SKOLR: Structured Koopman Operator Linear RNN for Time-Series Forecasting

Yitian Zhang¹²³ Liheng Ma¹²³ Antonios Valkanas¹²³ Boris N. Oreshkin⁴ Mark Coates¹²³

Abstract

Koopman operator theory provides a framework for nonlinear dynamical system analysis and timeseries forecasting by mapping dynamics to a space of real-valued measurement functions, enabling a linear operator representation. Despite the advantage of linearity, the operator is generally infinite-dimensional. Therefore, the objective is to learn measurement functions that yield a tractable finite-dimensional Koopman operator approximation. In this work, we establish a connection between Koopman operator approximation and linear Recurrent Neural Networks (RNNs), which have recently demonstrated remarkable success in sequence modeling. We show that by considering an extended state consisting of lagged observations, we can establish an equivalence between a structured Koopman operator and linear RNN updates. Building on this connection, we present SKOLR, which integrates a learnable spectral decomposition of the input signal with a multilayer perceptron (MLP) as the measurement functions and implements a structured Koopman operator via a highly parallel linear RNN stack. Numerical experiments on various forecasting benchmarks and dynamical systems show that this streamlined, Koopman-theorybased design delivers exceptional performance. Our code is available at: https://github. com/networkslab/SKOLR.

1. Introduction

Time-series prediction and analysis of nonlinear dynamical systems remain fundamental challenges across various domains. Koopman operator theory (Koopman, 1931) offers a promising mathematical framework that transforms nonlinear dynamics into linear operations in the space of measurement functions. Practical implementation of this theory faces significant challenges due to the infinite dimensionality of the resulting linear operator, necessitating finite dimensional approximations. The dynamic mode decomposition (DMD) and its variants are the most widely employed approximations (Rowley et al., 2009; Schmid, 2010; Williams et al., 2015), although alternative techniques have emerged (Bevanda et al., 2021; Khosravi, 2023), including ones employing learnable neural measurement functions (Li et al., 2017)

In parallel developments, linear Recurrent Neural Networks (linear RNNs) have emerged as a powerful architecture in deep learning and sequence modeling (Stolzenburg et al., 2018; Gu & Dao, 2023; Wang et al., 2024b). These models leverage the computational efficiency of linear recurrence while maintaining impressive modeling capabilities.

In this work, we consider the task of time-series forecasting and establish both an explicit connection and a direct architectural match between Koopman operator approximation and linear RNNs. In particular, we show that by representing the dynamic state using a collection of time-delayed observations, we can establish an equivalence between the application of an extended DMD-style approximation of the Koopman operator and the state update of a linear RNN.

Building on this connection, we introduce Structured Koopman Operator Linear RNN (SKOLR) for time-series forecasting. SKOLR implements a structured Koopman operator through a highly parallel linear RNN stack. Through a learnable spectral decomposition of the input signal, the RNN chains jointly attend to different dynamical patterns from different representation subspaces, creating a theoretically-grounded yet computationally efficient design that naturally aligns with Koopman principles.

Through extensive experiments on various forecasting benchmarks and dynamical systems, we demonstrate that

¹Department of Electrical and Computer Engineering, McGill University, Montreal, Canada ²Mila - Quebec Artificial Intelligence Institute, Montreal, Canada ³ILLS - International Laboratory on Learning Systems, Montreal, Canada ⁴Amazon Science. This work does not relate to the author's position at Amazon. Correspondence to: Yitian Zhang <yitian.zhang@mail.mcgill.ca>, Liheng Ma <liheng.ma@mail.mcgill.ca>.

Proceedings of the 42^{nd} International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).

this streamlined, Koopman-theory-based design delivers exceptional performance, while maintaining the simplicity of the linear RNN and its outstanding parameter efficiency.

2. Preliminary

This section provides foundational background for the proposed forecasting methodology. We define the discrete-time dynamical systems used to model observed time series and then introduce the Koopman operator.

Definition 2.1 (Discrete-time Dynamical Systems). We consider the (autonomous) discrete-time dynamical system:

$$\mathbf{x}_{k+1} = \mathbf{F}(\mathbf{x}_k) \tag{1}$$

where $\mathbf{x}_k \in \mathcal{M}$ denotes the system state at time $k \in \mathbb{Z}^+$; and $F : \mathcal{M} \to \mathcal{M}$ represents the underlying dynamics mapping the state forward in time. We assume a Euclidean state space $\mathcal{M} \subset \mathbb{R}^C$, although it can be more generally defined on an *n*-dimensional manifold (Bevanda et al., 2021).

The Koopman operator framework enables globally linear representations of nonlinear systems by applying a linear operator to measurement functions (observables) g of the state \mathbf{x}_k . The following theorem formalizes key properties of the Koopman operator.

Theorem 2.2 (Koopman Operator Theorem (Koopman, 1931; Brunton et al., 2022)). *Considering real-valued* measurement functions (*a.k.a.* observables) $g : \mathcal{M} \to \mathbb{R}$, the Koopman operator $\mathcal{K} : \mathcal{F} \to \mathcal{F}$ is an infinite-dimensional linear operator on the space of all possible measurement functions \mathcal{F} , which is an infinite-dimensional Hilbert space, satisfying:

$$\mathcal{K} \circ g = g \circ \mathbf{F},\tag{2}$$

where \circ is the composition operator.

In other words,

$$\mathcal{K}(g(\mathbf{x}_k)) = g(\mathbf{F}(\mathbf{x}_k)) = g(\mathbf{x}_{k+1}).$$
(3)

This is true for any measurement function g *and for any state* \mathbf{x}_k .¹

While this facilitates analysis via linear maps, Koopman operator is generally infinite-dimensional, acting on a Hilbert space of functions. For practical learning and inference, we seek effective finite-dimensional approximations. In this paper, we construct these approximations efficiently by leveraging a connection to linear recurrent neural networks (RNNs). For clarity, we now define a linear RNN.

Definition 2.3 (Linear Recurrent Neural Network (Stolzenburg et al., 2018)). Consider a hidden state space $\mathcal{H} \subseteq \mathbb{R}^{d_h}$

and input space $\mathcal{V} \subseteq \mathbb{R}^{d_v}$. For any sequence $(\mathbf{v}_k)_{k=1}^L \in \mathcal{V}$, the linear RNN defines a discrete-time dynamical system through the hidden state transition equation:

$$\mathbf{h}_k = \mathbf{W}\mathbf{h}_{k-1} + \mathbf{U}\mathbf{v}_k + \mathbf{b} \tag{4}$$

where $\mathbf{W} \in \mathbb{R}^{d_h \times d_h}$ is the hidden state transition matrix, $\mathbf{U} \in \mathbb{R}^{d_h \times d_v}$ is the weight matrix applied to the input, and $\mathbf{b} \in \mathbb{R}^{d_h}$ is the bias vector. The evolution of hidden states $\mathbf{h}_k \in \mathcal{H}$ is uniquely determined by this linear map.

In order to prepare the connection to our proposed Koopman operator learning strategy and forecasting method, let us introduce $\mathbf{v}_k := \psi(\mathbf{y}_k)$ and define $g(\mathbf{y}_k) := \mathbf{U}\psi(\mathbf{y}_k) + \mathbf{b}$ for a suitable function ψ . Then, if $g(\mathbf{y}_{k-s}) = 0$ for s > L, we can unroll Eq. 4 to the following form:

$$\mathbf{h}_{k} = g(\mathbf{y}_{k}) + \sum_{s=1}^{L} \mathbf{W}^{s} g(\mathbf{y}_{k-s}) \,. \tag{5}$$

Here \mathbf{W}^s indicates *s* applications of \mathbf{W} .

3. Methodology

3.1. Problem Statement

Let us denote L steps of the trajectory of a discrete-time dynamical system as $\mathbf{x}_1, \ldots, \mathbf{x}_L$. We focus on the setting where $\mathbf{x}_k \in \mathcal{X} \subseteq \mathbb{R}^C$. We do not directly observe \mathbf{x}_k , but instead observe $\mathbf{y}_k = h(\mathbf{x}_k)$ for some unknown function h.

We have available a set of training data consisting of multiple sequences of length L + T. The inference task is to forecast the values $\mathbf{y}_{L+1}, \ldots, \mathbf{y}_{L+T}$ given observations of the first L values of a sequence $\mathbf{y}_1, \ldots, \mathbf{y}_L$.

3.2. Strategy: Directly observable systems

We first consider the setting where we directly observe the state, i.e., h is the identity, and $\mathbf{y}_k = h(\mathbf{x}_k) = \mathbf{x}_k$. Assume that the dynamical system can be captured by $\mathbf{x}_{k+1} = F(\mathbf{x}_k)$. Then an appropriate forecasting approach is to learn a finite dimensional approximation to the Koopman operator associated with F, and then propagate it forward in time to construct the forecast. In this section, for notational simplicity, we describe a setting where learning is based on a single observation $\mathbf{y}_0, \ldots, \mathbf{y}_L$ (which, for now, $= \mathbf{x}_0, \ldots, \mathbf{x}_L$). The extension to learning using multiple observed series is straightforward.

Let us introduce measurement functions $g_1, g_2, \ldots g_{n_g} \in \mathcal{H}$, where \mathcal{H} is a Hilbert space containing real-valued functions defined on \mathcal{X} . Denoting $\mathcal{L}(\mathcal{H})$ as the space of bounded linear operators $T : \mathcal{H} \to \mathcal{H}$, Khosravi (2023) formulates the learning of the Koopman operator as a minimization

¹Koopman operators can also be defined on a continuous-time dynamical system, which we disregard here for now.

task with a Tikhonov-regularized empirical loss:

$$\min_{\mathbf{K}\in\mathcal{L}(\mathcal{H})}\sum_{k=1}^{L}\sum_{l=1}^{n_g} \left(\mathbf{x}_{kl} - (\mathbf{K}g_l)(\mathbf{x}_{k-1})\right)^2 + \lambda ||\mathbf{K}||^2.$$
(6)

Although this optimization is over an infinite-dimensional space of linear operators, Khosravi (2023) demonstrates that if the measurement functions satisfy certain conditions, then there is a unique solution \hat{K} , which can be derived by solving a finite-dimensional optimization problem.

A special case occurs when \hat{K} is invariant in the subspace spanned by the measurement functions, $\mathcal{G} =$ span $\{g_1, \ldots, g_{n_g}\}$, i.e. $\hat{K} \in \mathcal{L}(\mathcal{G})$. In this setting, we can follow the approach of the Extended Dynamic Mode Decomposition (EDMD) method (Li et al., 2017), which approximates the Koopman operator by a finite-dimensional linear map U : $\mathcal{G} \to \mathcal{G}$. We can represent the action of U on g_l using a matrix $\mathbf{M} \in \mathbb{R}^{n_g \times n_g}$. Because the dimension of \mathcal{G} is finite, we can identify an \mathbf{M} such that $Ug_l = \sum_{j=1}^{n_g} [M]_{j,l}g_j$. The matrix \mathbf{M} can then be estimated via the following minimization:

$$\min_{\mathbf{M}} ||\mathbf{P}_{G}\mathbf{M} - \mathbf{Y}||_{F}^{2}.$$
 (7)

where $\mathbf{P}_G = [g_l(\mathbf{x}_{k-1})]_{k=1,l=1}^{L,n_g}$ and $\mathbf{Y} = [g_l(\mathbf{x}_k)]_{k=1,l=1}^{L,n_g}$. When applying EDMD, we assume that \mathbf{P}_G is full rank so that a unique minimizer can be identified via the Moore-Penrose pseudoinverse.

This learning task can be made more flexible by allowing for learning of the measurement functions g_l . Li et al. (2017) propose a method that incorporates neural network based learning of the measurement functions and an L_1 regularizer to promote sparsity.

3.3. Unobserved states

In many settings, we do not observe the state of the system \mathbf{x}_k directly. Instead we observe $\mathbf{y}_k = h(\mathbf{x}_k)$ for some unknown observation function $h(\cdot)$. In this setting, \mathbf{y}_k may not provide sufficient information in isolation to recover the state \mathbf{x}_k . We can instead construct a state representation by considering the past L measurements, i.e., we define $\tilde{\mathbf{x}}_k = [\mathbf{y}_{k-L+1}, \dots, \mathbf{y}_k]^{\top}$. With this representation, we can model the dynamics as $\tilde{\mathbf{x}}_k = \tilde{\mathbf{F}}(\tilde{\mathbf{x}}_{k-1})$ and target learning of the Koopman operator \tilde{K} associated with $\tilde{\mathbf{F}}$. Note that this permits us to perform prediction, because we are interested in predicting \mathbf{y}_{k+1} , which can be recovered from $\tilde{\mathbf{x}}_{k+1}$.

Under the same assumptions of invariance of the Koopman operator with respect to \mathcal{G} , we can adopt the same approach as outlined above, learning a matrix **M**. Given the structure of the constructed $\tilde{\mathbf{x}}_k$, we are motivated to impose further structure on the Koopman operator matrix, with the goal of introducing an inductive bias that can facilitate learning and

make it more robust. In particular, we enforce a structure on M that allows us to write: L

$$g(\widetilde{\mathbf{x}}_k) = \mathbf{M}g(\widetilde{\mathbf{x}}_{k-1}) = g(\mathbf{y}_{k-1}) + \sum_{s=1}^{n} \mathbf{W}^s g(\mathbf{y}_{k-s}).$$
(8)

With this structure, we see that **M** is a blockwise diagonal matrix, where each block is a power of a learnable matrix **W**. Moreover, by comparing Eq. 8 with Eq. 5, we see that this structure, which represents the dynamic state using a collection of time-delayed observations (Arbabi & Mezic, 2017), can be implemented as a linear RNN.

3.4. SKOLR

Building on our analysis of Koopman operator approximation and the connection to the linear RNN, we present SKOLR, which integrates a learnable spectral decomposition of the input signal with a multilayer perceptron (MLP) for the measurement functions.

Inspired by multiresolution DMD (Kutz et al., 2016), instead of learning a single linear RNN acting on a high dimensional space, we propose to split the space into multiple subspaces, resulting in learning a structured Koopman operator via a highly parallel linear RNN stack. This structure also improves the parameter efficiency, as shown in Fig. 1

Encoder Let the input sequence be $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_L]$, where $\mathbf{y}_k \in \mathbb{R}^P$, with P being the dimension of the observation. The encoder performs learnable frequency decomposition via reconstruction of the soft-gated frequency spectrum via Fast Fourier Transform (FFT) and Inverse FFT (IFFT):

$$S = FFT(Y),$$

$$S_n = S \cdot Sigmoid(w_n),$$
(9)

$$Y_n = IFFT(S_n).$$

The reconstructed signals $\{\mathbf{Y}_n\}_{n=1}^N$ form N parallel branches, with each branch representing a frequency-based subspace, while $\{\mathbf{w}_n\}_{n=1}^N$ contain learnable parameters for frequency selection.

For each branch n we parameterize the measurement functions using non-linear feed-forward network:

$$\mathbf{z}_{n,k} = \text{FFN}_{\text{enc},n}(\mathbf{y}_k), \text{ for } k = 1, \dots, L$$
 (10)

where FFN : $\mathbb{R}^P \to \mathbb{R}^D$. We utilize multiple layer perceptrons (MLPs) for simplicity, generally with only one layer or two. Other FFNs like wiGLU (Shazeer, 2020) are also applicable. After this encoding is complete, we have constructed the $\mathbf{z}_k = g(\mathbf{y}_k)$ (and hence the $g(\tilde{\mathbf{x}}_k)$) that appear in Eq. 8 for $k = 1, \ldots, L$. By structuring the architecture into multiple branches and incorporating both frequency-domain filtering and time-domain encoding, we enhance flexibility in learning suitable measurement functions for diverse time-series patterns.



Figure 1: Architecture of SKOLR (Structured Koopman Operator Linear RNN) The input time series goes through an encoder with learnable frequency decomposition and a MLP that models the measurement functions. With the branch decomposition, the highly parallel linear RNN chains jointly attend to different dynamical patterns from different representation subspaces. Finally, a decoder reconstructs predictions by parameterizing the inverse measurement functions. This structured approach maintains computational efficiency while naturally aligning with Koopman principles.

RNN Stack Given the collection for each branch n and time step k: $\mathbf{Z}_n = [\mathbf{z}_{1,n}, \dots, \mathbf{z}_{L,n}] \in \mathbb{R}^{D \times L}$, we take it as input to a linear RNN and introduce learnable branch-specific weight matrices \mathbf{W}_n for each branch.

$$\mathbf{h}_{k+1,n} = \mathbf{W}_n \mathbf{h}_{k,n} + \mathbf{z}_{k,n} \tag{11}$$

Each branch weight matrix \mathbf{W}_n defines a matrix \mathbf{M}_n , as discussed above, which specifies a finite-dimensional approximation to a Koopman operator for $\tilde{\mathbf{x}}$ for the learned measurement functions on that branch.

Together, the branch matrices \mathbf{M}_n form a structured finitedimensional Koopman operator approximation $\widehat{\mathbf{K}}$ with block diagonal structure:

$$\widehat{\mathbf{K}} = \begin{bmatrix} \mathbf{M}_1 & 0 & \cdots & 0 \\ 0 & \mathbf{M}_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{M}_N \end{bmatrix}$$
(12)

By imposing this structure and using a stack of linear RNNs, we can learn local approximations to the evolution dynamics of different observables. For each branch n:

$$\mathbf{H}_{n} = [\mathbf{z}_{1,n}, \mathbf{z}_{2,n} + \mathbf{M}_{n}\mathbf{z}_{1,n}, \dots, \mathbf{z}_{L,n} + \sum_{s=0}^{L-1}\mathbf{M}_{n}^{s}\mathbf{z}_{s,n}]$$
(13)

For prediction of length T, we recursively apply the operator to predict the Koopman space for future steps per branch:

$$\mathbf{H}_{[L+1:L+T],n} = [\mathbf{M}_n \mathbf{h}_{L,n}, \dots, \mathbf{M}_n^T \mathbf{h}_{L,n}]$$
(14)

Decoder For reconstruction, we use mirrored feedforward networks to parameterize the inverse measurement functions g^{-1} . The decoder processes the hidden states as:

$$\hat{\mathbf{y}}_{k,n} = \text{FFN}_{\text{dec},n}(\mathbf{h}_{k,n}) \tag{15}$$

where $\text{FFN}_{\text{dec}} : \mathbb{R}^D \to \mathbb{R}^P$.

The decoder combines predictions from all branches to generate the final prediction $\hat{\mathbf{y}}_{[L+1,L+T]}$. The model is trained end-to-end using the loss function:

$$\mathcal{L} = \|\hat{\mathbf{y}}_{[L+1:L+T]} - \mathbf{y}_{[L+1:L+T]}\|_2^2.$$
(16)

The structured approach, induced by both the linear RNN and the branch decomposition, enables efficient parallel processing and reduces the parameter count. Since all architectural components are very simple (basic sigmoid frequency gating, one- or two-layer MLPs for encoding/decoding, linear RNN), the architecture is very fast to train and has low memory cost, as we illustrate in the experiments section.

Models	T	SKO	OLR	Ko	opa	iTrans	former	Patch	nTST	Time	esNet	Dli	near	MI	CN	KI	NF	Autof	ormer
		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ECL	48 96 144 192	0.137 0.132 0.143 0.149	0.229 0.225 0.236 0.244	0.130 0.136 0.149 0.156	0.234 0.236 0.247 0.254	$ \begin{array}{c c} 0.133 \\ \hline 0.134 \\ 0.145 \\ 0.154 \end{array} $	$\begin{array}{r} \textbf{0.225} \\ \underline{0.230} \\ \underline{0.240} \\ 0.249 \end{array}$	0.147 0.143 0.145 0.145 0.147	0.246 0.241 0.241 0.240	0.149 0.170 0.183 0.189	0.254 0.275 0.287 0.291	0.158 0.153 0.152 0.153	0.241 0.245 0.245 0.246	0.156 0.165 0.163 0.171	0.271 0.277 0.274 0.284	0.175 0.198 0.204 0.245	0.265 0.284 0.297 0.321	0.164 0.182 0.210 0.221	0.272 0.289 0.315 0.324
Traffic	48 96 144 192	0.400 0.368 0.375 0.377	0.258 0.248 0.255 0.256	$\begin{array}{c c} 0.415 \\ 0.401 \\ \underline{0.397} \\ 0.403 \end{array}$	0.274 0.275 0.276 0.284	0.369 0.388 0.375 0.373	$\begin{array}{r} \textbf{0.256} \\ \underline{0.270} \\ \underline{0.267} \\ \underline{0.267} \\ \underline{0.267} \end{array}$	0.426 0.413 0.405 0.404	0.286 0.283 0.278 0.277	0.567 0.611 0.603 0.604	0.306 0.337 0.322 0.321	0.488 0.485 0.452 0.438	0.352 0.336 0.317 0.309	0.496 0.511 0.498 0.494	0.301 0.312 0.309 0.312	0.621 0.645 0.683 0.699	0.382 0.376 0.402 0.405	0.640 0.668 0.681 0.692	0.361 0.367 0.379 0.385
Weather	48 96 144 192	0.131 0.154 0.172 0.193	$ \begin{array}{r} 0.170 \\ 0.202 \\ 0.220 \\ 0.241 \end{array} $	0.126 0.154 0.172 0.193	0.168 0.205 0.225 0.241	$ \begin{array}{r} 0.136 \\ 0.169 \\ \underline{0.192} \\ 0.204 \end{array} $	0.174 0.216 0.242 0.251	$ \begin{array}{c c} 0.140 \\ \underline{0.160} \\ 0.174 \\ 0.195 \end{array} $	$\begin{array}{r} 0.179 \\ 0.206 \\ \underline{0.221} \\ 0.243 \end{array}$	0.138 0.180 0.190 0.212	0.191 0.231 0.244 0.265	0.156 0.186 0.199 0.217	0.198 0.229 0.244 0.261	0.157 0.187 0.197 0.214	0.217 0.250 0.257 0.270	0.201 0.295 0.394 0.462	0.288 0.308 0.401 0.437	0.185 0.230 0.268 0.325	0.240 0.279 0.308 0.347
ETTm1	48 96 144 192	0.280 0.289 0.319 0.328	0.330 0.340 <u>0.361</u> 0.373	$ \begin{array}{c} \underline{0.283} \\ \underline{0.294} \\ \underline{0.322} \\ \underline{0.337} \end{array} $	$ \begin{array}{r} 0.333 \\ \underline{0.345} \\ 0.366 \\ 0.378 \end{array} $	0.313 0.302 0.331 0.343	0.356 0.353 0.374 0.381	0.286 0.299 0.325 0.343	0.336 0.346 0.363 0.375	0.308 0.329 0.358 0.462	0.354 0.370 0.387 0.441	0.322 0.309 0.327 0.337	0.355 0.346 0.359 <u>0.365</u>	0.294 0.306 0.342 0.386	0.353 0.364 0.390 0.415	1.026 0.957 0.921 0.896	0.792 0.782 0.760 0.731	0.592 0.493 0.735 0.592	0.419 0.469 0.569 0.506
ETTm2	48 96 144 192	0.134 0.171 0.209 0.241	0.228 0.255 0.283 0.304	0.134 0.171 0.206 0.226	0.226 0.254 0.280 0.298	0.139 0.177 0.216 0.237	0.234 0.268 0.296 0.310	0.135 0.171 0.205 0.221	$ \begin{array}{r} 0.231 \\ \underline{0.255} \\ 0.282 \\ \underline{0.294} \end{array} $	0.142 0.187 0.216 0.243	0.234 0.269 0.291 0.313	0.144 0.172 0.200 0.219	0.240 0.256 0.276 0.290	0.131 0.197 0.210 0.248	0.238 0.295 0.297 0.328	0.621 1.535 1.337 1.355	0.623 1.012 0.876 0.908	0.191 0.241 0.300 0.324	0.280 0.311 0.352 0.370
ETTh1	48 96 144 192	0.333 0.371 0.405 0.422	$\begin{array}{r} \underline{0.373}\\ 0.398\\ 0.417\\ 0.432 \end{array}$	0.336 0.371 0.405 0.416	0.377 0.405 0.418 <u>0.429</u>	0.342 0.393 0.425 0.456	0.380 0.412 0.430 0.454	$\begin{array}{c c} 0.337 \\ \hline 0.372 \\ \hline 0.394 \\ \hline 0.416 \end{array}$	0.375 0.393 <u>0.412</u> 0.439	0.365 0.411 0.442 0.469	0.399 0.430 0.447 0.470	0.343 0.379 0.393 0.407	0.371 0.393 0.403 0.416	0.375 0.406 0.437 0.518	0.406 0.429 0.448 0.496	0.876 0.975 0.801 0.941	0.709 0.744 0.662 0.744	0.442 0.634 0.522 0.525	0.438 0.523 0.491 0.501
ETTh2	48 96 144 192	$\begin{array}{c c} 0.238 \\ 0.299 \\ \underline{0.335} \\ 0.365 \end{array}$	0.306 0.352 0.377 <u>0.397</u>	0.226 0.297 0.333 0.356	0.300 0.349 0.381 0.393	0.243 0.306 0.347 0.375	0.313 0.358 0.385 0.403	0.223 0.300 0.346 0.383	0.297 0.353 0.390 0.406	0.241 0.325 0.374 0.394	0.319 0.376 0.408 0.434	0.226 0.294 0.354 0.385	$\begin{array}{r} 0.305 \\ \underline{0.351} \\ 0.397 \\ 0.418 \end{array}$	0.260 0.343 0.374 0.455	0.336 0.393 0.411 0.464	0.385 0.433 0.441 0.528	$\begin{array}{c} 0.376 \\ 0.446 \\ 0.456 \\ 0.503 \end{array}$	0.355 0.427 0.457 0.503	0.380 0.432 0.461 0.491
ILI	24 36 48 60	1.556 1.462 1.537 2.187	0.760 0.728 0.798 0.995	$ \begin{array}{r} \frac{1.621}{1.803} \\ \hline{1.768} \\ 1.743 \end{array} $	0.800 0.855 0.903 0.891	$ \begin{array}{r} 1.763 \\ 2.067 \\ \underline{1.667} \\ 2.011 \end{array} $	$0.843 \\ 0.919 \\ \underline{0.879} \\ 1.000$	2.063 2.178 1.916 <u>1.981</u>	0.881 0.943 0.896 <u>0.917</u>	2.464 2.388 2.370 2.193	1.039 1.007 1.040 1.003	2.624 2.693 2.852 2.554	1.118 1.156 1.229 1.144	4.380 3.314 2.457 2.379	1.558 1.313 1.085 1.040	3.722 3.941 3.287 2.974	1.432 1.448 1.377 1.301	2.831 2.801 2.322 2.470	1.085 1.088 1.006 1.061
Rank 1 st	#	17	15	10	7	3	2	3	3	0	0	5	7	1	0	0	0	0	0

Table 1: Prediction results on benchmark datasets, L = 2T and $T \in \{48, 96, 144, 192\}$ (ILI: $T \in \{24, 36, 48, 60\}$). Best results and second best results are highlighted in red and blue respectively.

4. Experiments

4.1. Benchmarking SKOLR

4.1.1. DATASETS

We evaluate SKOLR on widely-used public benchmark datasets. For long-term forecasting, we use Weather, Traffic, Electricity, ILI and four ETT datasets (ETTh1, ETTh2, ETTm1, ETTm2). We assess short-term performance on M4 dataset (Makridakis et al., 2020), which includes six subsets of periodically recorded univariate marketing data. For more information about the datasets see Appendix A.1.

4.1.2. BASELINES AND EXPERIMENTAL SETTINGS

We compare against state-of-the-art deep forecasting models. The comparison includes transformer-based models: Autoformer (Wu et al., 2021), PatchTST (Nie et al., 2023), iTransformer (Liu et al., 2024), TCN-based models: Times-Net (Wu et al., 2023), MICN (Wang et al., 2023a), linear model: DLinear (Zeng et al., 2023), and Koopman-based models: KNF (Wang et al., 2023b), Koopa (Liu et al., 2023). We select these representative baselines for their established performance and public implementations.

Following Koopa (Liu et al., 2023), we set the lookback window length L = 2T for prediction horizon $T \in$ {48, 96, 144, 192} for all datasets, except ILI, for which we use $T \in$ {24, 36, 48, 60}. This setting leverages more historical data for longer forecasting horizons. We report baseline results from Liu et al. (2023) except for iTransformer; we reproduce iTransformer results with L = 2Tusing the officially released code. Performance is measured using Mean Squared Error (MSE) and Mean Absolute Error (MAE). Appendix A.2 provides implementation details.

4.1.3. RESULTS AND ANALYSIS

Table 1 reports the experimental results for eight benchmarks. The performance is measured by MSE and MAE; the best and second-best results for each case (dataset, horizon, and metric) are highlighted in bold and underlined, respectively. The results are the average of 3 trials.

We rank the algorithms in Table 1 based on their MSE



Figure 2: Boxplot for ranks of the algorithms (based on their MSE) across seven datasets and four prediction horizons. The medians and means of the ranks are shown by the vertical lines and the black triangles respectively; whiskers extend to the minimum and maximum ranks.

Table 2: Model evaluation results (MSE/MAE) on nonlinear dynamical systems (NLDS)

	SKO	OLR	KooPA				
Dataset	MSE	MAE	MSE	MAE			
Pendulum Duffing Lotka-Volterra Lorenz '63	0.0001 0.0047 0.0018 0.9740	0.0083 0.0518 0.0354 0.7941	0.0039 0.0365 0.0178 1.0937	0.0470 0.1479 0.1050 0.8325			

and order them based on their average rank across eight datasets and four prediction horizons. Figure 2 shows the relative ranks. We observe that SKOLR achieves SOTA performance, with the best average rank across all settings.

The model shows strength in capturing complex patterns in the Weather dataset, matching Koopa's performance while surpassing other transformers, indicating effective handling of meteorological dynamics. For the ILI dataset, which features highly nonlinear epidemic patterns, SKOLR outperforms the baseline methods, with significant error reduction for the shorter horizons.

While SKOLR demonstrates strong performance in longterm forecasting, we also evaluate its effectiveness on shortterm predictions with M4 dataset. Results in Appendix B.1 show consistent improvements over both transformer-based forecasting methods and Koopman-based alternatives across different time scales.

4.2. State Prediction for Non-Linear Systems

Koopman operator-based approaches have gained attention for their ability to perform system identification in a fully data-driven manner. To evaluate SKOLR's performance in this context, we conducted a series of experiments on nonlinear dynamical systems (NLDS) (details in Appendix E).



Figure 3: Analysis of SKOLR's branch-wise behavior: (a) frequency decomposition and (b) prediction performance. We observe that different branches focus on different frequency components.

Table 2 demonstrates SKOLR's effectiveness across different dynamical systems. For periodic systems like Pendulum, SKOLR achieves substantial improvements, indicating superior capture of oscillatory patterns. In chaotic systems like Lorenz '63, SKOLR shows better stability with 10.9% reduction on MSE, suggesting robust handling of sensitive dependence on initial conditions. The model demonstrates particularly strong performance on mixed dynamics: Lotka-Volterra and Duffing oscillator. These results validate that SKOLR's structured operator design effectively captures both periodic motions and complex nonlinear dynamics.

Fig. 3 demonstrates SKOLR's multi-scale decomposition strategy. The FFT analysis reveals how different branches place more emphasis on some frequency bands. This natural frequency partitioning emerges from our structured Koopman design, enabling each branch to focus on specific temporal scales. The prediction visualization illustrates the complementary nature of these branches, where their combined forecasts reconstruct complex dynamics through principled superposition of simpler, frequency-specific predictions. More analysis can be found in Appendix D.3.

4.3. Analysis and Ablation Study

4.3.1. ANALYSIS: STRUCTURED KOOPMAN OPERATOR

We analyze the impact of branch configurations through two controlled experiments: (1) Fixed parameter count scenario, where total parameters remain constant (~1.6M) across configurations while varying the learnable frequency decomposition w, with N branches and lookback window L; (2) Fixed dimension scenario: We maintain constant Koopman operator approximation dimension dim(\hat{K}) = 512 while varying branch number N. As N increases, each branch's dimension D decreases proportionally (D = 512/N), leading to reduced parameter count.

We conduct experiments on the ETTm1 dataset. As we can see in Table 3 and Fig. 4, maintaining similar parameter



Figure 4: MSE comparison on ETTm1 dataset across different branch configurations and prediction horizons. Bars show MSE values for each configuration. All configurations maintain similar parameter counts (\sim 1.6M). Increasing branch number improves performance.

Table 3: Performance comparison with similar parameter counts on ETTm1

Config.	Ν	MSE for Different T									
(D, N)	48	96	144	192							
(512, 1) (256, 4)	0.284	0.292	0.326 0.317	0.330							
(128, 16)	0.281	0.291	0.323	0.329							

counts (~1.6M) and increasing branch numbers from N = 1 to N = 16 improves performance for most horizons.

More significantly, when keeping dim $(\hat{K}) = 512$ (Table 4), models with more branches maintain strong performance despite substantial parameter reduction. Notably, the configuration with D = 32, N = 16 achieves comparable performance to the 1-branch model while using only 0.25M parameters (85% reduction). This demonstrates that structured decomposition through multiple branches enables significantly more efficient parameter utilization while maintaining or improving forecasting accuracy.

4.3.2. Ablation: Impact of Frequency Decomposition

The improved performance with multiple branches motivates further analysis of our learnable frequency decomposition strategy. In Equation 9, the learnable matrix w enables adaptive frequency allocation across branches, in contrast to uniform decomposition (w = 1). Table 5 demonstrates that this learnable approach consistently outperforms uniform allocation across prediction horizons.

This adaptive capability is particularly beneficial for datasets with complex temporal patterns (Weather, ECL), where different frequency bands may carry varying importance at different time scales. The learned masks show distinct patterns across datasets, suggesting that the model successfully adapts its frequency decomposition strategy based on the un-

Table 4: Performance comparison with $\dim(\widehat{K}) = 512$ and varying branch numbers N on ETTm1. Parameter counts shown for horizon T = 192.

Config.	Params	Ν	MSE for Different T										
(D, N)	(M)	48	96	144	192								
(512, 1)	1.71	0.284	0.292	0.326	0.330								
(256, 2)	0.92	0.280	0.294	0.318	0.334								
(128, 4)	0.53	0.280	0.293	0.318	0.335								
(64, 8)	0.34	0.283	0.297	0.316	0.329								
(32, 16)	0.25	0.282	0.292	0.313	0.328								

Table 5: Ablation Study on frequency decomposition

Dataset	Т	SKOLF	R (learn)	SKOLF	R(w = 1)
	-	MSE	MAE	MSE	MAE
	48	0.137	0.229	0.150	0.239
ECI	96	0.132	0.225	0.134	0.227
ECL	144	0.143	0.236	0.144	0.237
	192	0.149	0.244	0.150	0.244
	48	0.131	0.170	0.134	0.173
Waathar	96	0.154	0.202	0.158	0.203
weather	144	0.172	0.220	0.175	0.221
	192	0.193	0.241	0.194	0.242
	48	0.333	0.373	0.329	0.371
ETTL 1	96	0.371	0.398	0.375	0.400
E1111	144	0.405	0.417	0.407	0.419
	192	0.422	0.432	0.429	0.433

derlying data characteristics. Notably, on the ETTh1 dataset, learnable decomposition occasionally underperforms uniform masking, particularly at shorter horizons (T = 48). This suggests potential overfitting on smaller datasets.

4.4. Model Efficiency

To demonstrate the computational efficiency of SKOLR, we analyze the model complexity in terms of parameter count, GPU Memory and Running Time. We compare these values with several baseline models on the ETTm1 and weather dataset with sequence length 96 and prediction length 48.

Fig. 5 demonstrates SKOLR's computational advantages. On ETTm1, SKOLR achieves the best MSE while using only 3.31 MiB GPU memory. The training speed is also notably faster than other methods. On Weather dataset, SKOLR maintains competitive accuracy, while using significantly less memory and training 4x faster compared to the best. This exceptional efficiency-performance trade-off stems from our structured linear operations in Koopman space, avoiding the quadratic complexity of selfattention while maintaining modeling capacity through parallel branch architecture. The computational efficiency for all datasets can be found in Appendix C.



Figure 5: Model comparison on error and training epoch time on P100 GPU. Memory consumption is proportional to circle radius. On ETTm1 (left) we observe that SKOLR is both the fastest and most accurate method while requiring the smallest memory footprint. On Weather (right) SKOLR is the second best method with much lower training memory and time consumption than the best. For an equivalent budget one could train an ensemble of our approach and obtain better results.

5. Related Work

5.1. Koopman Operator-based Time-series Forecasting

Koopman theory (Koopman, 1931) has been applied for modeling and analyzing complex dynamical systems for decades (Mezić, 2005; Brunton et al., 2022). The major advantage of the Koopman operator is that it can represent the dynamical system in the form of a linear operator acting on measurement functions (observables). However, learning the operator is challenging because it has infinite dimension. Researchers strive to develop effective strategies for performing finite dimensional approximations; key to this is the selection of good measurement functions. To address this, recent work has explored neural networks for learning the mapping and the approximate operator simultaneously (Li et al., 2017; Lusch et al., 2018; Takeishi et al., 2017; Yeung et al., 2019).

Three recent works address time-series forecasting using Koopman operators. K-Forecast (Lange et al., 2021) uses Koopman theory to handle the nonlinearity in temporal signals and proposes a data-dependent basis for long-term timeseries forecasting. By leveraging predefined measurement functions, KNF (Wang et al., 2023b) learns the Koopman operator and attention map to cope with time-series forecasting with changing temporal distributions. Koopa (Liu et al., 2023) introduces modular Koopman predictors that separately address time-variant and time-invariant components via a hierarchical architecture, using learnable operators for the latter and eDMD (Williams et al., 2015) for the former. These prior works rely on hierarchical architectures or complex spectral decompositions to approximate Koopman operators. Our work takes a different approach, drawing a connection with linear RNNs, paving the way to a very efficient and simple forecasting architecture. Our results demonstrate that this strategy leads to improved accuracy with reduced computational overhead and memory.

Although Orvieto et al. (2023) provided insights into the potential connections between the Koopman operator and a wide MLP + linear RNN for representing dynamical systems, this was not the primary focus of their work, and they did not provide equations demonstrating the connection or conduct empirical verification. In this work, building on similar insights, we establish an explicit connection by deriving equations that demonstrate a direct analogy between a structured approximation of a Koopman operator and an architecture consisting of an MLP encoder combined with a linear RNN.

5.2. Deep Learning for Time-Series Forecasting

Time-series forecasting has evolved from statistical models (Makridakis & Hibon, 1997; Hyndman et al., 2008) to deep learning approaches. Previous methods used RNNs (Salinas et al., 2020; Smyl, 2020; Mienye et al., 2024) and CNNs (Bai et al., 2018; Luo et al., 2024) for their ability to capture temporal dependencies. MLP-based architectures (Oreshkin et al., 2020; Challu et al., 2023; Vijay et al., 2023; Wang et al., 2024a) also demonstrated promising performance for forecasting. Recently, transformer architectures (Nie et al., 2023; Zhang et al., 2024; Hounie et al., 2024; Ilbert et al., 2024) introduced powerful attention mechanisms, with innovations in basis functions (Ni et al., 2024) and channel-wise processing (Liu et al., 2024). To address their quadratic complexity, sparse attention variants (Lin et al., 2024) were proposed, but these often struggle with capturing long-range dependencies due to information loss from pruned attention scores. Foundation models (Das et al., 2024; Darlow et al., 2024) and unified approaches (Woo et al., 2024) have recently emerged. These attempt to mitigate the limitations through pre-training and multi-task learning, but this comes at the cost of dramatically increased architectural complexity and computational overhead. To address the complexity challenges in timeseries forecasting, recent state space models (Gu & Dao, 2023) achieve linear complexity, while physics-informed approaches (Verma et al., 2024) enhance interpretability. However, these methods often require complex architectures or domain expertise. Our approach offers a balanced solution with a principled foundation based on Koopman theory, achieving excellent prediction performance with very low computation and memory requirements.

6. Conclusion

This work establishes a connection between Koopman operator approximation and linear RNNs, showing that timedelayed state representations yield an equivalence between structured Koopman operators and linear RNN updates. Based on this, we introduce SKOLR, which integrates learnable spectral decomposition with a parallelized linear RNN stack giving rise to the structured Koopman operator. By aligning deep learning with Koopman theory, this approach provides a principled and computationally efficient solution for nonlinear time-series modeling. Empirical evaluations on forecasting benchmarks and dynamical systems demonstrate that SKOLR achieves strong predictive performance while maintaining the efficiency of linear RNNs. Future work includes extending this framework to broader dynamical systems and exploring alternative spectral representations for enhanced expressivity.

Acknowledgement

This research was funded by the Natural Sciences and Engineering Research Council of Canada (NSERC), [reference number 260250]. Cette recherche a été financée par le Conseil de recherches en sciences naturelles et en génie du Canada (CRSNG), [numéro de référence 260250]. Ce projet de recherche #324302 est rendu possible grâce au financement du Fonds de recherche du Québec.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Arbabi, H. and Mezic, I. Ergodic theory, dynamic mode decomposition, and computation of spectral properties of the koopman operator. *SIAM J. Appl. Dyn. Syst.*, 16(4): 2096–2126, 2017.
- Bai, S., Kolter, J. Z., and Koltun, V. An empirical evaluation of generic convolutional and recurrent networks for

sequence modeling. *arXiv e-prints: arXiv 1803.01271*, 2018.

- Bevanda, P., Sosnowski, S., and Hirche, S. Koopman operator dynamical models: Learning, analysis and control. *Annu. Rev. Control.*, 52, January 2021.
- Brunton, S. L., Budišić, M., Kaiser, E., and Kutz, J. N. Modern koopman theory for dynamical systems. *SIAM Rev.*, 64(2), May 2022.
- Challu, C., Olivares, K. G., Oreshkin, B. N., Ramirez, F. G., Canseco, M. M., and Dubrawski, A. Nhits: Neural hierarchical interpolation for time series forecasting. In *Proc. AAAI Conf. Artif. Intell.*, 2023.
- Darlow, L., Deng, Q., Hassan, A., Asenov, M., Singh, R., Joosen, A., Barker, A., and Storkey, A. Dam: Towards a foundation model for time series forecasting. In *Proc. Int. Conf. Learn. Represent.*, 2024.
- Das, A., Kong, W., Sen, R., and Zhou, Y. A decoder-only foundation model for time-series forecasting. In *Proc. Int. Conf. Mach. Learn.*, 2024.
- Gu, A. and Dao, T. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- Hounie, I., Porras-Valenzuela, J., and Ribeiro, A. Transformers with loss shaping constraints for long-term time series forecasting. In *Proc. Int. Conf. Mach. Learn.*, 2024.
- Hyndman, R., Koehler, A. B., Ord, J. K., and Snyder, R. D. Forecasting with exponential smoothing: the state space approach. Springer Science & Business Media, 2008.
- Ilbert, R., Odonnat, A., Feofanov, V., Virmaux, A., Paolo, G., Palpanas, T., and Redko, I. Samformer: Unlocking the potential of transformers in time series forecasting with sharpness-aware minimization and channel-wise attention. In *Proc. Int. Conf. Mach. Learn.*, 2024.
- Khosravi, M. Representer theorem for learning Koopman operators. *IEEE Trans. Autom. Control*, 68(5):2995–3010, 2023.
- Kim, T., Kim, J., Tae, Y., Park, C., Choi, J.-H., and Choo, J. Reversible instance normalization for accurate timeseries forecasting against distribution shift. In *Proc. Int. Conf. Learn. Represent.*, 2022.
- Koopman, B. Hamiltonian systems and transformation in hilbert space. *Proc. Natl. Acad. Sci.*, 17(5):315–318, 1931.
- Kutz, J. N., Fu, X., and Brunton, S. L. Multiresolution dynamic mode decomposition. *SIAM J. Appl. Dyn. Syst.*, 15(2):713–735, 2016.

- Lange, H., Brunton, S. L., and Kutz, J. N. From fourier to koopman: Spectral methods for long-term time series prediction. J. Mach. Learn. Res., 22(1):1881–1918, 2021.
- Li, Q., Dietrich, F., Bollt, E. M., and Kevrekidis, I. G. Extended dynamic mode decomposition with dictionary learning: A data-driven adaptive spectral decomposition of the koopman operator. *Chaos*, 27(10), 2017.
- Lin, S., Weiwei, L., Wentai, W., Haojun, C., and Junjie, Y. Sparsetsf: Modeling long-term time series forecasting with 1k parameters. *Proc. Int. Conf. Mach. Learn.*, 2024.
- Liu, Y., Li, C., Wang, J., and Long, M. Koopa: Learning nonstationary time series dynamics with koopman predictors. In Adv. Neural Inf. Process. Syst., 2023.
- Liu, Y., Hu, T., Zhang, H., Wu, H., Wang, S., Ma, L., and Long, M. itransformer: Inverted transformers are effective for time series forecasting. In *Proc. Int. Conf. Learn. Represent.*, 2024.
- Loshchilov, I. Decoupled weight decay regularization. *arXiv* preprint arXiv:1711.05101, 2017.
- Luo, H., Wang, S., Zhang, T., Wang, J., Liu, W., and Lin, W. Modernten: A modern pure convolution structure for general time series analysis. In *Proc. Int. Conf. Learn. Represent.*, 2024.
- Lusch, B., Kutz, J. N., and Brunton, S. L. Deep learning for universal linear embeddings of nonlinear dynamics. *Nat. Commun.*, 9(1):4950, 2018.
- Makridakis, S. and Hibon, M. ARMA models and the Box– Jenkins methodology. J. Forecast., 16(3):147–163, 1997.
- Makridakis, S., Spiliotis, E., and Assimakopoulos, V. The m4 competition: 100,000 time series and 61 forecasting methods. *Int. J. Forecast.*, 36(1):54–74, 2020.
- Mezić, I. Spectral properties of dynamical systems, model reduction and decompositions. *Nonlinear Dyn.*, 41:309– 325, 2005.
- Mienye, I. D., Swart, T. G., and Obaido, G. Recurrent neural networks: A comprehensive review of architectures, variants, and applications. *Information*, 15(9):517, 2024.
- Ni, Z., Yu, H., Liu, S., Li, J., and Lin, W. Basisformer: Attention-based time series forecasting with learnable and interpretable basis. *Adv. Neural Inf. Process. Syst.*, 36, 2024.
- Nie, Y., Nguyen, N. H., Sinthong, P., and Kalagnanam, J. A time series is worth 64 words: Long-term forecasting with transformers. In *Proc. Int. Conf. Learn. Represent.*, 2023.

- Oreshkin, B. N., Carpov, D., Chapados, N., and Bengio, Y. N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. In *Proc. Int. Conf. Learn. Represent.*, 2020.
- Orvieto, A., Smith, S. L., Gu, A., Fernando, A., Gulcehre, C., Pascanu, R., and De, S. Resurrecting recurrent neural networks for long sequences. In *Proc. Int. Conf. Mach. Learn.*, 2023.
- Rowley, C. W., Mezić, I., Bagheri, S., Schlatter, P., and Henningson, D. S. Spectral analysis of nonlinear flows. *J. Fluid Mech.*, 641:115–127, 2009.
- Salinas, D., Flunkert, V., Gasthaus, J., and Januschowski, T. DeepAR: Probabilistic forecasting with autoregressive recurrent networks. *Int. J. Forecast.*, 36(3):1181–1191, 2020.
- Schmid, P. Dynamic mode decomposition of numerical and experimental data. J. Fluid Mech., 656:5–28, 2010.
- Shazeer, N. Glu variants improve transformer. *arXiv* preprint arXiv:2002.05202, 2020.
- Smyl, S. A hybrid method of exponential smoothing and recurrent neural networks for time series forecasting. *Int. J. Forecast.*, 36(1):75–85, 2020.
- Stolzenburg, F., Litz, S., Michael, O., and Obst, O. The power of linear recurrent neural networks. *arXiv preprint arXiv:1802.03308*, 2018.
- Takeishi, N., Kawahara, Y., and Yairi, T. Learning koopman invariant subspaces for dynamic mode decomposition. In *Adv. Neural Inf. Process. Syst.*, 2017.
- Tay, Y., Dehghani, M., Abnar, S., Shen, Y., Bahri, D., Pham, P., Rao, J., Yang, L., Ruder, S., and Metzler, D. Long range arena : A benchmark for efficient transformers. In *Proc. Int. Conf. Learning Representations (ICLR)*, 2021.
- Trindade, A. Electricity load diagrams 20112014. UCI Machine Learning Repository, 2015. DOI: https://doi.org/10.24432/C58C86.
- Verma, Y., Markus, H., and Vikas, G. Climode: Climate and weather forecasting with physics-informed neural odes. *Proc. Int. Conf. Learn. Represent.*, 2024.
- Vijay, E., Jati, A., Nguyen, N., Sinthong, G., and Kalagnanam, J. TSMixer: Lightweight MLP-mixer model for multivariate time series forecasting. In *Proc.* ACM SIGKDD Int. Conf. Knowl. Discov. Data Min., 2023.
- Wang, H., Peng, J., Huang, F., Wang, J., Chen, J., and Xiao, Y. MICN: Multi-scale local and global context modeling for long-term series forecasting. In *Proc. Int. Conf. Learn. Represent.*, 2023a.

- Wang, R., Dong, Y., Arik, S. Ö., and Yu, R. Koopman neural forecaster for time series with temporal distribution shifts. In *Proc. Int. Conf. Learn. Represent.*, 2023b.
- Wang, S., Wu, H., Shi, X., Hu, T., Luo, H., Ma, L., Zhang, J. Y., and Zhou, J. Timemixer: Decomposable multiscale mixing for time series forecasting. *Proc. Int. Conf. Learn. Represent.*, 2024a.
- Wang, Z., Kong, F., Feng, S., Wang, M., Yang, X., Zhao, H., Wang, D., and Zhang, Y. Is mamba effective for time series forecasting? *Neurocomputing*, pp. 129178, 2024b.
- Williams, M., Kevrekidis, I., and Rowley, C. A data–driven approximation of the koopman operator: Extending dynamic mode decomposition. *J. Nonlinear Sci.*, 25:1307– 1346, 2015.
- Woo, G., Liu, C., Kumar, A., Xiong, C., Savarese, S., and Sahoo, D. Unified training of universal time series forecasting transformers. *Proc. Int. Conf. Mach. Learn.*, 2024.
- Wu, H., Xu, J., Wang, J., and Long, M. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In *Adv. Neural Inf. Process. Syst.*, 2021.
- Wu, H., Hu, T., Liu, Y., Zhou, H., Wang, J., and Long, M. Timesnet: Temporal 2d-variation modeling for general time series analysis. In *Proc. Int. Conf. Learn. Represent.*, 2023.
- Yeung, E., Kundu, S., and Hodas, N. Learning deep neural network representations for koopman operators of nonlinear dynamical systems. In *Proc. Amer. Control Conf.* IEEE, 2019.
- Zeng, A., Chen, M., Zhang, L., and Xu, Q. Are transformers effective for time series forecasting? In *Proc. AAAI Conf. Artif. Intell.*, 2023.
- Zhang, Y., Ma, L., Pal, S., Zhang, Y., and Coates, M. Multiresolution time-series transformer for long-term forecasting. In *Int. Conf. Artif. Intell. Stat.*, 2024.
- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proc. AAAI Conf. Artif. Intell.*, 2021.
- Zhou, T., Ma, Z., Wen, Q., Sun, L., Yao, T., Yin, W., Jin, R., et al. Film: Frequency improved legendre memory model for long-term time series forecasting. In *Adv. Neural Inf. Process. Syst.*, 2022.

A. Experimental Details

A.1. Dataset

We evaluate the performance of our proposed SKOLR on eight widely-used public benchmark datasets, including Weather, Traffic, Electricity, ILI and four ETT datasets (ETTh1, ETTh2, ETTm1, ETTm2). *Weather* is a collection of 2020 weather data from 21 meteorological indicators, including air temperature and humidity, provided by the Max-Planck Institute for Biogeochemistry. ² *Traffic* is a dataset provided by Caltrans Performance Measurement System (PeMS), collecting hourly data of the road occupancy rates measured by different sensors on San Francisco Bay area freeways from California Department of Transportation. ³ *Electricity* contains hourly time series of the electricity consumption of 321 customers from 2012 to 2014 (Trindade, 2015; Wu et al., 2021). ⁴ *ILI* dataset⁵ contains the weekly time series of ratio of patients seen with ILI and the total number of the patients in the United States between 2002 and 2021. *ETT* datasets are a series of measurements, including load and oil temperature, from electricity transformers between 2016 and 2018, provided by Zhou et al. (2021). Following the standard pipelines, the dataset is split into training, validation, and test sets with the ratio of 6:2:2 for four ETT datasets and 7:1:2 for the remaining datasets.

The M4 dataset (Makridakis et al., 2020) consists of 100,000 real-world time series across six frequencies: yearly, quarterly, monthly, weekly, daily, and hourly. It includes data from diverse domains such as finance, economics, demographics, and industry, making it a comprehensive benchmark for evaluating forecasting models. Detailed statistics of the datasets are summarized in Table. 6.

Tasks	Dataset	Dim	Prediction Length	Dataset Size	Information (Frequency)
	ETTm1, ETTm2	7	{48, 96, 144, 192}	(34465, 11521, 11521)	Electricity (15 mins)
	ETTh1, ETTh2	7	{48, 96, 144, 192}	(8545, 2881, 2881)	Electricity (15 mins)
I and tam	Electricity	321	{48, 96, 144, 192}	(18317, 2633, 5261)	Electricity (Hourly)
Long-term	Traffic	862	{48, 96, 144, 192}	(12185, 1757, 3509)	Transportation (Hourly)
	Weather	21	{48, 96, 144, 192}	(36792, 5271, 10540)	Weather (10 mins)
	ILI	7	{24, 36, 48, 60}	(617, 74, 170)	Illness (Weekly)
	M4-Yearly	1	6	(23000, 0, 23000)	Demographic
	M4-Quarterly	1	8	(24000, 0, 24000)	Finance
Short-term	M4-Monthly	1	18	(48000, 0, 48000)	Industry
	M4-Weekly	1	13	(359, 0, 359)	Macro
	M4-Daily	1	14	(4227, 0, 4227)	Micro
	M4-Hourly	1	48	(414, 0, 414)	Other

Table 6: Forecasting dataset descriptions. The dataset size is organized in (Train, Validation, Test).

A.2. Implementation Details

We implement SKOLR in PyTorch, applying instance-normalization and denormalization (Kim et al., 2022) to inputs and predictions respectively. Following the protocol in the previous works (Zeng et al., 2023; Nie et al., 2023), SKOLR processes each $\mathbf{y}_{1:L,c}$ independently to generate the output $\hat{\mathbf{y}}_{L+1:L+T,c}$, and subsequently combine them to form a multivariate forecast. This technique is termed *channel-independence* and we omit the variate index *m* in order to simplify the notation in the Section 3.

To improve computational efficiency, we adopt non-overlapping patch tokenization (Nie et al., 2023) before feeding the frequency-decomposed signals $\{\mathbf{Y}_n\}_{n=1}^N$ into the linear RNN branches. This reduces the sequence length by a factor of P, significantly decreasing computation time while maintaining model effectiveness.

The model architecture employs a single-layer linear RNN to preserve linear state transitions between time steps, with

²https://www.bgc-jena.mpg.de/wetter

³https://pems.dot.ca.gov

 $^{^{4}}$ Wu et al. (2021) selected 321 of 370 customers from the original dataset in Trindade (2015). This version is widely used in the follow-up works.

⁵https://gis.cdc.gov/grasp/fluview/fluportaldashboard.html

MLPs using ReLU activation functions in both encoder and decoder components. The patch length adapts to the input window length as P = L/6. Through grid search, we optimize the branch number $N \in \{2, 3, 4, 8\}$, number of MLP layers $M \in \{1, 2, 3\}$ and dynamic dimension $D \in \{128, 256, 512\}$, while maintaining the hidden dimension at H = 2D. We train using AdamW optimizer (Loshchilov, 2017) with learning rate $1 \times e^{-4}$ and weight decay $5 \times e^{-4}$, using batch size 32 across all datasets. Complete hyperparameter configurations are detailed in Table 7.

Dataset	Traffic	ELC	Weather	ETTh1	ETTh2	ETTm1	ETTm2	ILI
Number of Branches N	2	2	2	2	2	2	2	2
Number of MLP hidden layers M	2	2	2	1	1	1	1	1
Dynamic Dimension D	256	256	128	256	128	256	256	256
Dropout	0.05	0.05	0.2	0.2	0.2	0.2	0.1	0.2

Table 7: Hyperparameters of SKOLR

B. Additional Experimental Results

B.1. Short-term Forecasting

B.1.1. EXPERIMENTAL SETTING

For the short-term forecasting, following the N-BEATS (Oreshkin et al., 2020), we adopt the symmetric mean absolute percentage error (sMAPE), mean absolute scaled error (MASE), and overall weighted average (OWA) as the metrics, where OWA is a special metric used in the M4 competition. These metrics can be calculated as follows:

$$\begin{split} \mathrm{sMAPE} &= \frac{200}{T} \sum_{i=1}^{T} \frac{|Y_i - \hat{Y}_i|}{|Y_i| + |\hat{Y}_i|},\\ \mathrm{MASE} &= \frac{1}{T} \sum_{i=1}^{T} \frac{|Y_i - \hat{Y}_i|}{\frac{1}{T - m} \sum_{j=m+1}^{T} |Y_j - Y_{j-m}|},\\ \mathrm{OWA} &= \frac{1}{2} \left(\frac{\mathrm{sMAPE}}{\mathrm{sMAPE_{Naïve2}}} + \frac{\mathrm{MASE}}{\mathrm{MASE_{Naïve2}}} \right), \end{split}$$

where m is the periodicity of the data. $Y, \hat{Y} \in \mathbb{R}^{T \times C}$ are the ground truth and prediction results of the future with T time points and C dimensions. Y_i represents the *i*-th future time point.

We compare SKOLR against state-of-the-art models on the M4 competition dataset, which contains six diverse domains (Yearly, Quarterly, Monthly, Weekly, Daily, Hourly) with prediction horizons ranging from 6 to 48 steps, as shown in Table 6. The baselines includes: N-BEATS (Oreshkin et al., 2020), N-HiTS (Challu et al., 2023), PatchTST (Nie et al., 2023), TimesNet (Wu et al., 2023), DLinear (Zeng et al., 2023), MICN (Wang et al., 2023a), KNF (Wang et al., 2023b), FiLM (Zhou et al., 2022), Autoformer (Wu et al., 2021), and KooPA (Liu et al., 2023). These models are not specifically designed for long-term forecasting, but also generalize well on short-term tasks.

B.1.2. RESULTS

Table 8 demonstrates SKOLR's effectiveness in handling diverse temporal patterns across M4 domains. At yearly, quarterly and monthly predictions, where data exhibits strong seasonality and multiple frequency components, SKOLR's parallel branch design allows simultaneous tracking of different temporal scales. Our structured Koopman operator design proves particularly powerful for the M4 dataset, which contains richer dynamic information than typical long-term forecasting benchmarks, resulting in consistently outperforming both traditional forecasting models (N-BEATS, N-HiTS) and other Koopman-based approaches (KooPA) across all evaluation metrics.

M4	Metric	SKOLR	KooPA	N-HiTS	N-BEATS	PatchTST	TimesNet	DLinear	MICN	KNF	FiLM	Autoformer
Year	sMAPE MASE OWA	13.291 2.996 0.784	$\frac{13.352}{\frac{2.997}{0.786}}$	13.371 3.025 0.790	13.466 3.059 0.797	13.517 3.031 0.795	13.394 3.004 0.787	13.866 3.006 0.802	14.532 3.359 0.867	13.986 3.029 0.804	14.012 3.071 0.815	14.786 3.349 0.874
Quarter	sMAPE MASE OWA	9.986 1.166 0.878	10.159 1.189 0.895	10.454 1.219 0.919	$\frac{10.074}{1.163}\\ \underline{0.881}$	10.847 1.315 0.972	10.101 1.183 0.890	10.689 1.294 0.957	11.395 1.379 1.020	10.343 1.202 0.965	10.758 1.306 0.905	12.125 1.483 1.091
Month	sMAPE MASE OWA	12.536 0.921 0.867	$\frac{12.730}{\frac{0.953}{0.901}}$	12.794 0.960 0.895	12.801 0.955 <u>0.893</u>	14.584 1.169 1.055	12.866 0.964 0.894	13.372 1.014 0.940	13.829 1.082 0.988	12.894 1.023 0.985	13.377 1.021 0.944	15.530 1.277 1.139
Others	sMAPE MASE OWA	4.652 3.233 0.999	4.861 3.124 1.004	$\frac{4.696}{3.130}$ 0.988	5.008 3.443 1.070	6.184 4.818 1.140	4.982 3.323 1.048	4.894 3.358 1.044	6.151 4.263 1.319	4.753 3.138 1.019	5.259 3.608 1.122	5.841 4.308 1.294
Average	sMAPE MASE OWA	11.704 1.572 0.843	$\frac{11.863}{1.595}\\ \underline{0.858}$	11.960 1.606 0.861	11.910 1.613 0.862	13.022 1.814 0.954	11.930 1.597 0.867	12.418 1.656 0.891	13.023 1.836 0.960	12.126 1.641 0.874	12.489 1.690 0.902	14.057 1.954 1.029

Table 8: Comparison of Models for short-term prediction. Best results and second best results are highlighted in **red** and blue respectively.

Table 9: Prediction results on benchmark datasets with L = 96. Best results and second best results are highlighted in red and blue respectively.

Dataset	T	SK MSE	OR MAE	Ko MSE	opa MAE	iTrans MSE	former MAE	Patch MSE	nTST MAE	Cross MSE	former MAE	TI MSE	DE MAE	Time MSE	esNet MAE	Dlin MSE	near MAE	SCI MSE	Net MAE	Statio MSE	onary MAE	Autof MSE	ormer MAE
ETTm1	96 192 336 720	0.313 0.359 0.389 0.449	0.356 0.384 0.406 0.443	0.330 0.385 0.402 0.472	0.368 0.395 0.413 0.449	0.334 0.377 0.426 0.491	0.368 0.391 0.420 0.459	$\frac{\frac{0.329}{0.367}}{\frac{0.399}{0.454}}$	$\frac{\frac{0.367}{0.385}}{\frac{0.410}{0.439}}$	0.404 0.450 0.532 0.666	0.426 0.451 0.515 0.589	0.364 0.398 0.428 0.487	0.387 0.404 0.425 0.461	0.338 0.374 0.410 0.478	0.375 0.387 0.411 0.450	0.345 0.380 0.413 0.474	0.372 0.389 0.413 0.453	0.418 0.439 0.490 0.595	0.438 0.450 0.485 0.550	0.386 0.459 0.495 0.585	0.398 0.444 0.464 0.516	0.505 0.553 0.621 0.671	0.475 0.496 0.537 0.561
ETTm2	96 192 336 720	0.173 0.239 0.302 0.398	0.256 0.300 0.341 0.396	$0.181 \\ 0.248 \\ \underline{0.303} \\ 0.403$	$\begin{array}{r} 0.263 \\ 0.308 \\ \underline{0.343} \\ 0.401 \end{array}$	0.180 0.250 0.311 0.412	0.264 0.309 0.348 0.407	$ \begin{array}{r} $	$\frac{\frac{0.259}{0.302}}{\frac{0.343}{0.400}}$	0.287 0.414 0.597 1.730	0.366 0.492 0.542 1.042	0.207 0.290 0.377 0.558	0.305 0.364 0.422 0.524	0.187 0.249 0.321 0.408	0.267 0.309 0.351 0.403	0.193 0.284 0.369 0.554	0.292 0.362 0.427 0.522	0.286 0.399 0.637 0.960	0.377 0.445 0.591 0.735	0.192 0.280 0.334 0.417	0.274 0.339 0.361 0.413	0.255 0.281 0.339 0.433	0.339 0.340 0.372 0.432
ETTh1	96 192 336 720	0.371 0.423 0.471 0.499	0.397 0.427 0.453 0.484	0.401 0.449 0.494 0.484	0.413 0.439 0.461 0.472	0.386 0.441 0.487 0.503	$0.405 \\ 0.436 \\ \underline{0.458} \\ 0.491$	0.414 0.460 0.501 0.500	0.419 0.445 0.466 0.488	0.423 0.471 0.570 0.653	0.448 0.474 0.546 0.621	0.479 0.525 0.565 0.594	0.464 0.492 0.515 0.558	$\begin{array}{c} 0.384\\ \hline 0.436\\ \hline 0.491\\ 0.521 \end{array}$	$0.402 \\ \underline{0.429} \\ 0.469 \\ 0.500$	$\begin{array}{c c} 0.386 \\ 0.437 \\ \underline{0.481} \\ 0.519 \end{array}$	$\begin{array}{r} \underline{0.400}\\ 0.432\\ 0.459\\ 0.516\end{array}$	0.654 0.719 0.778 0.836	0.599 0.631 0.659 0.699	0.513 0.534 0.588 0.643	0.491 0.504 0.535 0.616	0.449 0.500 0.521 0.514	0.459 0.482 0.496 0.512
ETTh2	96 192 336 720	0.293 0.370 0.410 0.431	0.344 0.384 0.428 0.446	$0.316 \\ 0.384 \\ \frac{0.423}{0.450}$	0.361 0.405 0.438 0.458	$\begin{array}{c} \underline{0.297} \\ \underline{0.380} \\ 0.428 \\ 0.427 \end{array}$	$0.349 \\ \underline{0.400} \\ \underline{0.432} \\ \underline{0.445}$	0.302 0.388 0.426 <u>0.431</u>	$\begin{array}{r} \underline{0.348} \\ \underline{0.400} \\ 0.433 \\ \underline{0.446} \end{array}$	0.745 0.877 1.043 1.104	0.584 0.656 0.731 0.763	0.400 0.528 0.643 0.874	0.440 0.509 0.571 0.679	0.340 0.402 0.452 0.462	0.374 0.414 0.452 0.468	0.333 0.477 0.594 0.831	0.387 0.476 0.541 0.657	0.707 0.860 1.000 1.249	0.621 0.689 0.744 0.838	0.476 0.512 0.552 0.562	0.458 0.493 0.551 0.560	0.346 0.456 0.482 0.515	0.388 0.452 0.486 0.511
ECL	96 192 336 720	$ \begin{array}{r} 0.153 \\ \underline{0.168} \\ \underline{0.189} \\ 0.230 \end{array} $	$\begin{array}{c} 0.246 \\ \underline{0.259} \\ \underline{0.282} \\ 0.318 \end{array}$	0.146 0.169 0.189 0.226	0.244 0.266 0.285 0.314	0.148 0.162 0.178 0.225	0.240 0.253 0.269 0.317	0.181 0.188 0.204 0.246	0.270 0.274 0.293 0.324	0.219 0.231 0.246 0.280	0.314 0.322 0.337 0.363	0.237 0.236 0.249 0.284	0.329 0.330 0.344 0.373	0.168 0.184 0.198 0.220	0.272 0.289 0.300 0.320	0.197 0.196 0.209 0.245	0.282 0.285 0.301 0.333	0.247 0.257 0.269 0.299	0.345 0.355 0.369 0.390	0.169 0.182 0.200 0.222	0.273 0.286 0.304 0.321	0.201 0.222 0.231 0.254	0.317 0.334 0.338 0.361
Traffic	96 192 336 720	$\frac{0.427}{0.455}\\ \frac{0.472}{0.519}$	$\frac{\frac{0.270}{0.289}}{\frac{0.298}{0.326}}$	0.462 0.566 0.514 0.552	0.290 0.386 0.331 0.346	0.395 0.417 0.433 0.467	0.268 0.276 0.283 0.302	0.462 0.466 0.482 <u>0.514</u>	0.295 0.296 0.304 <u>0.322</u>	0.522 0.530 0.558 0.589	0.290 0.293 0.305 0.328	0.805 0.756 0.762 0.719	0.493 0.474 0.477 0.449	0.593 0.617 0.629 0.640	0.321 0.336 0.336 0.350	0.650 0.598 0.605 0.645	0.396 0.370 0.373 0.394	0.788 0.789 0.797 0.841	0.499 0.505 0.508 0.523	0.612 0.613 0.618 0.653	0.338 0.340 0.328 0.355	0.613 0.616 0.622 0.660	0.388 0.382 0.337 0.408
Weather	96 192 336 720	0.162 0.208 0.266 0.344	0.207 0.249 0.292 0.343	0.157 0.209 0.266 0.350	0.202 0.251 0.290 0.348	0.174 0.221 0.278 0.358	0.214 0.254 0.296 0.347	0.177 0.225 0.278 0.354	0.218 0.259 0.297 0.348	0.158 0.206 0.272 0.398	0.230 0.277 0.335 0.418	0.202 0.242 0.287 0.351	0.261 0.298 0.335 0.386	0.172 0.219 0.280 0.365	0.220 0.261 0.306 0.359	0.196 0.237 0.283 0.345	0.255 0.296 0.335 0.381	0.221 0.261 0.309 0.377	0.306 0.340 0.378 0.427	0.173 0.245 0.321 0.414	0.223 0.285 0.338 0.410	0.266 0.307 0.359 0.419	0.336 0.367 0.395 0.428

B.2. Results for a shorter look-back window

For a fair comparison, we also conduct the experiments under the L = 96 setting that is the default for iTransformer, TimesNet, and other transformer-based models, except PatchTST. For the shorter look-back window, we use the patch length P = 16 to obtain the tokens for SKOLR. By default, SKOLR follow the hyperparameters in Table 7. As shown in Table 9, SKOLR emerges as the leading performer, achieving Rank 1 or 2 in 24 cases out of 28 cases in terms of MSE and 25 cases in terms of MAE.

Dataset	Models	MSE	MAE
	48 96	$\begin{array}{c} 0.137 \pm 0.0003 \\ 0.132 \pm 0.0005 \end{array}$	0.229 ± 0.0003 0.225 ± 0.0004
ECL	144	0.143 ± 0.0001	0.236 ± 0.0001
	192	0.149 ± 0.0001	0.244 ± 0.0001
	48	0.400 ± 0.0003	0.258 ± 0.0040
Traffic	96	0.368 ± 0.0007	0.248 ± 0.0007
	144	0.375 ± 0.0003	0.255 ± 0.0002
	192	0.377 ± 0.0003	0.256 ± 0.0002
	48	0.131 ± 0.0009	0.170 ± 0.0008
Weather	96	0.154 ± 0.0015	0.202 ± 0.0015
,, outilor	144	0.172 ± 0.0009	0.220 ± 0.0006
	192	0.193 ± 0.0004	0.241 ± 0.0005
	48	0.280 ± 0.0013	0.330 ± 0.0015
ETTm1	96	0.287 ± 0.0003	0.340 ± 0.0001
	144	0.313 ± 0.0020	0.361 ± 0.0023
	192	0.328 ± 0.0019	0.373 ± 0.0018
	48	0.134 ± 0.0011	0.228 ± 0.0007
ETTm2	96	0.171 ± 0.0015	0.255 ± 0.0013
EIIIIZ	144	0.209 ± 0.0014	0.283 ± 0.0014
	192	0.241 ± 0.0013	0.304 ± 0.0015
	48	0.333 ± 0.0009	0.373 ± 0.0007
ETTh 1	96	0.371 ± 0.0011	0.398 ± 0.0008
LIIII	144	0.405 ± 0.0019	0.417 ± 0.0020
	192	0.422 ± 0.0030	0.432 ± 0.0034
	48	0.238 ± 0.0012	0.306 ± 0.0004
ETTL2	96	0.299 ± 0.0034	0.352 ± 0.0042
EIIIIZ	144	0.335 ± 0.0042	0.377 ± 0.0048
	192	0.365 ± 0.0033	0.397 ± 0.0040
	24	1.556 ± 0.0213	0.760 ± 0.0159
пт	36	1.462 ± 0.0711	0.728 ± 0.0676
ILI	48	1.537 ± 0.0038	0.798 ± 0.0030
	60	2.187 ± 0.0435	0.995 ± 0.0498

Table 10: Model performance across different datasets with mean \pm standard deviation for MSE and MAE metrics.

B.3. Experimental Variability

Table 10 reports standard deviation (std) across 3 independent runs for all datasets and forecast horizons. The low stds (<0.003 for most datasets) demonstrate the consistency of SKOLR's performance.

B.4. Comparison with Orvieto et al. (2023)

The Linear Recurrent Unit (LRU) presented by Orvieto et al. (2023) is derived from vanilla recurrent neural networks (RNNs) through a sequence of principled modifications including linearization of the recurrence, diagonalization, exponential parametrization for stability, and forward-pass normalization. These changes yield a model that can match the performance of recent deep state-space models (SSMs) such as S4 and S5 on benchmarks like the Long Range Arena (Tay et al., 2021), without relying on discretization of continuous-time dynamics.

We implement our code in PyTorch. Our implementation follows from the JAX pseudocode presented in the original paper's appendix (Orvieto et al., 2023). Additionally, we consulted a community implementation in PyTorch⁶. The LRU is trained using the AdamW optimizer with no weight decay applied to the recurrent parameters. Learning rates are selected via grid search on a logarithmic scale. All experiments use networks of 6 LRU layers with residual and normalization layers between blocks and a final linear output layer and a 64-dimensional hidden state.

⁶https://github.com/Gothos/LRU-pytorch

	SKO	OLR	Ko	oPA	LRU		
Dataset	MSE	MAE	MSE	MAE	MSE	MAE	
Pendulum	0.0001	0.0083	0.0039	0.0470	0.0572	0.0242	
Duffing	0.0047	0.0518	0.0365	0.1479	0.0573	0.5970	
Lotka-Volterra	0.0018	0.0354	0.0178	0.1050	0.2058	0.3779	
Lorenz '63	0.9740	0.7941	1.0937	0.8325	1.1905	0.8932	

Table 11: Performance comparison of LRU, Koopa and SKOLR on non-linear dynamical systems (NLDS)

Whereas our focus in SKOLR is time-series forecasting, Orvieto et al. (2023) target long-range reasoning. Although it is possible to convert their architecture to address forecasting, performance suffers because it is not the design goal, as shown in Table 11.

C. Model Efficiency

C.1. Theoretical Complexity Analysis

SKOLR achieves computational efficiency through its structured design and linear operations. For a time series of length L with patch length P, embedding dimension D, and N branches, we analyze both time and space complexity.

The time complexity of SKOLR consists of several components: spectral decomposition, encoder/decoder MLPs, and linear RNN computation. If we perform a single FFT operation $O(L \log L)$ followed by branch-specific frequency filtering, the main computational cost comes from encoder/decoder MLPs $O(N \times (L/P) \times D^2)$ and linear RNN computation $O(N \times (L/P) \times D^2)$, giving a total time complexity of $O(N \times (L/P) \times D^2)$. The memory complexity includes model parameters $O(N \times D^2)$ and activation memory $O(N \times (L/P) \times D)$.

Our structured approach with N branches provides substantial efficiency gains compared to a non-structured approach with equivalent representational capacity. For a non-structured model with dimension $D' = N \times D$, the time complexity would be $O((L/P) \times N^2 D^2)$ and memory complexity $O(N^2 D^2)$. This represents an N-fold increase in computational requirements. For example, with N = 16 branches, our structured approach requires approximately $16 \times$ fewer parameters and operations while maintaining equivalent or better modeling capacity, as shown in Section 4.3.1. Compared to transformer-based approaches with time complexity $O((L/P)^2 \times D + (L/P) \times D^2)$ and memory complexity $O((L/P)^2 + (L/P) \times D)$, SKOLR demonstrates a fundamental advantage for long sequences by avoiding the quadratic scaling with sequence length.

C.2. Parallel Computing

SKOLR further benefits from parallel processing capabilities. The N separate branches can be processed completely independently, reducing the effective time complexity to $O((L/P) \times D^2)$ with sufficient parallel resources. This branch-level parallelism is implemented in our current code.

In future work, we plan to implement additional parallelization of the linear RNN computation itself. Since our RNN has no activation functions, we can express the hidden state evolution for each branch with sequence length L/P in closed form: $h_k = g(y_k) + \sum_{s=1}^{L/P} W^s g(y_{k-s})$, where W^s indicates s applications of W (the state transition matrix). This formulation allows us to compute all hidden states simultaneously through efficient matrix operations, potentially reducing the time complexity further to $O(D^3 \log(L/P) + (L/P)^2 \times D)$ per branch.

For time series with patching where $L/P \ll D$, this approach achieves significant speedups by eliminating the sequential dependency in RNN computation. With both branch and RNN parallelism implemented, SKOLR can achieve greater computational efficiency while maintaining its forecasting performance.

C.3. Computational Efficiency

We provided efficiency results on the ETTm2 and Traffic datasets in Fig. 5. We have expanded our evaluation across additional datasets to offer a more comprehensive assessment in Table 12. In all datasets, the proposed architecture provides a compelling trade-off between efficiency and accuracy compared to baselines.

Table 12: Model Efficiency and Performance Comparison for Different Datasets with T = 96. Parameters (Params) are measured in millions (M), GPU memory (GPU) in MiB, computation time per epoch in seconds (s) on NVIDIA V100 GPU with batch size 32.

(b) Electricity

Model	Params (M)	GPU(MiB)	Time (s)	MSE	Mod	el	Params (M)	GPU(MiB)	Time (s)	MSE
Autoformer	14.914	18.811	51.0	0.668	Auto	former	11.214	17.373	68.7	0.182
iTransformer	6.405	62.710	126.0	0.388	iTran	sformer	4.957	86.478	58.6	0.134
PatchTST	3.755	22.132	1042.0	0.413	Patch	nTST	6.904	73.517	1231.0	0.143
MICN	236.151	32.310	84.0	0.511	MIC	N	6.635	32.668	18.0	0.165
TimesNet	30.170	111.998	6563.0	0.611	Time	sNet	15.037	33.435	11351.0	0.170
DLinear	0.009	12.861	7.7	0.485	DLin	near	0.019	76.016	6.8	0.153
Koopa	5.429	50.335	25.5	0.401	Koop	ba	4.076	31.067	33.1	0.136
SKOLR	1.479	5.915	216.0	0.368	SKO	LR	1.541	6.163	99.1	0.132
	(c)	ETTh1					(d) E	TTm2		
Model	Params (M)	GPU(MiB)	Time (s)	MSE	Mode	l	Params (M)	GPU(MiB)	Time (s)	MSE
Autoformer	10.536	16.523	29.5	0.634	Autof	ormer	10.536	14.599	152.6	0.241
iTransformer	0.237	27.245	4.1	0.393	iTrans	former	0.237	27.245	13.1	0.177
PatchTST	3.752	22.018	8.5	0.372	Patch	ГЅТ	10.056	39.910	980.0	0.171
MICN	252.001	65.974	21.1	0.406	MICN	1	252.001	65.974	84.2	0.197
TimesNet	0.605	26.053	22.1	0.411	Times	Net	1.192	34.783	113.0	0.187
DLinear	0.140	26.440	0.6	0.379	DLine	ear	18.291	9.312	1.9	0.172
Koopa	0.135	31.951	10.1	0.371	Koopa	a	0.135	31.951	48.2	0.171
SKOLR	0.429	1.717	2.8	0.371	SKOL	R	0.429	1.717	12.6	0.171

D. Additional Analysis

D.1. Scaling Up Forecast Horizon

(a) Traffic

We have conducted experiments to explore performance in the setting where the forecast horizon is increased at test-time. In this experiment, SKOLR and Koopa were evaluated by scaling up from the training horizon (T_{tr}) to a larger test horizon (T_{te}) . Unlike Koopa (Liu et al., 2023), SKOLR does not incorporate an operator adaptation (OA) mechanism to update its Koopman operator using incoming ground truth. Instead, our architecture possesses a natural recursive structure that enables straightforward extension to longer horizons. Even when weights are trained to minimize a loss function specified over a given horizon, the algorithm can be recursively applied to predict over extended periods.

As demonstrated in Table 13, SKOLR maintains competitive performance without requiring additional adaptation mechanisms. The structured Koopman operator and linear RNN design enable robust long-term predictions, with error percentages remaining comparable to Koopa OA across various datasets. This demonstrates SKOLR's inherent capability to handle extended forecast horizons efficiently through its recursive architecture.

D.2. Ablation Study

We have also conducted a more comprehensive ablation study on the design elements of SKOLR. As shown in Table 14, we compare our full SKOLR model with two ablated variants: (1) "w/o Structure": no structured decomposition, using a single branch with dimension (N×D); (2) "w/o Spectral Encoder": no learnable frequency decomposition, while maintaining the multi-branch structure.

The results show that both components contribute meaningfully. Removing the structured decomposition leads to performance degradation on 27/32 tasks, with notable declines on ETTh1 and ILI, while increasing computational overhead. Similarly, removing the spectral encoder impacts performance on 23/32 tasks, though with a smaller overall effect.

	ETTh2		ILI		ECL		Traffic		Weather	
	(ADF -4.135)		(ADF -5.406)		(ADF -8.483)		(ADF -15.046)		(ADF -26.661)	
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Koopa (T_tr)	0.226	0.300	1.621	0.800	0.130	0.234	0.415	0.274	0.126	0.168
Koopa(T_te)	0.437	0.429	2.836	1.065	0.199	0.298	0.709	0.437	0.237	0.276
Error(+ %)	93%	43%	75%	33%	53%	27%	71%	59%	88%	64%
Koopa OA(T_te)	0.372	0.404	2.427	0.907	0.182	0.271	0.699	0.426	0.225	0.264
Error(+ %)	65%	35%	50%	13%	40%	16%	68%	55%	79%	57%
SKOLR (T_tr)	0.238	0.306	1.556	0.760	0.137	0.229	0.400	0.258	0.131	0.170
SKOLR (T_te)	0.393	0.402	2.392	0.958	0.204	0.289	0.612	0.383	0.222	0.257
Error(+ %)	65%	31%	54%	26%	49%	26%	53%	48%	69%	51%

Table 13: Scaling up forecast horizon: $(T_tr, T_te) = (24, 48)$ for ILI and $(T_tr, T_te) = (48, 144)$ for others. Koopa and SKOLR conducts vanilla rolling forecast and Koopa OA has operator adaptation.

Table 14: Ablation study comparing SKOLR with versions without structure and without spectral encoder

Dataset	Т	SKO	OLR	w/o St	ructure	w/o Spectral Encoder		
		MSE	MAE	MSE	MAE	MSE	MAE	
ECL	48	0.137	0.229	0.148	0.238	0.149	0.238	
	96	0.132	0.225	0.135	0.228	0.133	0.227	
	144	0.143	0.236	0.146	0.241	0.142	0.235	
	192	0.149	0.244	0.150	0.245	0.148	0.243	
Traffic	48	0.400	0.258	0.395	0.255	0.397	0.257	
	96	<u>0.368</u>	0.248	0.367	<u>0.249</u>	0.369	0.249	
	144	0.375	0.255	0.375	0.255	0.375	0.255	
	192	0.377	0.256	0.378	0.256	0.377	0.256	
Weather	48	0.131	0.170	<u>0.134</u>	0.173	0.134	<u>0.172</u>	
	96	0.154	0.202	<u>0.157</u>	0.203	0.158	0.202	
	144	0.172	0.220	0.177	0.225	<u>0.175</u>	0.221	
	192	0.193	0.241	0.195	0.242	0.197	0.244	
ETTm1	48	0.280	0.330	0.284	0.334	0.282	<u>0.332</u>	
	96	0.287	0.340	0.292	0.343	<u>0.291</u>	0.342	
	144	0.313	0.361	0.325	0.365	<u>0.319</u>	0.361	
	192	0.328	0.373	0.332	0.372	0.332	0.372	
ETTm2	48	0.134	0.228	0.135	0.229	0.162	0.259	
	96	<u>0.171</u>	0.255	0.174	0.259	0.169	0.253	
	144	<u>0.209</u>	0.283	0.206	0.280	<u>0.209</u>	<u>0.282</u>	
	192	0.241	0.304	0.241	0.305	0.230	0.299	
ETTh1	48	0.333	0.373	0.338	0.377	0.336	<u>0.374</u>	
	96	0.371	0.398	0.387	0.408	<u>0.373</u>	<u>0.399</u>	
	144	0.405	0.417	0.414	0.423	<u>0.410</u>	<u>0.420</u>	
	192	0.422	0.432	0.409	0.421	<u>0.413</u>	<u>0.422</u>	
ETTh2	48	0.238	0.306	0.233	0.304	0.239	<u>0.305</u>	
	96	0.299	0.352	<u>0.301</u>	<u>0.350</u>	0.303	0.350	
	144	0.335	0.377	0.341	0.382	<u>0.337</u>	<u>0.381</u>	
	192	0.365	0.397	0.370	<u>0.398</u>	0.370	0.401	
ILI	24	1.556	0.760	1.795	0.842	1.522	0.741	
	36	1.462	0.728	1.990	0.889	<u>1.496</u>	<u>0.734</u>	
	48	1.537	0.798	1.875	0.909	<u>1.571</u>	<u>0.810</u>	
	60	2.187	0.995	2.407	1.056	<u>2.263</u>	<u>0.999</u>	



Figure 6: Analysis of SKOLR's branch-wise behavior on ETTm2 feature 6: (a) frequency decomposition and (b) prediction performance.



Figure 7: Analysis of SKOLR's branch-wise behavior on Electricity feature 11: (a) frequency decomposition and (b) prediction performance.

D.3. Branch-wise Visualization

We add examples from ETTm2 and Electricity in order to analyze SKOLR's branch-wise behavior. The figures show distinct frequency specializations. In Fig. 6, SKOLR decomposes the time series into complementary frequency components, with Branch 1 focusing on lower, distributed frequencies, and Branch 2 capturing more specific dominant frequency peaks. In Fig. 7, Branch 1 shows higher amplitudes across most of the very low $(0 - 20\mu Hz)$ frequency components compared to Branch 2. The time-domain plots demonstrate how these spectral differences translate into signal reconstruction; Branch 2 focuses more on prediction of the higher-frequency components of the time series.

D.4. Analysis of Error Accumulation

Time-series forecasting has two main prediction methods: direct prediction, which forecasts the entire horizon at once but is parameter-inefficient and cannot extend the prediction horizon after training, and recursive prediction, which iteratively uses predictions as inputs but may suffer from error accumulation over long sequences.

In SKOLR, we use a patching approach (Appendix A.2) to create an effective middle ground between these methods. Instead of operating at the individual timestep level, we work with patches of multiple timesteps, directly predicting all values within each patch while only applying recursion between patches. This dramatically reduces the number of recursive steps (e.g., from 720 to just 5 with patch length 144), controlling error accumulation. Additionally, this method also reduces complexity to $O(\frac{L}{P})$) from RNN standard timestamp-based approaches (O(L)) while maintaining the core principle of Koopman operator theory.



Figure 8: Error progression across 720 time steps: SKOLR (multiple patch sizes) vs. iTransformer

To better address this question, we add an experiment on a longer horizon L = 720, T = 720 on dataset ETTm2 with varying iteration steps. We vary the patch length $P = \{16, 24, 48, 144, 240\}$ of the SKOLR model to see the difference performance caused by the number of iterations. SKOLR with P = 16 requires 45 recursive steps, while P = 240 needs only 3, yet they maintain comparable error profiles in Fig. 8. This empirically demonstrates that SKOLR's patch-based approach effectively controls error accumulation, even with increased recursion.

We also compare SKOLR with iTransformer (Liu et al., 2024), which performs direct prediction without recurrence. Both models show similar patterns of error increase with longer horizons, suggesting that this modest increase is inherent to all forecasting approaches when extending the prediction range, rather than being caused by recurrent error accumulation.

D.5. Analysis on Koopman operator eigenvalue

The Koopman modes are derived through eigendecomposition of the RNN weight matrices M_i . These modes represent dynamical patterns in the data. Each mode captures specific components of the time series. The stability and oscillatory behavior of each mode is determined by the corresponding eigenvalue's position in the complex plane. The eigenvalue plots in Fig. 9 show that each branch learns complementary spectral properties, with all eigenvalues within the unit circle, indicating stable dynamics. Branch 1 shows concentration at magnitude 0.4, while Branch 2 exhibits a more uniform distribution.

Moreover, this observation provides some reassurance against error accumulation concerns, as error divergence is more likely for unstable systems. Our learned system's stability encourages error effects to naturally decay over time during forward prediction steps rather than compounding.

E. State Prediction for Nonlinear Dynamical Systems (NLDS)

We generated datasets for four nonlinear dynamical systems (NLDS) to evaluate the performance of Koopman-based models in state prediction tasks. Each dataset contains a trajectory with a total of 20000 time steps. The first 14000 steps were designated for training, 2000 steps were used for validation, while the remaining 4,000 steps were used for inference. Below, we describe the generation process for each system:

• **Pendulum:** A simple nonlinear pendulum system described by its angular displacement and velocity. The trajectory was initialized with random a angle and angular velocity, with updates computed using the equations of motion under gravity. The system being simulated is a simple pendulum, consisting of a mass m attached to the end of a rigid, massless rod of length l. The pendulum swings in a two-dimensional plane under the influence of gravity, with the gravitational acceleration g. The motion of the pendulum is governed by the equation of motion:



Koopman Operator Eigenvalues: traffic

Figure 9: Koopman operator eigenvalue analysis for SKOLR on the Traffic dataset

$$\theta''(t) + \frac{g}{l}\sin(\theta(t)) = 0 \tag{17}$$

where:

- $\theta(t)$ is the angular displacement of the pendulum at time t,
- $\theta''(t)$ is the angular acceleration,
- g is the acceleration due to gravity (9.81 m/s^2) ,
- l is the length of the pendulum (set to 1.0 m).

The system is further characterized by its initial conditions:

- The initial angle θ_0 , which is randomly chosen from a uniform distribution between $-\pi$ and π ,
- The initial angular velocity ω_0 , which is randomly chosen from a uniform distribution between -1 rad/s and 1 rad/s.

The motion of the pendulum is modeled using numerical methods, specifically the Euler method, which approximates the solution of the system of equations over discrete time steps.

• **Duffing Oscillator:** A nonlinear oscillator characterized by damping and cubic stiffness terms. Trajectories were generated using randomized initial positions and velocities, with dynamics influenced by an external periodic driving force. The system modeled by the code is a Duffing oscillator, a type of nonlinear second-order differential equation

commonly used to describe systems with nonlinear restoring forces and damping. The equation of motion for the Duffing oscillator is given by:

$$\ddot{x} + \delta \dot{x} + \alpha x + \beta x^3 = \gamma \cos(\omega t)$$

where x(t) represents the displacement of the oscillator, $y(t) = \dot{x}(t)$ represents its velocity, and t is time. The parameters of the system are: $\alpha = 1.0$ (linear stiffness), $\beta = 5.0$ (nonlinear stiffness), $\delta = 0.3$ (damping coefficient), $\gamma = 8.0$ (driving force amplitude), and $\omega = 0.5$ (angular frequency of the driving force). The system undergoes periodic driving forces, and its motion is influenced by both the nonlinear restoring force and damping. The motion of the oscillator is simulated by numerically integrating the equations of motion using a simple time-stepping method, where dt is the time step, and the initial conditions for x and y are randomly selected within a small range. The system's behavior is characterized by chaotic dynamics for the chosen parameter values.

• Lotka-Volterra: A predator-prey population model, where the prey and predator populations interact dynamically. Trajectories were initialized with random population sizes, and updates followed the Lotka-Volterra equations. The equations governing the Lotka-Volterra predator-prey model are given by:

$$\frac{dN_{\text{prey}}}{dt} = \alpha N_{\text{prey}} - \beta N_{\text{prey}} N_{\text{predator}}$$
$$\frac{dN_{\text{predator}}}{dt} = \delta N_{\text{prey}} N_{\text{predator}} - \gamma N_{\text{predator}}$$

where:

- N_{prey} is the population of the prey species,
- N_{predator} is the population of the predator species,
- α is the natural growth rate of the prey,
- β is the predation rate (rate at which predators kill prey),
- δ is the rate at which predators increase due to consuming prey,
- γ is the natural death rate of the predator.

In this model, the prey species grows exponentially in the absence of predators, and the predator species declines exponentially in the absence of prey. The interaction between the species causes cyclical fluctuations in their populations.

We implement this model by numerically integrating the differential equations using a simple Euler method. The process starts by initializing the prey and predator populations randomly within a given range. The parameters $\alpha = 1.1$, $\beta = 0.4$, $\delta = 0.1$, and $\gamma = 0.4$ are then used to update the populations at each time step.

• Lorenz '63: A chaotic system described by three variables: x, y, and z. Each trajectory is started with randomized initial conditions, and updated using the Lorenz equations with standard parameters. The equations governing the Lorenz system are given by:

$$\frac{dx}{dt} = \sigma(y - x)$$
$$\frac{dy}{dt} = x(\rho - z) - y$$
$$\frac{dz}{dt} = xy - \beta z$$

where:

- -x, y, and z represent the state variables of the system, typically interpreted as the variables describing the convection rolls in the atmosphere,
- σ is the Prandtl number, a measure of the fluid's viscosity, set to 10.0,

- ρ is the Rayleigh number, representing the temperature difference between the top and bottom of the fluid, set to 28.0,
- β is a geometric factor, set to $\frac{8}{3}$.

The Lorenz system exhibits chaotic behavior for these parameter values, meaning that small differences in initial conditions can lead to vastly different outcomes over time. In the simulation, the system of differential equations is solved using the Euler method over a series of time steps. A visualization of the system is shown in Fig. 10.



Figure 10: Lorentz '63 system plotted in 3D