

# **Are Transformers Effective for Time Series Forecasting?**

**Ailing Zeng<sup>1</sup>, Muxi Chen<sup>1</sup>, Lei Zhang<sup>2</sup>, Qiang Xu<sup>1</sup>**

**<sup>1</sup>The Chinese University of Hong Kong**

**<sup>2</sup>International Digital Economy Academy (IDEA)**

# The Discussion Around Transformers for TSF

## Do We Really Need Deep Learning Models for Time Series Forecasting?

Shereen Elsayed\*, Daniela Thyssens\*, Ahmed Rashed, Hadi Samer Jomaa, and Lars Schmidt-Thieme

Department of Computer Science  
University of Hildesheim  
31141 Hildesheim, Germany

## Large Language Models Are Zero-Shot Time Series Forecasters

Nate Gruver\* NYU Marc Finzi\* CMU Shikai Qiu\* NYU Andrew Gordon Wilson NYU

Valeriy M., PhD, MBA, CQF ✅ @predict\_addict · Jan 31 · ...  
Exciting News! 🎉

My recent tweet on "Are Transformers Effective for Time Series Forecasting?" has taken off, reaching over 130K views in a very short time! The response has been nothing short of spectacular.

Inspired by this, I've updated my repo...

Show more

## Transformers\_Are\_What\_You\_Dont\_Need

The best repository showing why transformers don't work in time series forecasting

1. [Are Transformers Effective for Time Series Forecasting?](#) by Ailing Zeng, Muxi Chen, Lei Zhang, Qiang Xu (The Chinese University of Hong Kong, International Digital Economy Academy (IDEA), 2022) [code](#) 🔥🔥🔥
2. [LLMs and foundational models for time series forecasting: They are not \(yet\) as good as you may hope](#) by Christoph Bergmeir (2023) 🔥🔥🔥
3. [Transformers Are What You Do Not Need](#) by Valeriy Manokhin (2023) 🔥🔥🔥
4. [Deep Learning is What You Do Not Need](#) by Valeriy Manokhin (2022) 🔥🔥🔥
5. [Why do Transformers suck at Time Series Forecasting](#) by Devansh (2023)
6. [Frequency-domain MLPs are More Effective Learners in Time Series Forecasting](#) by Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Defu Lian, Ning An, Longbing Cao, Zhendong Niu (Beijing Institute of Technology, Tongji University, University of Oxford, University of Technology Sydney, University of Macau, HeFei University of Technology, Macquarie University) (2023) 🔥🔥🔥
7. [Forecasting CPI inflation under economic policy and geo-political uncertainties](#) by Shovon Sengupta, Tanujit Chakraborty, Sunny Kumar Singh (Fidelity Investments, Sorbonne University, BITS Pilani Hyderabad). (2024) 🔥🔥🔥
8. [Revisiting Long-term Time Series Forecasting: An Investigation on Linear Mapping](#) by Zhe Li, Shiyi Qi, Yiduo Li, Zenglin Xu (Harbin Institute of Technology, Shenzhen, 2023) [code](#)
9. [SCI-Net: Time Series Modeling and Forecasting with Sample Convolution and Interaction](#) by Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qixia Lai, Lingna Ma, Qiang Xu (The Chinese University of Hong Kong)

5

59

323

35K

## TSMixer: An All-MLP Architecture for Time Series Forecasting

# THE HYPE FOR TRANSFORMERS IN TIME SERIES FORECASTING

## TRANSFORMERS IN TSF, IRL



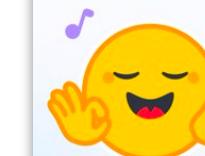
## Yes, Transformers are Effective for Time Series Forecasting (+ Autoformer)

Published June 16, 2023

Update on GitHub

elisim Eli Simhayev guest kashif KashifRasul nielsr Niels Rogge

Open in Colab



Transformers are **effective** for Time Series Forecasting

## Are Transformers Effective for Time Series Forecasting?

Ailing Zeng<sup>1\*</sup>, Muxi Chen<sup>1\*</sup>, Lei Zhang<sup>2</sup>, Qiang Xu<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong

<sup>2</sup>International Digital Economy Academy (IDEA)

{alzeng, mxchen21, qxu}@cse.cuhk.edu.hk  
(leizhang)@idea.edu.cn

### Abstract

Recently, there has been a surge of Transformer-based solutions for the long-term time series forecasting (LTSF) task. Despite the growing performance over the past few years, we question the validity of this line of research in this work. Specifically, Transformers is arguably the most successful solution to extract the semantic correlations among the elements in a long sequence. However, in time series modeling, we are to extract the temporal relations in an ordered set of continuous points. While employing positional encoding and using tokens to embed sub-series in Transformers facilitate preserving some ordering information, the nature of the permutation-invariant self-attention mechanism inevitably results in temporal information loss.

ergy management, and financial investment. Over the past several decades, TSF solutions have undergone a progression from traditional statistical methods (e.g., ARIMA [1]) and machine learning techniques (e.g., GBRT [11]) to deep learning-based solutions, e.g., Recurrent Neural Networks [15] and Temporal Convolutional Networks [3, 17].

Transformer [26] is arguably the most successful sequence modeling architecture, demonstrating unparalleled performances in various applications, such as natural language processing (NLP) [7], speech recognition [8], and computer vision [19, 29]. Recently, there has also been a surge of Transformer-based solutions for time series analysis, as surveyed in [27]. Most notable models, which focus on the less explored and challenging long-term time series forecasting (LTSF) problem, include LogTrans [16]



Simon Judge @mrprediction3 · Jun 23, 2023

The paper 'Are Transformers Effective for Time Series Forecasting?' asked whether **Transformers** really **are** required when much simpler models **are** better. There's now a rebuttal showing that **Transformers** are better than DLinear.

[mrprediction.com/transformers-f...](#)

#machinelearning

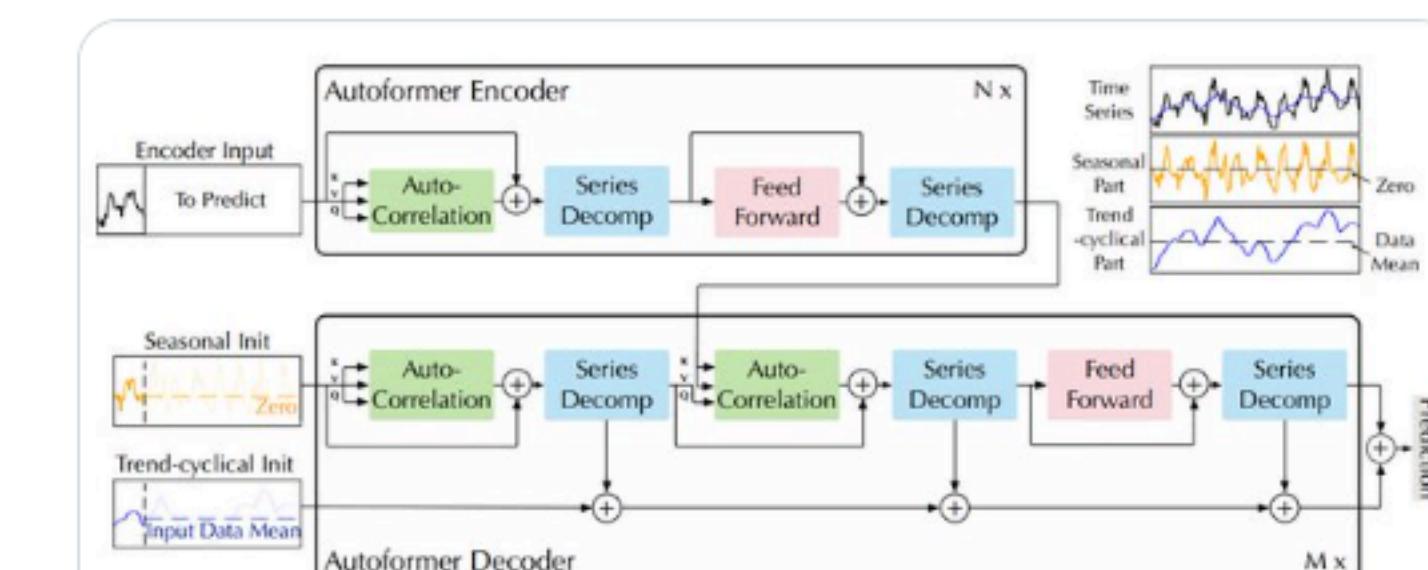


Figure 1: Autoformer architecture. The encoder eliminates the long-term trend-cyclical part by series decomposition blocks (blue blocks) and focuses on seasonal patterns modeling. The decoder accumulates the trend part extracted from hidden variables progressively. The past seasonal information from encoder is utilized by the encoder-decoder Auto-Correlation (center green block) in decoder.

# Paper Preview

## Are Transformers Effective for Time Series Forecasting?

Ailing Zeng<sup>1\*</sup>, Muxi Chen<sup>1\*</sup>, Lei Zhang<sup>2</sup>, Qiang Xu<sup>1</sup>

<sup>1</sup>The Chinese University of Hong Kong

<sup>2</sup>International Digital Economy Academy (IDEA)

{alzeng, mxchen21, qxu}@cse.cuhk.edu.hk

{leizhang}@idea.edu.cn

### Abstract

Recently, there has been a surge of Transformer-based solutions for the long-term time series forecasting (LTSF) task. Despite the growing performance over the past few years, we question the validity of this line of research in this work. Specifically, Transformers is arguably the most successful solution to extract the semantic correlations among the elements in a long sequence. However, in time series modeling, we are to extract the temporal relations in an ordered set of continuous points. While employing positional encoding and using tokens to embed sub-series in Transformers facilitate preserving some ordering information, the nature of the permutation-invariant self-attention mechanism inevitably results in temporal information loss.

ergy management, and financial investment. Over the past several decades, TSF solutions have undergone a progression from traditional statistical methods (e.g., ARIMA [1]) and machine learning techniques (e.g., GBRT [11]) to deep learning-based solutions, e.g., Recurrent Neural Networks [15] and Temporal Convolutional Networks [3, 17].

Transformer [26] is arguably the most successful sequence modeling architecture, demonstrating unparalleled performances in various applications, such as natural language processing (NLP) [7], speech recognition [8], and computer vision [19, 29]. Recently, there has also been a surge of Transformer-based solutions for time series analysis, as surveyed in [27]. Most notable models, which focus on the less explored and challenging long-term time series forecasting (LTSF) problem, include LogTrans [16]

With the above, we conclude that the temporal modeling capabilities of Transformers for time series are exaggerated, at least for the existing LTSF benchmarks. At the

**Claim:** Transformers' MHA results in temporal information loss, impacting LTSF performance.

**Their Model(s): LTSF-Linear**

**Experiments: 9X datasets**

**Additional Study: TSF Transformers design elements**

# Time Series Forecasting (TSF)

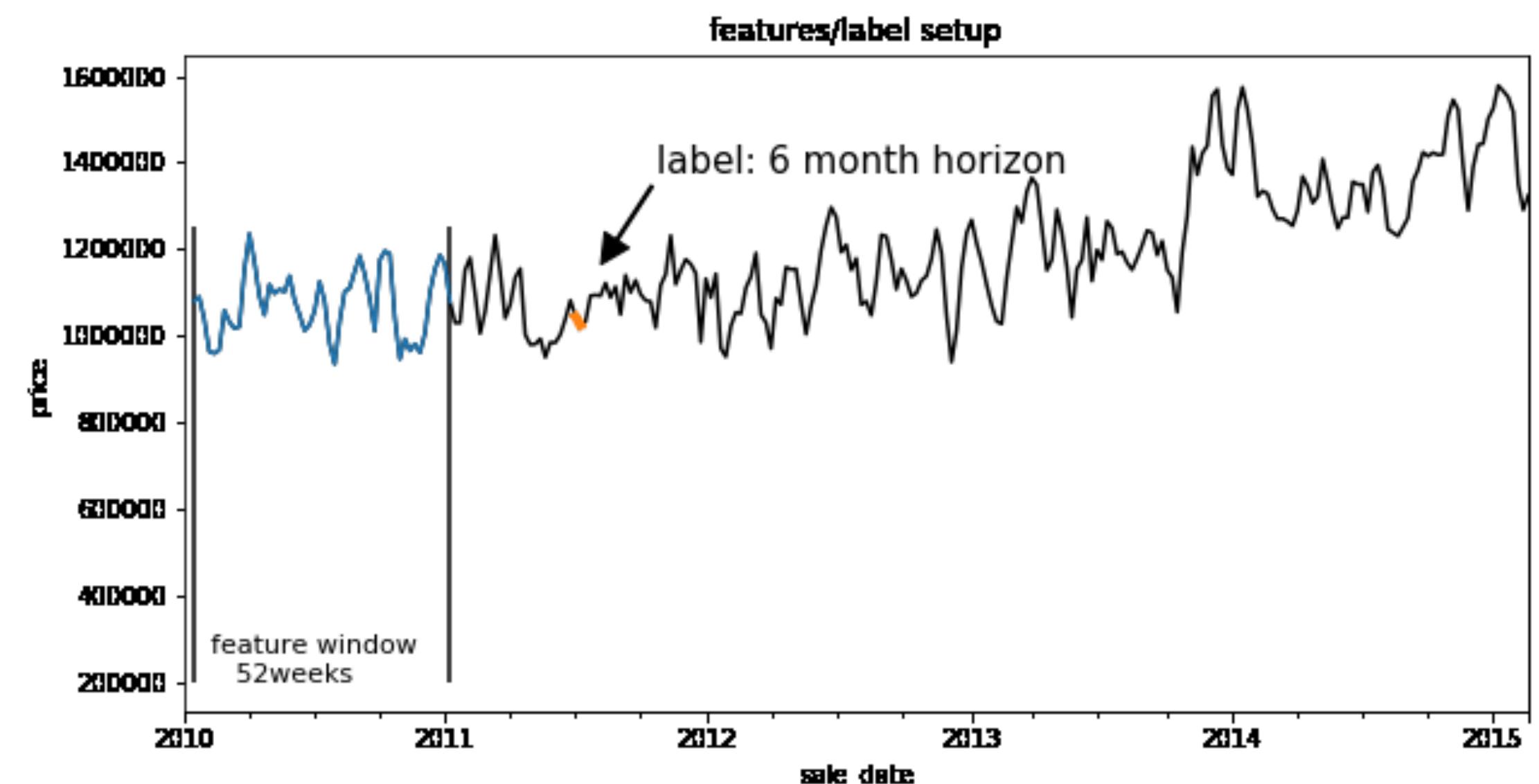
- Time series - ubiquitous in today's world
- TSF Solutions:
  - Statistical Methods -
    - ARIMA (1960s)
  - Machine Learning Techniques -
    - GBRT (2000s)
  - Deep Learning Techniques -
    - RNNs (1980s), Temporal CNNs (2010s)
- Transformers
  - Sequence modeling
    - NLP, speech recognition, computer vision
- Time series analysis
  - LogTrans, Informer, Pyraformer, Triformer, FEDformer

Model Type	Short Description	Date	Reference
BaselineHA	Historical Average Baseline	.	.
BaselineRW	Random Walk Baseline - guess the previous value	.	.
Regression4TS	Time-Series Regression with Lagged Variables	1949	[19]
ARIMA	AutoRegressive, Integrated, Moving-Average	1951	[9, 125]
SES	Simple Exponential Smoothing	1957	[36]
VAR	Vector AutoRegressive	1957	[96]
VARMA	Vector AutoRegressive, Moving-Average	1957	[96]
VECM	Vector Error Correction Model	1957	[96]
ES	Exponential Smoothing (Holt-Winters)	1960	[127]
SARIMA	Seasonal ARIMA	1967	[11]
SARIMAX	SARIMA, with eXogenous variables	1970	[10]
NAR/NARX	Nonlinear AutoRegressive, Exogenous	1978	[46]
Neural Network	Neural Network	1988	[108]
RNN	Recurrent Neural Network	1989	[126]
FDA/FTSA	Functional Data/Time-Series Analysis	1991	[98]
CNN	Convolutional Neural Network	1995	[60]
SVR	Support Vector Regression	1997	[85]
LSTM	Long, Short-Term Memory	1997	[35]
ELM	Extreme Learning Machine	2007	[107]
GRU	Gated Recurrent Unit	2014	[16]
Encoder-Decoder	Encoder-Decoder with Attention	2014	[17]
MGU	Minimal Gated Unit	2016	[149]
TCN	Temporal Convolutional Network	2016	[59]
GNN/GCN	Graph Neural/Convolutional Network	2016	[56]
vTRF	(Vanilla) Transformer	2017	[118]

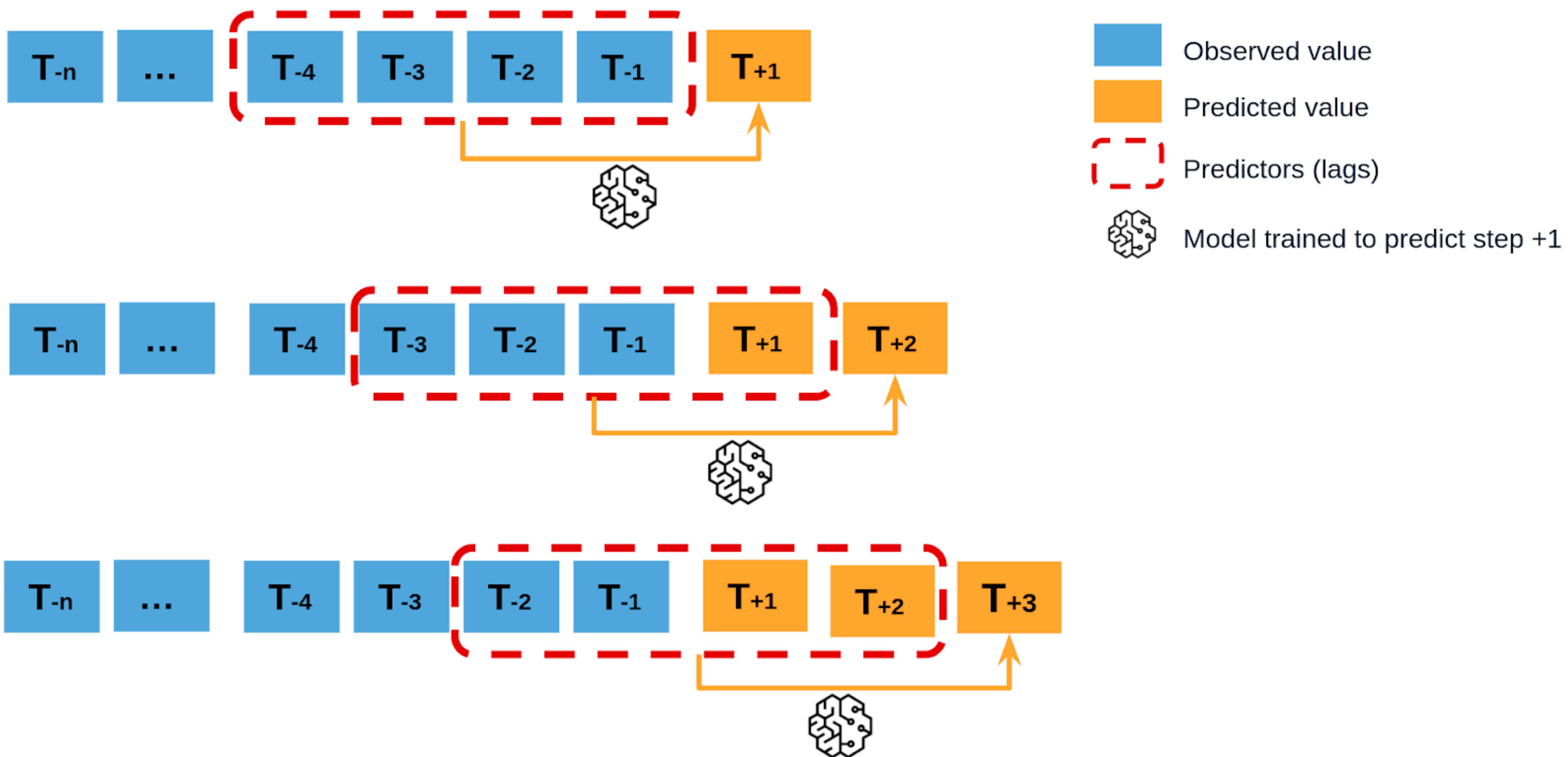
# TSF Problem Formulation

- Given:
  - Time series containing  $C$  variates, denoted as  $X = \{X_1^t, \dots, X_C^t\}_{t=1}^L$
  - $L$  is the look-back window size
  - $X_i^t$  is the value of the  $i_{th}$  variate at the  $t_{th}$  time step.
- Find:
  - Predicted values  $\hat{X} = \{\hat{X}_1^t, \dots, \hat{X}_C^t\}_{t=L+1}^{L+T}$  at the  $T$  future time steps.

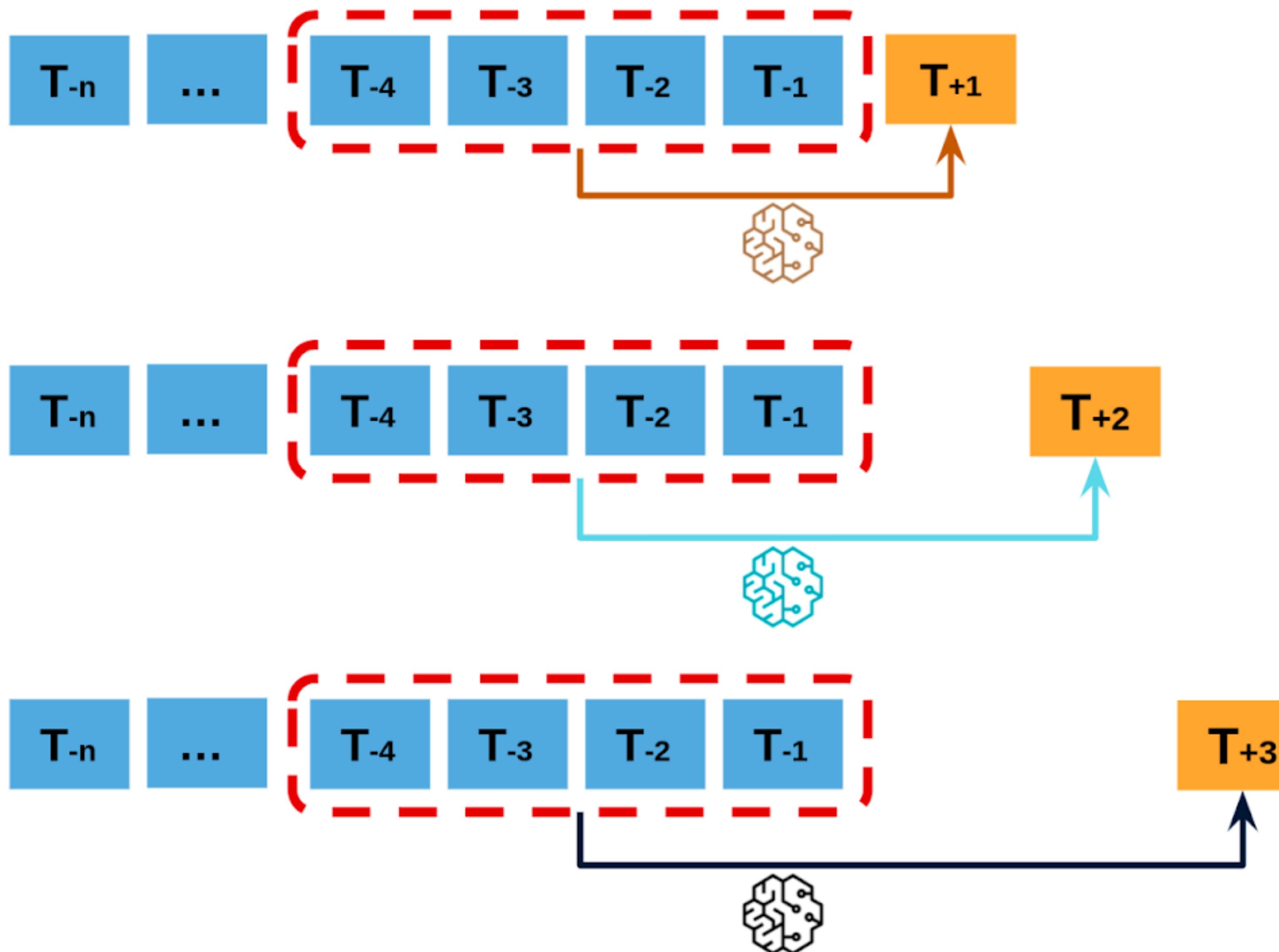
- When  $T > 1$ , can use iterated multi-step (IMS) forecasting.
- Alternatively, direct multi-step (DMS) is used.



# IMS Forecasting



# DMS Forecasting



- Observed value
- Predicted value
- Predictors (lags)
- Model trained to predict step +1
- Model trained to predict step +2
- Model trained to predict step +3

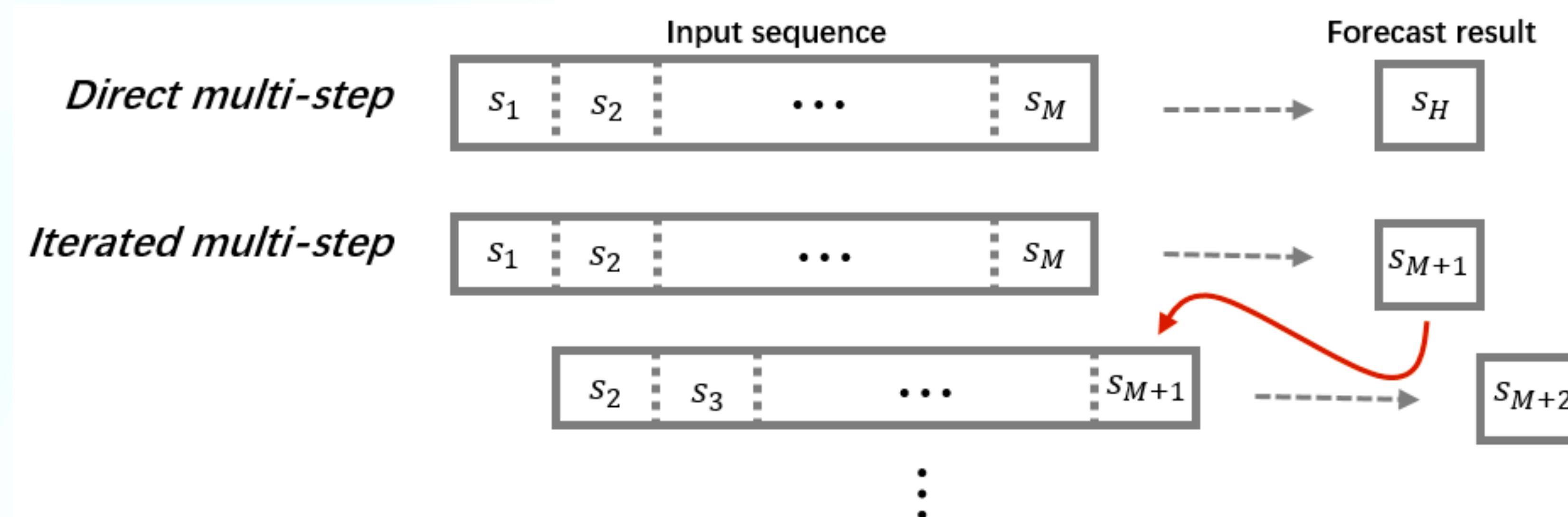
# IMS vs DMS Forecasting

## IMS

- Autoregressive estimation procedure
- Smaller variance
- Error accumulation effects
- Preferred when  $T$  is relatively small
- Needs highly accurate single-step forecaster

## DMS

- More accurate predictions
- Preferred when  $T$  is small
- Doesn't need unbiased single-step forecaster

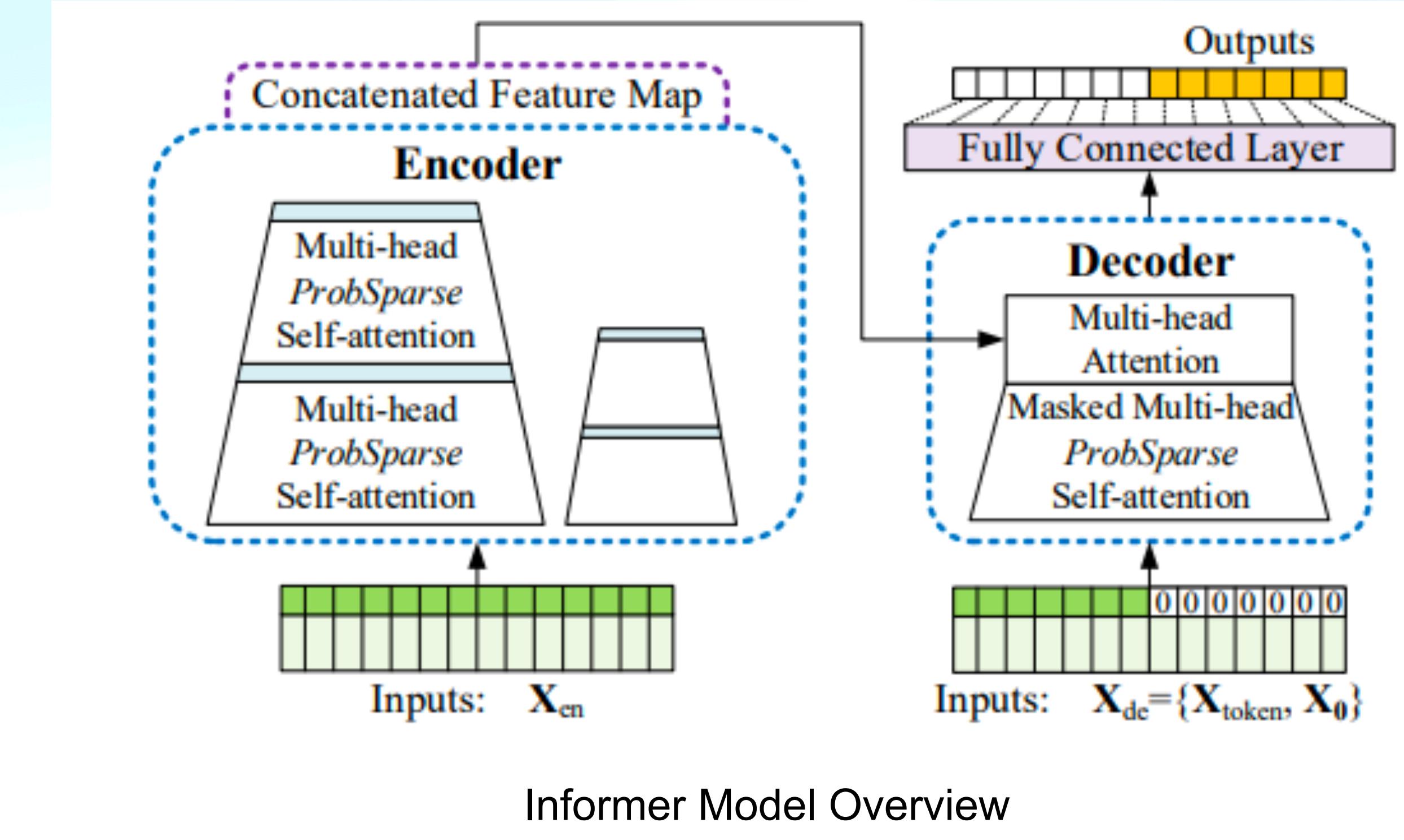


# Transformers as LTSF Solutions

- Unparalleled performance in AI tasks
- Triggered research into Transformers
- Lots of research into Transformers for LTSF tasks

## Vanilla Transformer

- time/memory complexity
  - Autoregressive decoder design
- Informer - novel architecture
  - Reduced complexity
  - DMS forecasting strategy



# Existing Transformer-based TSF Pipelines

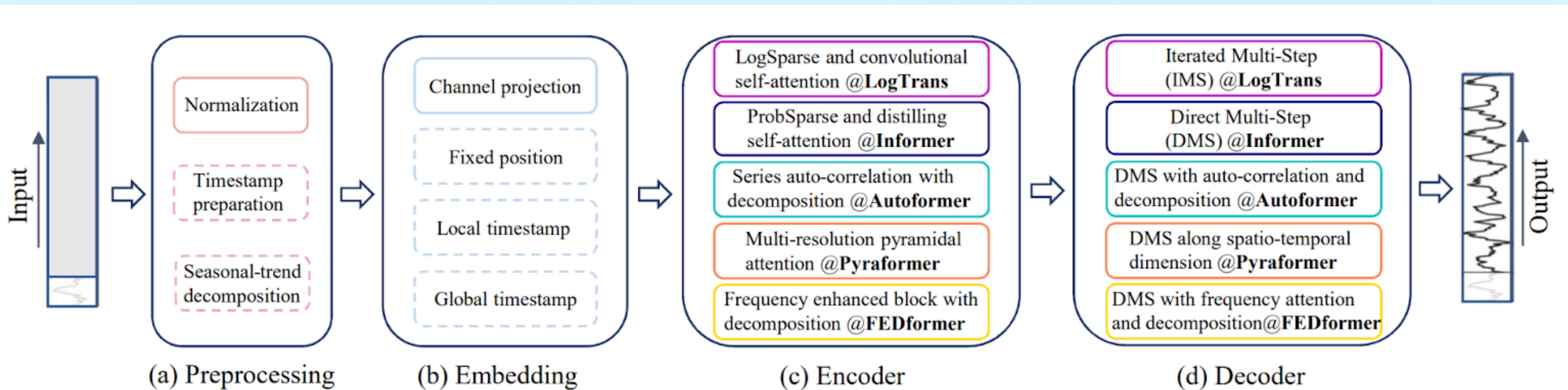


Figure 1. The pipeline of existing Transformer-based TSF solutions. In (a) and (b), the solid boxes are essential operations, and the dotted boxes are applied optionally. (c) and (d) are distinct for different methods [16, 18, 28, 30, 31].

# An Embarrassingly Simple Baseline

- LTSF Transformers compared to IMS forecasting
  - IMS - known for error accumulation effects
- Hypothesis:
  - Performance improvement mainly due to DMS forecasting in LTSF Trans.
- Baseline: Simplest DMS model
  - Temporal linear layer
  - Directly regresses historical TS for future prediction via weighted sum
  - $\hat{X}_i = WX_i, W \in \mathbb{R}^{T \times L}$

W - linear layer along temporal axis

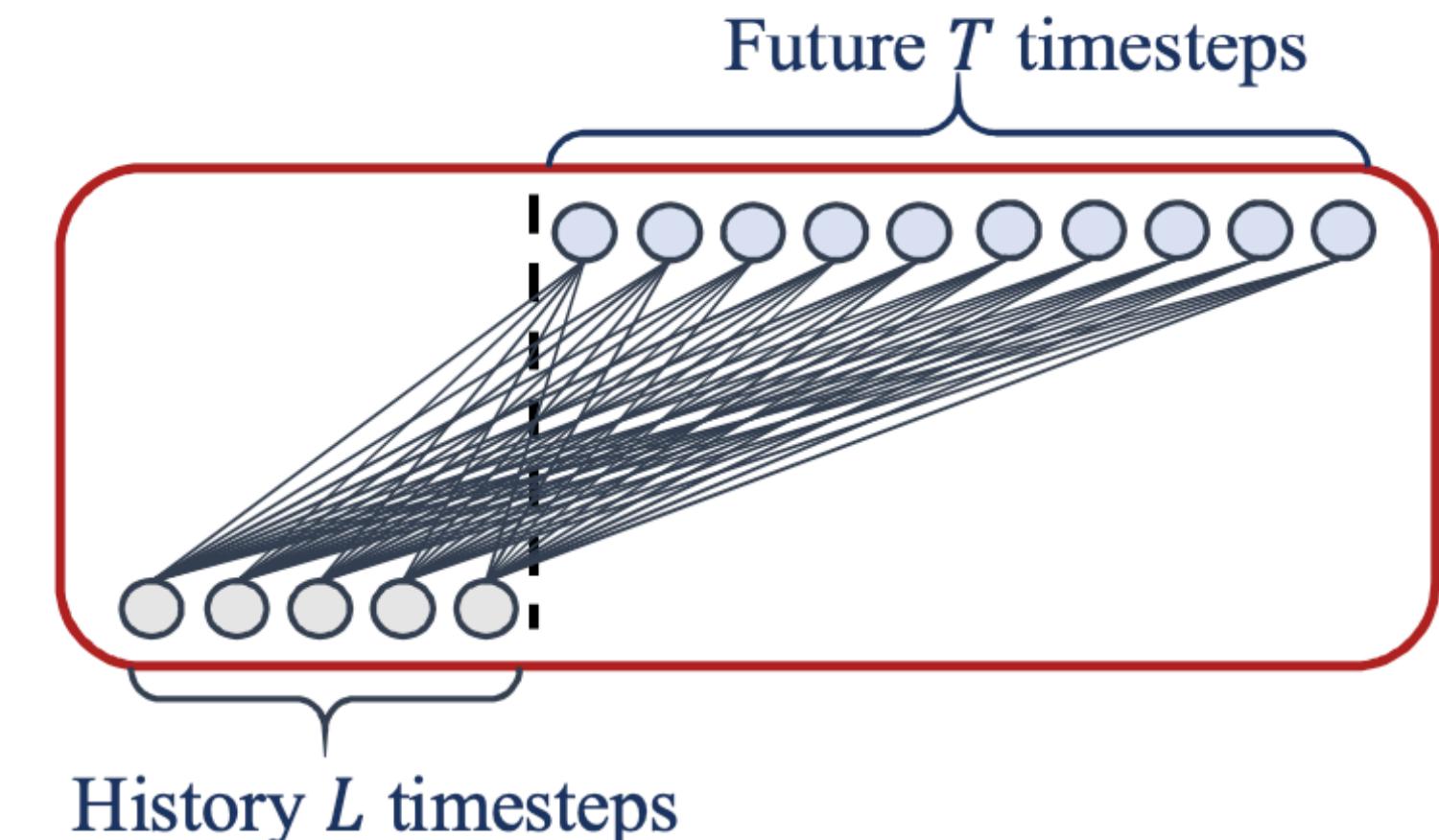
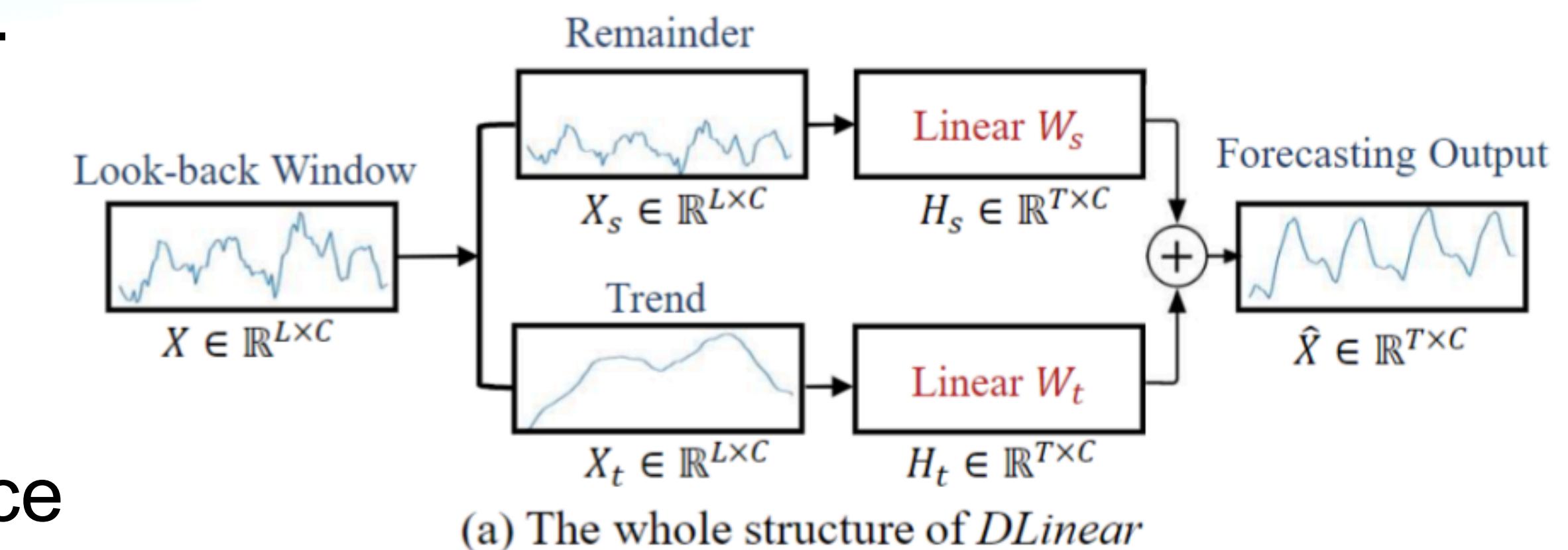


Figure 2. Illustration of the basic linear model.

# LTSF Linear (Set of Linear Models)

- Vanilla Linear
  - One layer linear model
- DLinear (Handles Trend)
  - Based on Autoformer & FEDformer decomposition scheme
  - Decomp. input data to a trend component (Mv. Avg. Kernel)
  - Two 1-layer linear layers applied to each comp.
  - Sum features to get final result



- NLinear (Handles Shift)
  - Subtracts input by the last value of the sequence
  - Puts through Linear Layer
  - Adds subtracted part back before final prediction

# Experimental Settings

- Datasets
  - Electricity Transformer Temperature (ETT)
    - ETTh1, ETTh2, ETTm1, ETTm2
  - Traffic, Electricity, Weather, ILI, Exchange-Rate
- Evaluation Metric
  - Mean Squared Error (MSE)
  - Mean Absolute Error (MAE)

## Compared Methods

- Transformers
  - FEDformer-f
    - Fourier TX variant
  - Autoformer
  - Informer
  - Pyraformer
  - LogTrans
- DMS Methods
  - Closest Repeat (Repeat)

Datasets	ETTh1&ETTh2	ETTm1 &ETTm2	Traffic	Electricity	Exchange-Rate	Weather	ILI
Variates	7	7	862	321	8	21	7
Timesteps	17,420	69,680	17,544	26,304	7,588	52,696	966
Granularity	1hour	5min	1hour	1hour	1day	10min	1week

Table 1. The statistics of the nine popular datasets for the LTSF problem.

# Comparison with Transformers - Quantitative results

Methods		IMP.	Linear*		NLinear*		DLinear*		FEDformer		Autoformer		Informer		Pyraformer*		LogTrans		Repeat*	
Metric		MSE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE								
Electricity	96	27.40%	<b>0.140</b>	<b>0.237</b>	0.141	<b>0.237</b>	<b>0.140</b>	<b>0.237</b>	<u>0.193</u>	<u>0.308</u>	0.201	0.317	0.274	0.368	0.386	0.449	0.258	0.357	1.588	0.946
	192	23.88%	<b>0.153</b>	0.250	0.154	<b>0.248</b>	<b>0.153</b>	0.249	<u>0.201</u>	<u>0.315</u>	0.222	0.334	0.296	0.386	0.386	0.443	0.266	0.368	1.595	0.950
	336	21.02%	<b>0.169</b>	0.268	0.171	<b>0.265</b>	<b>0.169</b>	0.267	<u>0.214</u>	<u>0.329</u>	0.231	0.338	0.300	0.394	0.378	0.443	0.280	0.380	1.617	0.961
	720	17.47%	<b>0.203</b>	0.301	0.210	<b>0.297</b>	<b>0.203</b>	0.301	<u>0.246</u>	<u>0.355</u>	0.254	0.361	0.373	0.439	0.376	0.445	0.283	0.376	1.647	0.975
Exchange	96	45.27%	0.082	0.207	0.089	0.208	<b>0.081</b>	0.203	<u>0.148</u>	<u>0.278</u>	0.197	0.323	0.847	0.752	0.376	1.105	0.968	0.812	<b>0.081</b>	<b>0.196</b>
	192	42.06%	0.167	0.304	0.180	0.300	<b>0.157</b>	0.293	<u>0.271</u>	<u>0.380</u>	0.300	0.369	1.204	0.895	1.748	1.151	1.040	0.851	0.167	<b>0.289</b>
	336	33.69%	0.328	0.432	0.331	0.415	<b>0.305</b>	0.414	<u>0.460</u>	<u>0.500</u>	0.509	0.524	1.672	1.036	1.874	1.172	1.659	1.081	<b>0.305</b>	<b>0.396</b>
	720	46.19%	0.964	0.750	1.033	0.780	<b>0.643</b>	<b>0.601</b>	<u>1.195</u>	<u>0.841</u>	1.447	0.941	2.478	1.310	1.943	1.206	1.941	1.127	0.823	0.681
Traffic	96	30.15%	<b>0.410</b>	0.282	<b>0.410</b>	<b>0.279</b>	<b>0.410</b>	0.282	<u>0.587</u>	<u>0.366</u>	0.613	0.388	0.719	0.391	2.085	0.468	0.684	0.384	2.723	1.079
	192	29.96%	<b>0.423</b>	0.287	<b>0.423</b>	<b>0.284</b>	<b>0.423</b>	0.287	<u>0.604</u>	<u>0.373</u>	0.616	0.382	0.696	0.379	0.867	0.467	0.685	0.390	2.756	1.087
	336	29.95%	0.436	0.295	<b>0.435</b>	<b>0.290</b>	0.436	0.296	<u>0.621</u>	0.383	0.622	<u>0.337</u>	0.777	0.420	0.869	0.469	0.734	0.408	2.791	1.095
	720	25.87%	0.466	0.315	<b>0.464</b>	<b>0.307</b>	0.466	0.315	<u>0.626</u>	<u>0.382</u>	0.660	0.408	0.864	0.472	0.881	0.473	0.717	0.396	2.811	1.097
Weather	96	18.89%	<b>0.176</b>	0.236	0.182	<b>0.232</b>	<b>0.176</b>	0.237	<u>0.217</u>	<u>0.296</u>	0.266	0.336	0.300	0.384	0.896	0.556	0.458	0.490	0.259	0.254
	192	21.01%	<b>0.218</b>	0.276	0.225	<b>0.269</b>	0.220	0.282	<u>0.276</u>	<u>0.336</u>	0.307	0.367	0.598	0.544	0.622	0.624	0.658	0.589	0.309	0.292
	336	22.71%	<b>0.262</b>	0.312	0.271	<b>0.301</b>	0.265	0.319	<u>0.339</u>	<u>0.380</u>	0.359	0.395	0.578	0.523	0.739	0.753	0.797	0.652	0.377	0.338
	720	19.85%	0.326	0.365	0.338	<b>0.348</b>	<b>0.323</b>	0.362	<u>0.403</u>	<u>0.428</u>	0.419	0.428	1.059	0.741	1.004	0.934	0.869	0.675	0.465	0.394
ILI	24	47.86%	1.947	0.985	<b>1.683</b>	<b>0.858</b>	2.215	1.081	<u>3.228</u>	<u>1.260</u>	3.483	1.287	5.764	1.677	1.420	2.012	4.480	1.444	6.587	1.701
	36	36.43%	2.182	1.036	<b>1.703</b>	<b>0.859</b>	1.963	0.963	<u>2.679</u>	<u>1.080</u>	3.103	1.148	4.755	1.467	7.394	2.031	4.799	1.467	7.130	1.884
	48	34.43%	2.256	1.060	<b>1.719</b>	<b>0.884</b>	2.130	1.024	<u>2.622</u>	<u>1.078</u>	2.669	1.085	4.763	1.469	7.551	2.057	4.800	1.468	6.575	1.798
	60	34.33%	2.390	1.104	<b>1.819</b>	<b>0.917</b>	2.368	1.096	<u>2.857</u>	<u>1.157</u>	<u>2.770</u>	<u>1.125</u>	5.264	1.564	7.662	2.100	5.278	1.560	5.893	1.677
ETTh1	96	0.80%	0.375	0.397	<b>0.374</b>	<b>0.394</b>	0.375	0.399	<u>0.376</u>	<u>0.419</u>	0.449	0.459	0.865	0.713	0.664	0.612	0.878	0.740	1.295	0.713
	192	3.57%	0.418	0.429	0.408	<b>0.415</b>	<b>0.405</b>	0.416	<u>0.420</u>	<u>0.448</u>	0.500	0.482	1.008	0.792	0.790	0.681	1.037	0.824	1.325	0.733
	336	6.54%	0.479	0.476	<b>0.429</b>	<b>0.427</b>	0.439	0.443	<u>0.459</u>	<u>0.465</u>	0.521	0.496	1.107	0.809	0.891	0.738	1.238	0.932	1.323	0.744
	720	13.04%	0.624	0.592	<b>0.440</b>	<b>0.453</b>	0.472	0.490	<u>0.506</u>	<u>0.507</u>	0.514	0.512	1.181	0.865	0.963	0.782	1.135	0.852	1.339	0.756
ETTh2	96	19.94%	0.288	0.352	<b>0.277</b>	<b>0.338</b>	0.289	0.353	<u>0.346</u>	<u>0.388</u>	0.358	0.397	3.755	1.525	0.645	0.597	2.116	1.197	0.432	0.422
	192	19.81%	0.377	0.413	<b>0.344</b>	<b>0.381</b>	0.383	0.418	<u>0.429</u>	<u>0.439</u>	0.456	0.452	5.602	1.931	0.788	0.683	4.315	1.635	0.534	0.473
	336	25.93%	0.452	0.461	<b>0.357</b>	<b>0.400</b>	0.448	0.465	<u>0.496</u>	<u>0.487</u>	0.482	0.486	4.721	1.835	0.907	0.747	1.124	1.604	0.591	0.508
	720	14.25%	0.698	0.595	<b>0.394</b>	<b>0.436</b>	0.605	0.551	<u>0.463</u>	<u>0.474</u>	0.515	0.511	3.647	1.625	0.963	0.783	3.188	1.540	0.588	0.517
ETTm1	96	21.10%	0.308	0.352	0.306	0.348	<b>0.299</b>	<b>0.343</b>	<u>0.379</u>	<u>0.419</u>	0.505	0.475	0.672	0.571	0.543	0.510	0.600	0.546	1.214	0.665
	192	21.36%	0.340	0.369	0.349	0.375	<b>0.335</b>	<b>0.365</b>	<u>0.426</u>	<u>0.441</u>	0.553	0.496	0.795	0.669	0.55					

# Comparison with Transformers - Qualitative results

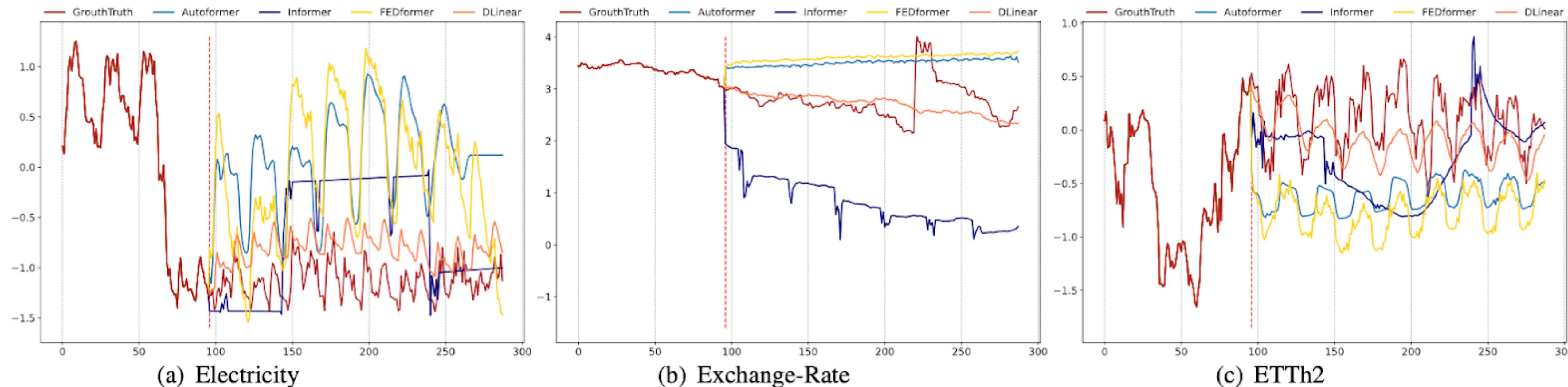


Figure 3. Illustration of the long-term forecasting output (Y-axis) of five models with an input length  $L=96$  and output length  $T=192$  (X-axis) on Electricity, Exchange-Rate, and ETTh2, respectively.

- Fail to capture scale & bias of future data (Electricity, ETTh2)
- Fail to predict trends on aperiodic data (Exchange-Rate)

# Can LTSF-Transformers extract temporal relations well from longer input sequences?

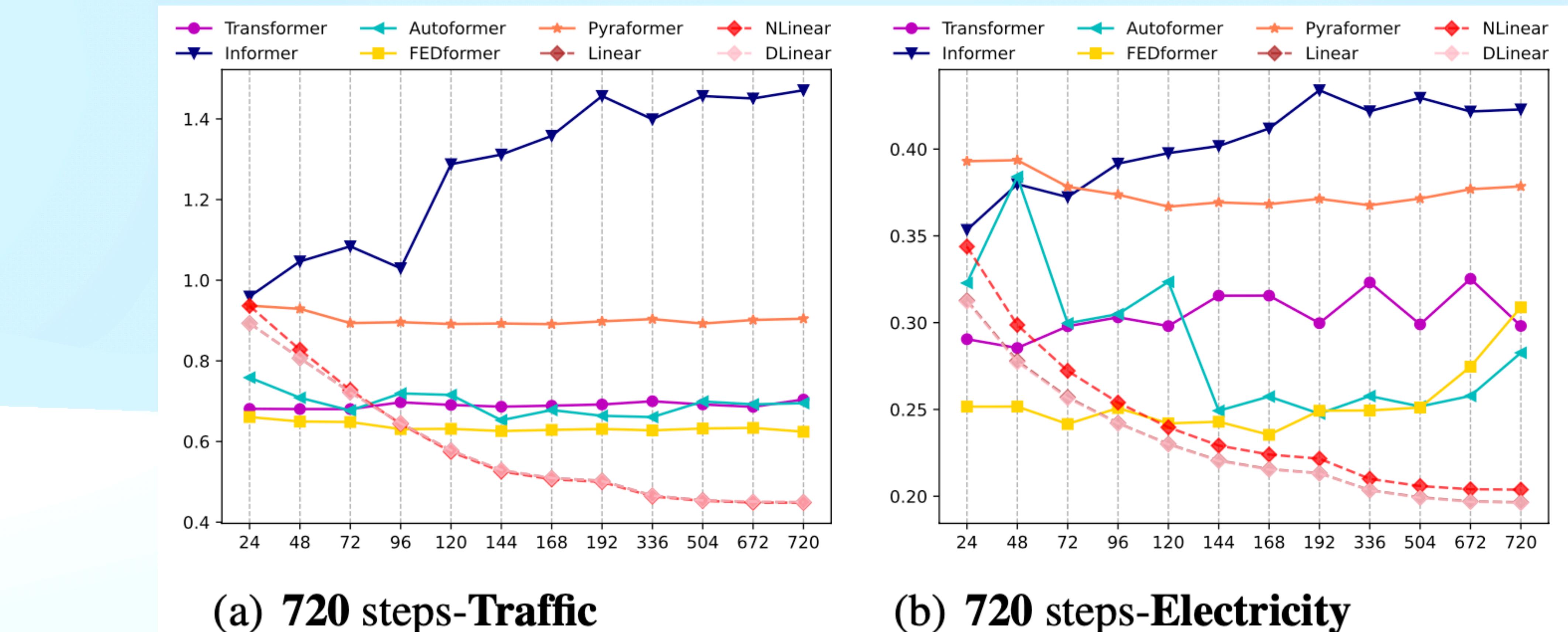


Figure 4. The MSE results (Y-axis) of models with different look-back window sizes (X-axis) of long-term forecasting ( $T=720$ ) on the Traffic and Electricity datasets.

# What can be learned for long-term forecasting?

- Short term forecasting
  - Temporal dynamics in look back window
- Long term forecasting
  - Hypothesis: trend & periodicity
  - Farther forecasting horizon, less impact look back window has

Methods	FEDformer		Autoformer	
Input	<i>Close</i>	<i>Far</i>	<i>Close</i>	<i>Far</i>
Electricity	0.251	0.265	0.255	0.287
Traffic	0.631	0.645	0.677	0.675

Table 3. Comparison of different input sequences under the MSE metric to explore what LTSF-Transformers depend on. If the input is *Close*, we use the  $96_{th}, \dots, 191_{th}$  time steps as the input sequence. If the input is *Far*, we use the  $0_{th}, \dots, 95_{th}$  time steps. Both of them forecast the  $192_{th}, \dots, (192 + 720)_{th}$  time steps.

# Are the self-attention scheme effective for LTSF?

- Gradually transformed Informer to Linear
  - Replace self-attention layers w/ linear layer (Att.-Linear)
  - Discard other auxiliary designs (e.g. FFN)
    - Embedding and linear layers remaining (Embed + Linear)
  - Simplify model to one linear layer
- Informer Performance Grows w/ Simplification

Methods	Informer	Att.-Linear	Embed + Linear	Linear
Exchange	96	0.847	1.003	0.173
	192	1.204	0.979	0.443
	336	1.672	1.498	1.288
	720	2.478	2.102	2.026
ETTh1	96	0.865	0.613	0.454
	192	1.008	0.759	0.686
	336	1.107	0.921	0.821
	720	1.181	0.902	1.051

Table 4. The MSE comparisons of gradually transforming Informer to a Linear from the left to right columns. *Att.-Linear* is a structure that replaces each attention layer with a linear layer. *Embed + Linear* is to drop other designs and only keeps embedding layers and a linear layer. The look-back window size is 96.

# Can existing LTSF-Transformers preserve temporal order well?

Methods		Linear			FEDformer			Autoformer			Informer		
Predict Length		Ori.	Shuf.	Half-Ex.	Ori.	Shuf.	Half-Ex.	Ori.	Shuf.	Half-Ex.	Ori.	Shuf.	Half-Ex.
Exchange	96	0.080	0.133	0.169	0.161	0.160	0.162	0.152	0.158	0.160	0.952	1.004	0.959
	192	0.162	0.208	0.243	0.274	0.275	0.275	0.278	0.271	0.277	1.012	1.023	1.014
	336	0.286	0.320	0.345	0.439	0.439	0.439	0.435	0.430	0.435	1.177	1.181	1.177
	720	0.806	0.819	0.836	1.122	1.122	1.122	1.113	1.113	1.113	1.198	1.210	1.196
Average Drop		N/A	27.26%	46.81%	N/A	-0.09%	0.20%	N/A	0.09%	1.12%	N/A	-0.12%	-0.18%
ETTh1	96	0.395	0.824	0.431	0.376	0.753	0.405	0.455	0.838	0.458	0.974	0.971	0.971
	192	0.447	0.824	0.471	0.419	0.730	0.436	0.486	0.774	0.491	1.233	1.232	1.231
	336	0.490	0.825	0.505	0.447	0.736	0.453	0.496	0.752	0.497	1.693	1.693	1.691
	720	0.520	0.846	0.528	0.468	0.720	0.470	0.525	0.696	0.524	2.720	2.716	2.715
Average Drop		N/A	81.06%	4.78%	N/A	73.28%	3.44%	N/A	56.91%	0.46%	N/A	1.98%	0.18%

Table 5. The MSE comparisons of models when shuffling the raw input sequence. *Shuf.* randomly shuffles the input sequence. *Half-EX.* randomly exchanges the first half of the input sequences with the second half. Average Drop is the average performance drop under all forecasting lengths after shuffling. All results are the average test MSE of five runs.

- Transformer performance not fluctuating after shuffling
- LTSF-Linear performance damaged significantly

# How effective are different embedding strategies?

Methods	Embedding	Traffic			
		96	192	336	720
FEDformer	All	0.597	0.606	0.627	0.649
	wo/Pos.	<b>0.587</b>	<b>0.604</b>	<b>0.621</b>	<b>0.626</b>
	wo/Temp.	0.613	0.623	0.650	0.677
	wo/Pos.-Temp.	0.613	0.622	0.648	0.663
Autoformer	All	0.629	0.647	0.676	<b>0.638</b>
	wo/Pos.	<b>0.613</b>	<b>0.616</b>	<b>0.622</b>	0.660
	wo/Temp.	0.681	0.665	0.908	0.769
	wo/Pos.-Temp.	0.672	0.811	1.133	1.300
Informer	All	<b>0.719</b>	<b>0.696</b>	<b>0.777</b>	<b>0.864</b>
	wo/Pos.	1.035	1.186	1.307	1.472
	wo/Temp.	0.754	0.780	0.903	1.259
	wo/Pos.-Temp.	1.038	1.351	1.491	1.512

Table 6. The MSE comparisons of different embedding strategies on Transformer-based methods with look-back window size 96 and forecasting lengths {96, 192, 336, 720}.

# Is training data size a limiting factor for LTSF- Transformers?

Unlike CV/NLP -

- TSF is performed on collected time series
- Difficult to scale up the training data size
- Size of the training data significant impact on performance

Experiments on Traffic

- Model trained on full dataset
  - 17,544\*0.7 hours
  - Namd Ori.
- Model trained on shortened dataset
  - 8,760 hours, i.e., 1 year
  - Named Short

Result: Prediction errors w/ reduced training data

- Lower in most cases
- Why?
  - whole-year data more clear temporal features
    - Opp. to longer but incomplete data size

Methods	FEDformer		Autoformer		
	Dataset	Ori.	Short	Ori.	Short
	96	0.587	<b>0.568</b>	0.613	<b>0.594</b>
	192	0.604	<b>0.584</b>	<b>0.616</b>	0.621
	336	0.621	<b>0.601</b>	0.622	<b>0.621</b>
	720	0.626	<b>0.608</b>	0.660	<b>0.650</b>

Table 7. The MSE comparison of two training data sizes

# Is efficiency really a top-level priority?

Existing LTSF- Transformers claim

- Vanilla Trans. quad. complexity too high for LTSF
- Prove able to improve theoretical complexity
- Unclear if
  - Actual cost on devices improved
  - Memory issue unacceptable for today's GPU

Compared average practical efficiencies with 5 runs

- Most Transformer variants
  - incur similar or worse in practice
    - inference time
    - parameters
- These follow-ups
  - add more design elements
  - increase practical costs
- vanilla Trans. Memory cost
  - practically acceptable
  - even for output length  $L = 720$ ,
    - weakens importance of developing a mem. efficient Transformer

Method	MACs	Parameter	Time	Memory
DLinear	<b>0.04G</b>	<b>139.7K</b>	<b>0.4ms</b>	<b>687MiB</b>
Transformer <sup>×</sup>	4.03G	13.61M	26.8ms	6091MiB
Informer	3.93G	14.39M	49.3ms	3869MiB
Autoformer	4.41G	14.91M	164.1ms	7607MiB
Pyraformer	0.80G	241.4M*	3.4ms	7017MiB
FEDformer	4.41G	20.68M	40.5ms	4143MiB

- <sup>×</sup> is modified into the same one-step decoder, which is implemented in the source code from Autoformer.

- \* 236.7M parameters of Pyraformer come from its linear decoder.

Table 8. Comparison of practical efficiency of LTSF-Transformers under  $L=96$  and  $T=720$  on the Electricity. MACs are the number of multiply-accumulate operations. We use Dlinear for comparison since it has the double cost in *LTSF-Linear*. The inference time averages 5 runs.

# Comparison of Univariate Forecasting

Methods		Linear		NLinear		DLinear		FEDformer-f		FEDformer-w		Autoformer		Informer		LogTrans	
Metric		MSE	MAE	MSE	MSE	MAE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
<i>ETTh1</i>	96	0.189	0.359	<b>0.053</b>	<b>0.177</b>	0.056	0.180	0.079	0.215	0.080	0.214	<u>0.071</u>	<u>0.206</u>	0.193	0.377	0.283	0.468
	192	0.078	0.212	<b>0.069</b>	<b>0.204</b>	0.071	0.204	<u>0.104</u>	<u>0.245</u>	0.105	0.256	0.114	0.262	0.217	0.395	0.234	0.409
	336	0.091	0.237	<b>0.081</b>	<b>0.226</b>	0.098	0.244	0.119	0.270	0.120	0.269	<u>0.107</u>	<u>0.258</u>	0.202	0.381	0.386	0.546
	720	0.172	0.340	<b>0.080</b>	<b>0.226</b>	0.189	0.359	0.142	0.299	0.127	0.280	<u>0.126</u>	<u>0.283</u>	0.183	0.355	0.475	0.629
<i>ETTh2</i>	96	0.133	0.283	<b>0.129</b>	<b>0.278</b>	0.131	0.279	<u>0.128</u>	<u>0.271</u>	0.156	0.306	0.153	0.306	0.213	0.373	0.217	0.379
	192	0.176	0.330	<b>0.169</b>	<b>0.324</b>	0.176	0.329	<u>0.185</u>	<u>0.330</u>	0.238	0.380	0.204	0.351	0.227	0.387	0.281	0.429
	336	0.213	0.371	<b>0.194</b>	<b>0.355</b>	0.209	0.367	<u>0.231</u>	<u>0.378</u>	0.271	0.412	0.246	0.389	0.242	0.401	0.293	0.437
	720	0.292	0.440	<b>0.225</b>	<b>0.381</b>	0.276	0.426	0.278	0.420	0.288	0.438	<u>0.268</u>	<u>0.409</u>	0.291	0.439	0.218	0.387
<i>ETTm1</i>	96	0.028	0.125	<b>0.026</b>	<b>0.122</b>	0.028	0.123	<u>0.033</u>	<u>0.140</u>	0.036	0.149	0.056	0.183	0.109	0.277	0.049	0.171
	192	0.043	0.154	<b>0.039</b>	<b>0.149</b>	0.045	0.156	<u>0.058</u>	<u>0.186</u>	0.069	0.206	0.081	0.216	0.151	0.310	0.157	0.317
	336	0.059	0.180	<b>0.052</b>	<b>0.172</b>	0.061	0.182	0.084	0.231	<u>0.071</u>	<u>0.209</u>	0.076	0.218	0.427	0.591	0.289	0.459
	720	0.080	0.211	<b>0.073</b>	<b>0.207</b>	0.080	0.210	<u>0.102</u>	<u>0.250</u>	<u>0.105</u>	<u>0.248</u>	0.110	0.267	0.438	0.586	0.430	0.579
<i>ETTm2</i>	96	0.066	0.189	<b>0.063</b>	<b>0.182</b>	0.063	0.183	0.067	0.198	<u>0.063</u>	<u>0.189</u>	0.065	0.189	0.088	0.225	0.075	0.208
	192	0.094	0.230	<b>0.090</b>	<b>0.223</b>	0.092	0.227	<u>0.102</u>	<u>0.245</u>	0.110	0.252	0.118	0.256	0.132	0.283	0.129	0.275
	336	0.120	0.263	<b>0.117</b>	<b>0.259</b>	0.119	0.261	<u>0.130</u>	<u>0.279</u>	0.147	0.301	0.154	0.305	0.180	0.336	0.154	0.302
	720	0.175	0.320	<b>0.170</b>	<b>0.318</b>	0.175	0.320	0.178	<u>0.325</u>	0.219	0.368	0.182	0.335	0.300	0.435	0.160	0.321

Table 9. Univariate long sequence time-series forecasting results on ETT full benchmark. The **best results** are highlighted in **bold** and the **best results of Transformers** are highlighted with a underline.

# Distribution of ETTh1, ETTh2, Electricity, and ILI dataset

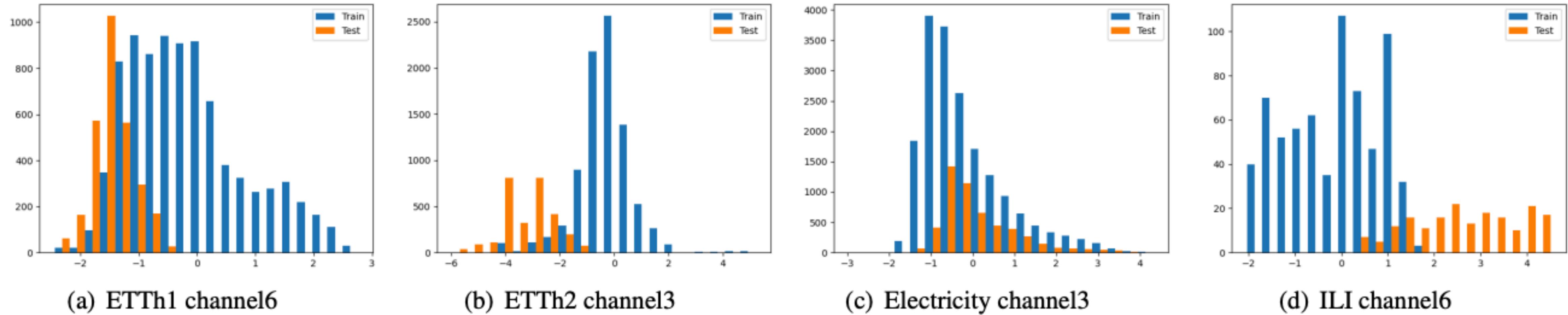


Figure 5. Distribution of ETTh1, ETTh2, Electricity, and ILI dataset. A clear distribution shift between training and testing data can be observed in ETTh1, ETTh2, and ILI.

# Visualizing Weights of LSTF-Linear on Different Datasets

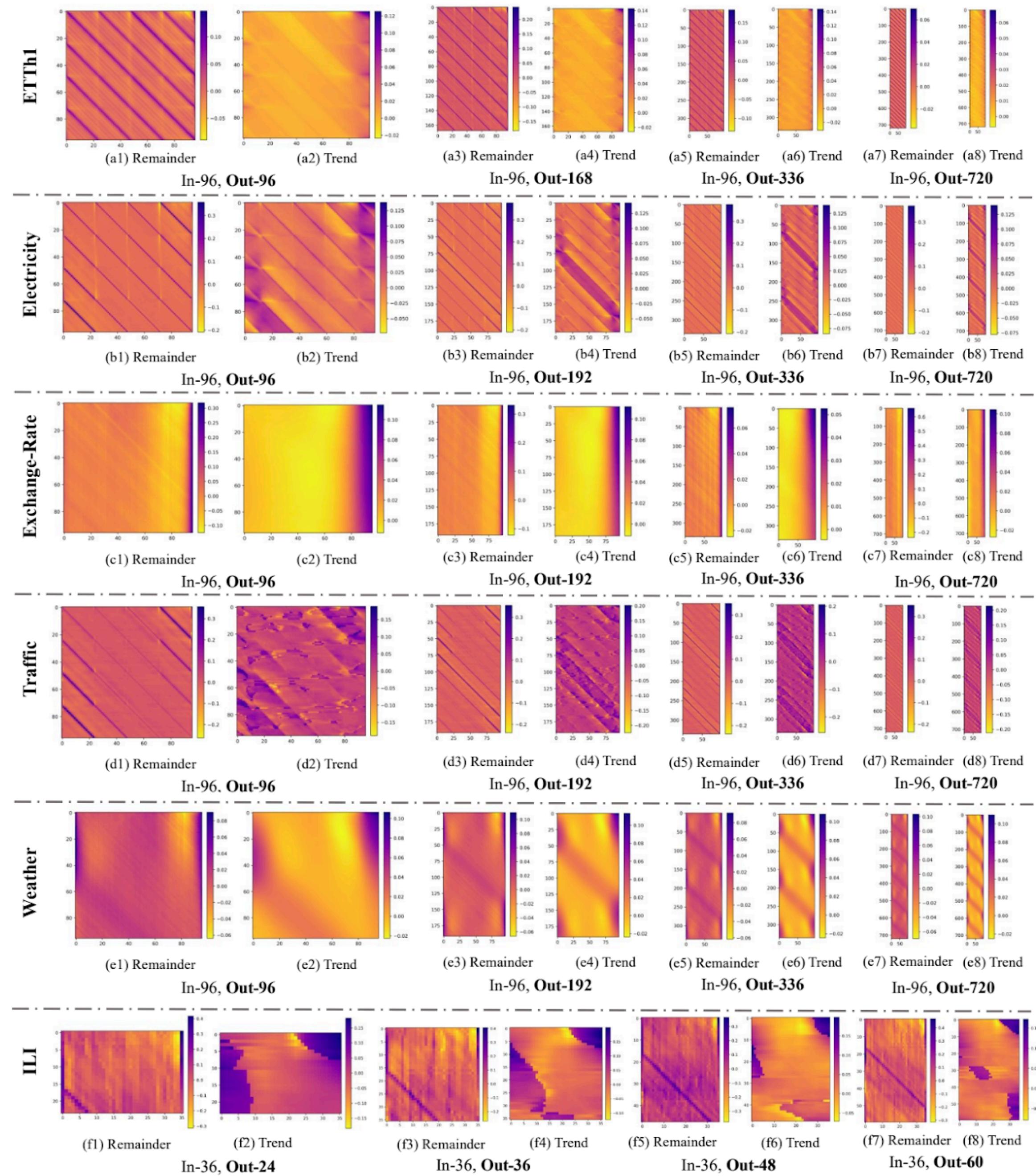


Figure 7. Visualization of the weights( $T^*L$ ) of *LTSF-Linear* on several benchmarks. Models are trained with a look-back window  $L$  (X-axis) and different forecasting time steps  $T$  (Y-axis). We show weights in the remainder and trend layer.

# The MSE results (Y-axis) of models with different look-back window sizes (X-axis) of the long-term forecasting

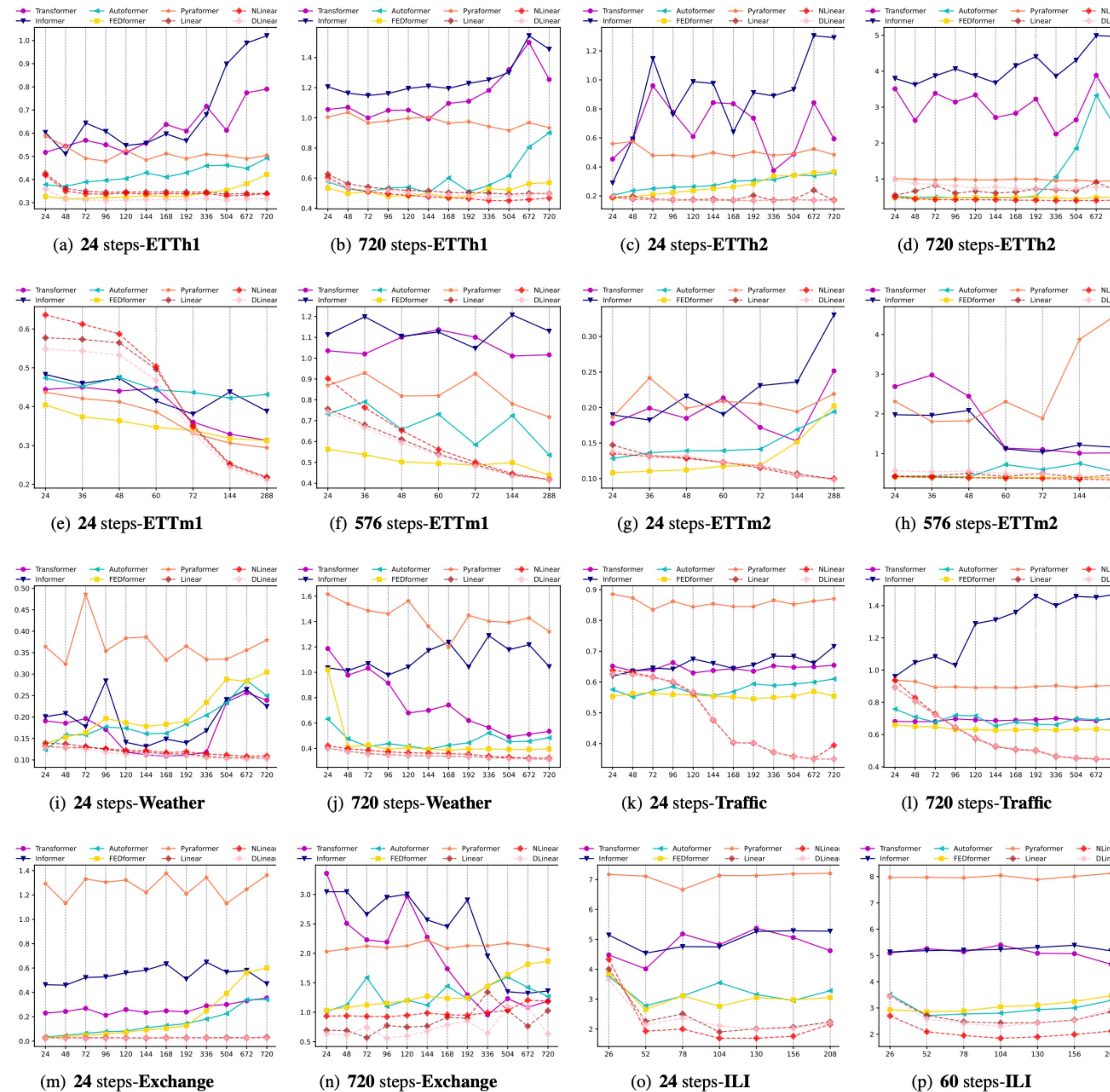


Figure 6. The MSE results (Y-axis) of models with different look-back window sizes (X-axis) of the long-term forecasting (e.g., 720-time steps) and the short-term forecasting (e.g., 24 time steps) on different benchmarks.