)\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |\widehat{g}_s(\omega)|^2 |X(\omega)|^2 d\omega.$$

1209 Using the semigroup property,

$$1210 g_{s_i} = g_{s_j}^{*m} := \underbrace{g_{s_j} * \cdots * g_{s_j}}_{m \text{ times}},$$

1211 the frequency responses satisfy

$$1212 \widehat{g}_{s_i}(\omega) = (\widehat{g}_{s_j}(\omega))^m, \quad |\widehat{g}_{s_i}(\omega)|^2 = |\widehat{g}_{s_j}(\omega)|^{2m}.$$

1213 Since $|\widehat{g}_{s_j}(\omega)| \leq 1$ for all ω and $m > 1$, we have

$$1214 |\widehat{g}_{s_i}(\omega)|^2 = |\widehat{g}_{s_j}(\omega)|^{2m} \leq |\widehat{g}_{s_j}(\omega)|^2 \quad \text{for all } \omega.$$

1215 Therefore, for any \mathbf{x} ,

$$1216 \|f(\mathbf{x}|s_i)\|_2^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |\widehat{g}_{s_i}(\omega)|^2 |X(\omega)|^2 d\omega$$

$$1217 \leq \frac{1}{2\pi} \int_{-\pi}^{\pi} |\widehat{g}_{s_j}(\omega)|^2 |X(\omega)|^2 d\omega$$

$$1218 = \|f(\mathbf{x}|s_j)\|_2^2,$$

1219 which proves the *energy reduction inequality*

$$1220 \|f(\mathbf{x}|s_i)\|_2 \leq \|f(\mathbf{x}|s_j)\|_2 \quad \text{for all } \mathbf{x} \quad \text{whenever } s_i = ms_j, m > 1.$$

1221 To see that the inequality is strict for some \mathbf{x} , note that for $s_j > 0$ we have

$$1222 |\widehat{g}_{s_j}(\omega)| < 1 \quad \text{for all } \omega \neq 0,$$

1223 and hence

$$1224 |\widehat{g}_{s_i}(\omega)|^2 = |\widehat{g}_{s_j}(\omega)|^{2m} < |\widehat{g}_{s_j}(\omega)|^2 \quad \text{for all } \omega \neq 0.$$

1225 Choose any \mathbf{x} whose spectrum is not supported solely at $\omega = 0$, i.e., such that $|X(\omega)|^2 > 0$ on a set
 1226 of nonzero measure in $\{\omega \in [-\pi, \pi] : \omega \neq 0\}$. Then on a set of positive measure,

$$1227 |\widehat{g}_{s_i}(\omega)|^2 |X(\omega)|^2 < |\widehat{g}_{s_j}(\omega)|^2 |X(\omega)|^2,$$

1228 so the inequality above is strict:

$$1229 \|f(\mathbf{x}|s_i)\|_2 < \|f(\mathbf{x}|s_j)\|_2.$$

1230 Thus the LDG kernel family satisfies both the universal inequality and the existence of a strict
 1231 example required by the energy reduction condition in Definition 3.1.

1241 □

A.5 PROOF OF THEOREM 5.3

A.5.1 DISCRETE SCALE-SPACE AXIOMS

Axioms (discrete, 1D, symmetric case). For each scale $s \geq 0$, let $\mathcal{T}_s : \ell^2(\mathbb{Z}) \rightarrow \ell^2(\mathbb{Z})$ denote the smoothing operator at scale s . A family $\{\mathcal{T}_s\}_{s \geq 0}$ satisfies the *discrete scale-space axioms* if:

- (A1) **Linearity and shift invariance.** \mathcal{T}_s is linear and commutes with shifts: there exists $K_s \in \ell^1(\mathbb{Z})$ with $(\mathcal{T}_s x)[n] = (K_s * x)[n]$.
- (A2) **Semigroup over scale and identity at 0.** $\mathcal{T}_0 = \text{Id}$ and $\mathcal{T}_{s+t} = \mathcal{T}_s \circ \mathcal{T}_t$, i.e., $K_{s+t} = K_s * K_t$ for all $s, t \geq 0$.
- (A3) **Regularity in scale.** The map $s \mapsto K_s$ is C^1 in s in the ℓ^1 topology.
- (A4) **DC normalization.** $\sum_{n \in \mathbb{Z}} K_s[n] = 1$ for all $s \geq 0$ (constants are preserved).
- (A5) **Non-enhancement of local extrema (NELE).** If $y(\cdot; s) = \mathcal{T}_s x$ has a local maximum at index n , then $\partial_s y[n; s] \leq 0$ (resp. $\partial_s y[n; s] \geq 0$ at local minima).
- (A6) **Symmetry.** $K_s[n] = K_s[-n]$ for all $n \in \mathbb{Z}$ and $s \geq 0$.

A.5.2 UNIQUENESS OF THE DISCRETE GAUSSIAN

Lemma A.3. *Let A be a translation-invariant, symmetric operator on $\ell^2(\mathbb{Z})$ with impulse response $a[\cdot] \in \ell^1(\mathbb{Z})$ such that*

$$\sum_{n \in \mathbb{Z}} a[n] = 0, \quad a[n] \geq 0 \text{ for } n \neq 0, \quad a[-n] = a[n].$$

Then its Fourier symbol $\phi(\omega) = \hat{A}(\omega) = \sum_n a[n] e^{-in\omega}$ admits the Lévy-Khintchine form

$$\phi(\omega) = \sum_{m=1}^{\infty} c_m (\cos(m\omega) - 1), \quad c_m \geq 0,$$

with real coefficients $c_m = 2a[m]$. Conversely, given any sequence $\{c_m\}_{m \geq 1}$ with $c_m \geq 0$ and $\sum_{m \geq 1} c_m < \infty$, defining

$$a[0] = -\sum_{m \geq 1} c_m, \quad a[\pm m] = \frac{1}{2} c_m \quad (m \geq 1),$$

yields a symmetric Toeplitz operator A with the above properties and symbol $\phi(\omega) = \sum_{m \geq 1} c_m (\cos m\omega - 1)$.

Proof. Since $a \in \ell^1(\mathbb{Z})$ and $a[-n] = a[n]$,

$$\phi(\omega) = \sum_{n \in \mathbb{Z}} a[n] e^{-in\omega} = a[0] + 2 \sum_{m=1}^{\infty} a[m] \cos(m\omega).$$

The zero-sum condition gives $a[0] = -2 \sum_{m \geq 1} a[m]$, hence

$$\phi(\omega) = 2 \sum_{m=1}^{\infty} a[m] (\cos(m\omega) - 1) = \sum_{m=1}^{\infty} c_m (\cos(m\omega) - 1), \quad c_m := 2a[m] \geq 0.$$

Absolute summability of a implies $\sum_m c_m < \infty$, ensuring uniform convergence and continuity of ϕ . The converse follows by reversing the construction. \square

Theorem A.4 (Uniqueness of the discrete Gaussian scale space). *Among symmetric kernel families, the discrete Gaussian*

$$K_s[n] = e^{-\alpha s} I_{|n|}(\alpha s), \quad \alpha > 0,$$

is the unique extended scaling operator family that satisfies (A1)–(A6). Equivalently,

$$\hat{K}_s(\omega) = \exp(\alpha s (\cos \omega - 1)) = \exp(-2\alpha s \sin^2(\omega/2)), \quad \omega \in [-\pi, \pi],$$

and any other family satisfying the axioms coincides with K_s up to reparameterizing scale $s \mapsto \alpha s$.

1296 *Proof.* By **(A1)**, $(\mathcal{T}_s x)[n] = (K_s * x)[n]$ for some $K_s \in \ell^1(\mathbb{Z})$. Let $\widehat{K}_s(\omega)$ be its DTFT. By
 1297 **(A2)–(A3)**, for each fixed ω the map $s \mapsto \widehat{K}_s(\omega)$ is a continuous one-parameter semigroup with
 1298 $\widehat{K}_0(\omega) = 1$, hence there exists a real, even *generator* $\phi(\omega)$ such that
 1299

$$1300 \quad \widehat{K}_s(\omega) = \exp(s\phi(\omega)), \quad s \geq 0.$$

1301 From **(A4)**, $\phi(0) = 0$. From **(A6)**, ϕ is real and even.

1302 Let $y(\cdot; s) = \mathcal{T}_s x$. Differentiating in s ,

$$1303 \quad \partial_s y = A * x, \quad A := \partial_s K_s \Big|_{s=0},$$

1304 so that $\widehat{A}(\omega) = \phi(\omega)$. By **(A4)** and **(A6)**, A is symmetric with zero row sum. The NELE axiom
 1305 **(A5)** at vanishing scale enforces the discrete maximum principle for A : off-diagonal coefficients are
 1306 nonnegative while the diagonal is nonpositive, and rows sum to zero. Thus A fits the hypotheses of
 1307 Lemma A.3, and
 1308

$$1309 \quad \phi(\omega) = \sum_{m=1}^{\infty} c_m (\cos(m\omega) - 1), \quad c_m \geq 0.$$

1310 We now show that **(A5)** forces $c_m = 0$ for all $m \geq 2$. Consider signals supported on three consec-
 1311 utive sites and apply **(A5)** at $s = 0$ to both local maxima and minima. A standard extremum test
 1312 yields that any positive coefficient at distance $m \geq 2$ would produce, for sufficiently small $s > 0$, an
 1313 increase at a newly formed off-center extremum before nearest neighbors equilibrate. Hence $c_m = 0$
 1314 for $m \geq 2$, and

$$1315 \quad \phi(\omega) = c_1 (\cos \omega - 1) =: \alpha (\cos \omega - 1), \quad \alpha := c_1 > 0.$$

1316 Substituting into equation A.5.2 gives

$$1317 \quad \widehat{K}_s(\omega) = \exp(\alpha s (\cos \omega - 1)).$$

1318 Taking the inverse DTFT yields $K_s[n] = e^{-\alpha s} I_{|n|}(\alpha s)$, the discrete Gaussian kernel. This family
 1319 satisfies **(A1)–(A6)**; conversely, any other family obeying the axioms has the same generator up to
 1320 the multiplicative constant α , i.e., a reparameterization of scale. The result follows. \square
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

Table 5: Dataset descriptions for long-term and short-term forecasting benchmarks. The dataset size is organized as (Train, Validation, Test). For long-term tasks, we adopt multivariate datasets covering diverse domains. For short-term tasks, we follow the official M4 benchmark, which consists of univariate series at different temporal frequencies. Following TimeMixer++ (Wang et al., 2025), forecastability is measured as one minus spectral entropy (Goerg, 2013), with higher values indicating better predictability.

Tasks	Dataset	Dim	Series Length	Dataset Size	Domain	Frequency	Forecast.
Forecasting (Long-term)	ETTh1	7	{96, 192, 336, 720}	(34465, 11521, 11521)	Temperature	15 min	0.46
	ETTh2	7	{96, 192, 336, 720}	(34465, 11521, 11521)	Temperature	15 min	0.55
	ETTm1	7	{96, 192, 336, 720}	(8545, 2881, 2881)	Temperature	Hourly	0.38
	ETTm2	7	{96, 192, 336, 720}	(8545, 2881, 2881)	Temperature	Hourly	0.45
	Electricity	321	{96, 192, 336, 720}	(18317, 2633, 5261)	Electricity	Hourly	0.77
	Traffic	862	{96, 192, 336, 720}	(12185, 1757, 3509)	Transportation	Hourly	0.68
	Exchange	8	{96, 192, 336, 720}	(5120, 665, 1422)	Exchange rate	Daily	0.41
	Weather	21	{96, 192, 336, 720}	(36792, 5271, 10540)	Weather	10 min	0.75
Forecasting (Short-term)	M4-Yearly	1	6	(23000, 0, 23000)	Demographic	Yearly	0.43
	M4-Quarterly	1	8	(24000, 0, 24000)	Finance	Quarterly	0.47
	M4-Monthly	1	18	(48000, 0, 48000)	Industry	Monthly	0.44
	M4-Weekly	1	13	(359, 0, 359)	Macro	Weekly	0.43
	M4-Daily	1	14	(4227, 0, 4227)	Micro	Daily	0.44
	M4-Hourly	1	48	(414, 0, 414)	Other	Hourly	0.46

Table 6: Experiment configurations of SIGMA across datasets. All experiments adopt the ADAM (Kinga et al., 2015) optimizer. We report the number of blocks (K), model dimension (d_{model}), initial learning rate (LR), loss function, batch size, and training epochs.

Dataset	K	d_{model}	LR	Loss	Batch Size	Epochs
ETTh1	1	32	0.0005	MSE	32	10
ETTh2	1	32	0.0005	MSE	32	10
ETTm1	1	8	0.02	MSE	32	10
ETTm2	1	16	0.0005	MSE	32	10
Exchange	1	16	0.0001	MSE	32	10
Electricity	1	16	0.01	MSE	32	10
Traffic	1	16	0.01	MSE	32	10
Weather	1	8	0.005	MSE	32	10
M4	1	16	0.01	SMAPE	32	10

B EXPERIMENTAL DETAILS

For the evaluation of forecasting models, we follow the protocol used in TimesNet Wu et al. (2023). For **long-term forecasting** (see Table 2), we report the mean square error (MSE) and mean absolute error (MAE). For **short-term forecasting** (see Table 3), we adopt symmetric mean absolute percentage error (SMAPE), mean absolute scaled error (MASE), and overall weighted average (OWA).

These metrics are defined as follows:

$$\text{SMAPE} = \frac{200}{H} \sum_{i=1}^T \frac{|\mathbf{X}_i - \hat{\mathbf{X}}_i|}{|\mathbf{X}_i| + |\hat{\mathbf{X}}_i|}, \quad (6)$$

$$\text{MASE} = \frac{1}{T} \sum_{i=1}^T \frac{|\mathbf{X}_i - \hat{\mathbf{X}}_i|}{\frac{1}{T-m} \sum_{j=m+1}^T |\mathbf{X}_j - \mathbf{X}_{j-m}|}, \quad (7)$$

$$\text{OWA} = \frac{1}{2} \left[\frac{\text{SMAPE}}{\text{SMAPE}_{\text{Naive2}}} + \frac{\text{MASE}}{\text{MASE}_{\text{Naive2}}} \right], \quad (8)$$

where m denotes the seasonal periodicity of the data, and $\mathbf{X} \in \mathbb{R}^{T \times C}$ and $\hat{\mathbf{X}} \in \mathbb{R}^{T \times C}$ represent the ground truth and predictions for T future time steps with C dimensions.

For methods that do not originally provide results on the M4 dataset, we rely on their official implementations and conduct a controlled hyperparameter search to ensure a fair comparison. All implementations are based on PyTorch (Paszke et al., 2019), with the modified Bessel function implemented via SciPy, and all experiments are executed on a single **NVIDIA RTX A6000** GPU. We summarize the datasets used in our experiments in Table 5, covering both long-term and short-term

Table 7: Standard deviation for SIGMA and the second-best method (AMD) on long-term forecasting datasets.

Method	SIGMA		AMD (2025a)	
Dataset	MSE	MAE	MSE	MAE
Weather	0.247 \pm 0.003	0.273 \pm 0.002	0.283 \pm 0.001	0.282 \pm 0.001
Electricity	0.175 \pm 0.001	0.269 \pm 0.001	0.208 \pm 0.000	0.303 \pm 0.000
Traffic	0.458 \pm 0.001	0.302 \pm 0.001	0.546 \pm 0.000	0.344 \pm 0.000
Exchange	0.353 \pm 0.001	0.400 \pm 0.001	0.358 \pm 0.009	0.401 \pm 0.002
ETTh1	0.443 \pm 0.004	0.433 \pm 0.002	0.447 \pm 0.004	0.434 \pm 0.003
ETTh2	0.376 \pm 0.003	0.402 \pm 0.002	0.376 \pm 0.022	0.400 \pm 0.005
ETTM1	0.383 \pm 0.003	0.397 \pm 0.002	0.395 \pm 0.003	0.399 \pm 0.001
ETTM2	0.276 \pm 0.001	0.322 \pm 0.001	0.285 \pm 0.043	0.328 \pm 0.012

Table 8: Standard deviation for SIGMA and the second-best method (MultiPatchFormer) on the short-term forecasting dataset (M4).

Method	SIGMA			MultiPatchFormer (2025)		
Dataset	SMAPE	MASE	OWA	SMAPE	MASE	OWA
Yearly	13.314 \pm 0.022	2.989 \pm 0.015	0.783 \pm 0.002	13.296 \pm 0.012	3.009 \pm 0.008	0.785 \pm 0.001
Quarterly	10.060 \pm 0.052	1.177 \pm 0.009	0.886 \pm 0.005	10.166 \pm 0.008	1.178 \pm 0.003	0.892 \pm 0.002
Monthly	12.750 \pm 0.019	0.936 \pm 0.006	0.882 \pm 0.003	12.810 \pm 0.016	0.942 \pm 0.001	0.887 \pm 0.001
Others	4.867 \pm 0.046	3.316 \pm 0.050	1.037 \pm 0.011	4.849 \pm 0.037	3.271 \pm 0.023	1.028 \pm 0.007
Averaged	11.840 \pm 0.016	1.585 \pm 0.009	0.868 \pm 0.002	11.889 \pm 0.010	1.591 \pm 0.001	0.872 \pm 0.000

forecasting tasks, and provide the detailed experimental settings, including model hyper-parameters and training setups, in Table 6.

C PROPERTIES OF SCALING OPERATOR FAMILY

Definition 3.1 characterizes a valid scaling operator through two essential properties: (1) non-expansiveness, requiring that applying the operator does not amplify differences between nearby inputs, and (2) energy reduction, requiring that coarser scales retain strictly less signal energy over multiplicable scales. These properties ensure that scaling progressively smooths the input while preserving stability.

Figure 6 reports the empirical behavior of all operator families in Table 1 across seven datasets. For every dataset, the induced output differences remain strictly smaller than the input differences, confirming non-expansiveness. Likewise, the average energy decreases monotonically over dyadic scales, indicating that each operator consistently removes high-frequency variation as the scale increases. Taken together, these results demonstrate that all operator families satisfy both conditions of Definition 3.1 not only in theory but also in practice, across diverse real-world time series.

D ERROR BARS

To assess the robustness of our experiments, we report the mean performance and standard deviation for SIGMA and the second-best baselines in Tables 7 and 8. On the long-term forecasting benchmarks, SIGMA achieves lower average MSE and MAE than AMD on nearly all datasets, while maintaining sufficiently small standard deviations.

On the short-term M4 benchmark, the averaged SMAPE, MASE, and OWA scores of SIGMA remain consistently better, and the corresponding standard deviations are of similar or smaller magnitude. Taken together, these error-bar analyses confirm that the gains reported by SIGMA are stable across runs and not attributable to random fluctuations.

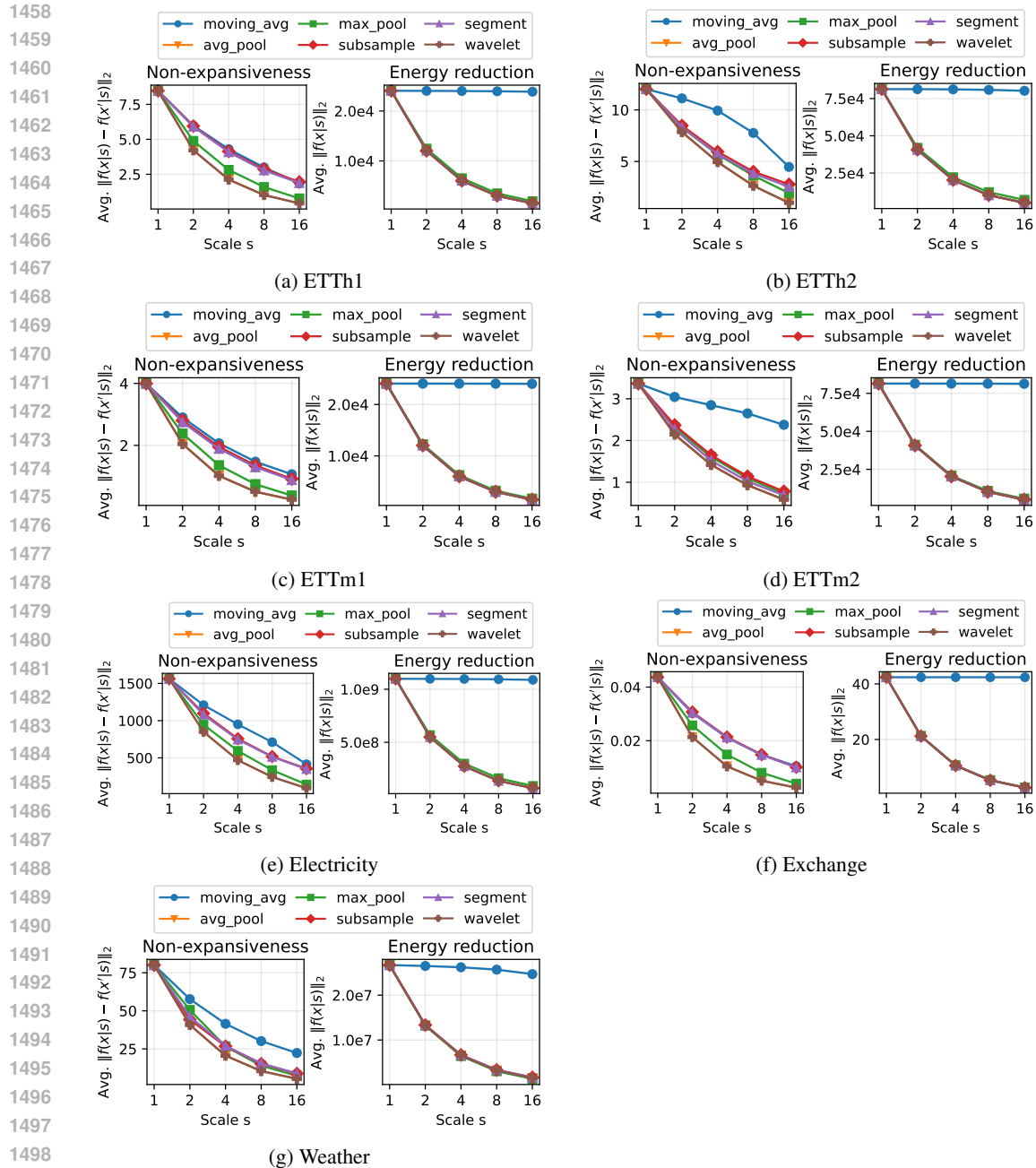
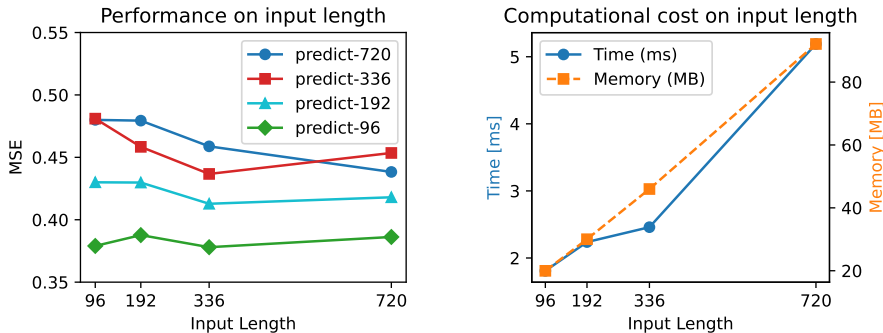


Figure 6: Energy reduction and non-expansiveness of the scaling operator families across all datasets.

E HYPERPARAMETER SENSITIVITY

Although we follow the evaluation protocol in TimesNet, the input length L remains a key hyperparameter in multi-scale forecasting models. To understand its impact on both accuracy and efficiency, we measure forecasting performance (MSE) and computational overhead on ETTh1 under varying input lengths $L \in \{96, 192, 336, 720\}$. As shown in Figure 7, longer lookback windows generally lead to better forecasting performance, as the model can leverage richer historical patterns. We observe that the growth in both time and memory is approximately linear, indicating that the model scales efficiently with the window size.

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526



(a) Performance across input lengths on ETTh1. The overall trend shows improved MSE with longer lookback windows.

(b) Computational cost across of input lengths on ETTh1. Both metrics increase approximately linearly with respect to L .

1527 **Figure 7:** Hyperparameter sensitivity with respect to the input length L on ETTh1. We report
1528 forecasting performance (MSE) and computational cost to analyze the trade-off between accuracy
1529 and efficiency under varying lookback windows.

1530
1531
1532
1533
1534
1535
1536
1537

This efficiency arises because the LDG kernel is defined by a symmetric Toeplitz matrix whose entries depend only on the relative distance $|i - j|$. This structure implies that applying the LDG kernel to an input sequence is equivalent to performing a 1D convolution with a distance-dependent kernel vector, rather than multiplying by a full $L \times L$ matrix. Consequently, the operation can be computed in $O(L)$ time instead of $O(L^2)$, which explains why runtime and memory footprint grow approximately linearly with the input length.

1538 F DESIGN GUIDE FOR EXTENDED SCALING OPERATOR FAMILIES

1539
1540
1541
1542
1543

The extended scaling operator family proposed in Section 5.1 provides a flexible framework for constructing continuous, input-dependent multi-scale transformations. Here we summarize practical guidelines for designing new instances of this family and illustrate how these principles are instantiated in our implementation.

1544
1545
1546
1547

Classical downsampling operators are defined only at integer scales. To make these operators learnable and input-adaptive, we extend the discrete index $j \in \{1, \dots, J\}$ to a continuous scale variable $s_t \in [1, J]$ at each position t . Instead of selecting a single integer operator, we blend nearby integer scales through a convex combination

1548
1549
1550

$$f(\mathbf{x}|s_t) = \sum_{j=1}^J w_j(s_t) f_j(\mathbf{x})_t,$$

1551
1552

where f_j is the original operator at scale j . This convex construction preserves non-expansiveness whenever each f_j is non-expansive, and it ensures that $f(\cdot|s)$ varies smoothly with respect to s .

1553
1554
1555

To satisfy Definition 3.1, the extended operator must reproduce the original discrete operator exactly when s_t is an integer. This is achieved by choosing weights $\{w_j(s)\}$ such that

1556

$$w_j(k) = \mathbf{1}\{j = k\} \quad \text{for all integers } k \in \{1, \dots, J\}.$$

1557

A convenient construction uses C^∞ bump functions:

1558
1559

$$w_j(s) = \frac{\phi(s-j)}{\sum_{k=1}^J \phi(s-k)}, \quad \phi(u) = \begin{cases} \exp(-1/(1-u^2)), & |u| < 1, \\ 0, & \text{otherwise.} \end{cases}$$

1560
1561
1562

At integer $s = k$, the only nonzero bump is $\phi(0)$, which yields $w_k(k) = 1$ and $w_{j \neq k}(k) = 0$. Thus the extended operator reduces exactly to the classical operator at integer scales, ensuring consistency.

1563
1564
1565

Because the bump function $\phi(\cdot)$ is C^∞ and the normalization preserves smoothness, the weights $w_j(s)$ are differentiable in s for all non-integer values. Since the output is a convex combination of the $f_j(\mathbf{x})$, the entire operator $f(\mathbf{x}|s)$ is differentiable in s , enabling backpropagation through the scale field and supporting input-adaptive scale prediction via a learnable scale head.

G COMPARISON WITH NON-MULTI-SCALE BASELINES

We evaluate SIGMA against strong non-multi-scale baselines, including TimeFilter (Hu et al., 2025b), PatchTST (Nie et al., 2023), and DLinear (Zeng et al., 2023) for long-term forecasting and short-term forecasting. Across the long-term forecasting benchmarks, SIGMA generally achieves the best accuracy on datasets with fewer variables (e.g., ETT). As the number of variables increases (e.g., Traffic), TimeFilter obtains the strongest overall performance, owing to its patch-wise filtration mechanism that explicitly captures spatiotemporal relationships. This suggests that explicitly modeling multi-scale patterns in multivariate settings can provide additional benefits. For short-term forecasting, SIGMA also performs strongly across all prediction horizons, demonstrating that SIGMA remains effective under diverse prediction-length settings.

Table 9: Long-term forecasting results with four methods: SIGMA, TimeFilter, PatchTST, and DLinear.

Method		SIGMA (Ours)		TimeFilter (2025b)		PatchTST (2023)		DLinear (2023)	
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Weather	96	<u>0.160</u>	0.204	0.159	<u>0.205</u>	0.174	0.216	0.196	0.258
	192	<u>0.209</u>	0.248	0.207	<u>0.249</u>	0.220	0.257	0.237	0.296
	336	<u>0.270</u>	<u>0.293</u>	0.264	0.291	0.278	0.298	0.283	0.333
	720	0.348	<u>0.345</u>	0.343	0.343	0.354	0.347	<u>0.346</u>	0.383
	Avg	<u>0.247</u>	<u>0.273</u>	0.243	0.272	0.256	0.279	0.266	0.318
Electricity	96	<u>0.146</u>	<u>0.241</u>	0.138	0.235	0.180	0.273	0.210	0.301
	192	<u>0.163</u>	<u>0.257</u>	0.156	0.253	0.187	0.279	0.210	0.305
	336	<u>0.179</u>	<u>0.273</u>	0.170	0.268	0.204	0.295	0.223	0.319
	720	<u>0.213</u>	<u>0.304</u>	0.191	0.289	0.246	0.328	0.258	0.350
	Avg	<u>0.175</u>	<u>0.269</u>	0.164	0.261	0.204	0.294	0.225	0.319
Traffic	96	<u>0.431</u>	<u>0.288</u>	0.391	0.260	0.460	0.298	0.696	0.429
	192	<u>0.444</u>	<u>0.296</u>	0.413	0.269	0.467	0.302	0.646	0.407
	336	<u>0.461</u>	<u>0.303</u>	0.429	0.277	0.483	0.308	0.653	0.410
	720	<u>0.494</u>	<u>0.320</u>	0.462	0.296	0.516	0.325	0.694	0.429
	Avg	<u>0.458</u>	<u>0.302</u>	0.424	0.276	0.482	0.308	0.672	0.419
Exchange	96	<u>0.084</u>	<u>0.204</u>	0.083	0.202	0.087	0.204	0.095	0.227
	192	0.174	0.297	<u>0.178</u>	<u>0.299</u>	0.188	0.308	0.184	0.323
	336	0.322	0.411	0.332	<u>0.416</u>	0.341	0.423	<u>0.328</u>	0.434
	720	0.833	0.687	<u>0.785</u>	<u>0.669</u>	0.921	0.720	0.762	0.667
	Avg	0.353	<u>0.400</u>	<u>0.345</u>	0.397	0.384	0.414	0.342	0.413
ETTh1	96	0.379	0.393	<u>0.384</u>	<u>0.395</u>	0.379	0.398	0.396	0.410
	192	0.430	<u>0.425</u>	<u>0.438</u>	0.424	0.430	0.434	0.445	0.441
	336	<u>0.481</u>	0.446	<u>0.481</u>	0.446	0.473	<u>0.460</u>	0.493	0.471
	720	<u>0.480</u>	<u>0.468</u>	0.479	0.466	0.523	0.506	0.515	0.512
	Avg	0.443	0.433	<u>0.446</u>	0.433	0.451	<u>0.450</u>	0.462	0.459
ETTh2	96	0.289	0.339	<u>0.290</u>	<u>0.340</u>	0.300	0.351	0.346	0.399
	192	0.369	0.394	<u>0.377</u>	0.394	0.380	0.400	0.478	0.477
	336	0.415	0.427	<u>0.424</u>	<u>0.435</u>	0.431	0.441	0.597	0.543
	720	0.431	0.446	0.464	0.464	<u>0.447</u>	<u>0.461</u>	0.841	0.661
	Avg	0.376	0.402	<u>0.389</u>	<u>0.408</u>	0.390	0.413	0.566	0.520
ETTm1	96	<u>0.323</u>	<u>0.359</u>	0.320	0.357	0.330	0.368	0.345	0.372
	192	0.360	0.381	<u>0.362</u>	0.381	0.370	0.391	0.382	<u>0.390</u>
	336	<u>0.392</u>	<u>0.404</u>	0.391	0.402	0.402	0.411	0.414	0.414
	720	0.455	<u>0.442</u>	0.461	0.438	<u>0.459</u>	0.446	0.474	0.451
	Avg	0.383	<u>0.397</u>	0.383	0.394	<u>0.390</u>	0.404	0.404	0.407
ETTm2	96	<u>0.174</u>	0.257	0.172	<u>0.258</u>	0.183	0.265	0.195	0.295
	192	<u>0.239</u>	0.299	0.237	<u>0.300</u>	0.246	0.308	0.282	0.359
	336	0.296	0.337	0.296	<u>0.338</u>	<u>0.312</u>	0.350	0.363	0.414
	720	0.394	0.394	0.394	<u>0.396</u>	<u>0.419</u>	0.412	0.547	0.519
	Avg	<u>0.276</u>	0.322	0.275	<u>0.323</u>	0.290	0.334	0.347	0.397

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

Table 10: Short-term forecasting results in the M4 benchmark dataset using four methods: SIGMA, TimeFilter, PatchTST, and DLinear.

	Method	SIGMA (Ours)	TimeFilter (2025b)	PatchTST (2023)	DLinear (2023)
Yearly	SMAPE	13.314	18.836	<u>14.311</u>	14.343
	MASE	2.989	4.153	3.240	<u>3.123</u>
	OWA	0.783	1.099	<u>0.845</u>	0.832
Quarterly	SMAPE	10.060	10.660	<u>10.242</u>	10.502
	MASE	1.177	1.228	<u>1.211</u>	1.240
	OWA	0.886	0.932	<u>0.907</u>	0.929
Monthly	SMAPE	12.750	13.477	<u>12.889</u>	13.373
	MASE	0.936	1.025	<u>0.955</u>	1.004
	OWA	0.882	0.949	<u>0.896</u>	0.935
Others	SMAPE	<u>4.867</u>	6.136	4.986	5.110
	MASE	<u>3.316</u>	4.042	3.231	3.655
	OWA	<u>1.037</u>	1.257	1.023	1.132
Weighted Average	SMAPE	11.840	13.666	<u>12.186</u>	12.494
	MASE	1.585	1.944	<u>1.656</u>	1.681
	OWA	0.868	0.995	<u>0.893</u>	0.920