

# APT: Affine Prototype-Timestamp For Time Series Forecasting Under Distribution Shift

Yujie Li,<sup>1,2</sup> Zezhi Shao,<sup>1</sup> Chengqing Yu,<sup>1,2</sup> Yisong Fu,<sup>1,2</sup> Tao Sun,<sup>1</sup> Yongjun Xu,<sup>1,2</sup> Fei Wang<sup>1,2\*</sup>

<sup>1</sup>State Key Laboratory of AI Safety, Institute of Computing Technology, Chinese Academy of Sciences

<sup>2</sup>University of Chinese Academy of Sciences

{liyujie23s, shaozezhi, yuchengqing22b, fuyisong24s, suntao, xyj, wangfei}@ict.ac.cn

## Abstract

Time series forecasting under distribution shift remains challenging, as existing deep learning models often rely on local statistical normalization (e.g., mean and variance) that fails to capture global distribution shift. Methods like RevIN and its variants attempt to decouple distribution and pattern but still struggle with missing values, noisy observations, and invalid channel-wise affine transformation. To address these limitations, we propose Affine Prototype-Timestamp (APT), a lightweight and flexible plug-in module that injects global distribution features into the normalization-forecasting pipeline. By leveraging timestamp-conditioned prototype learning, APT dynamically generates affine parameters that modulate both input and output series, enabling the backbone to learn from self-supervised, distribution-aware clustered instances. APT is compatible with arbitrary forecasting backbones and normalization strategies while introducing minimal computational overhead. Extensive experiments across six benchmark datasets and multiple backbone-normalization combinations demonstrate that APT significantly improves forecasting performance under distribution shift.

**Code** — <https://github.com/blisky-li/APT>

## Introduction

Time series reflect the artificial or natural regularities of complex dynamic systems across transportation (Li et al. 2024; Shao et al. 2022), health (Ferté et al. 2024) and weather (Yu et al. 2025). However, external factors commonly induce distributional shift and non-stationarity, making forecasting models struggle to effectively capture patterns under changing statistical properties.

Reversible Instance Normalization (RevIN) (Kim et al. 2021) introduces a two-stage paradigm for mitigating temporal distribution shift: instance normalization that removes instance-specific distributions and affine transformation that attempts to restore channel-wise distributional features. This decoupling of time-varying distributions from learnable temporal patterns has laid the foundation for many modern forecasting models, and subsequent advancements—including DishTS (Fan et al. 2023), SAN (Liu et al.

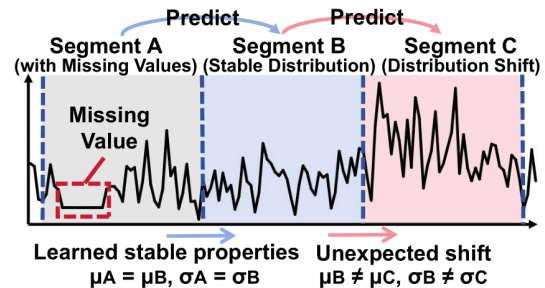


Figure 1: Local statistical normalization fails to handle distribution shift on ECL. The model retains outdated statistics when predicting from flawed Segment A to B but faces unseen shifts in Segment B and C. Shifts across all three segments exceed prior inter and intra shift issues, posing a broader global shift challenge for APT.

2023), SIN (Han, Ye, and Zhan 2024), and FAN (Ye et al. 2024)—have further refined this framework by modeling future distributions adaptively to address intra-instance shift.

Despite their success, RevIN-like methods still face inherent limitations. First, Local statistics such as mean and variance are sensitive to missing values and noise, which are common in time series. Second, their instance-specific normalization fails to capture global shifts, as shown in Figure 1, where features learned from Segments A and B collapse in Segment C under abrupt change.

Moreover, the channel-wise static affine transformation (Bebis et al. 1999) in RevIN yields only limited improvements as shown in Table 1, and is entirely omitted in follow-up works such as DishTS and SAN. This suggests static transformations fail to address inter-channel distribution variance and are ineffective against distribution shift.

To overcome the limitations of local statistical normalization, we propose **A**ffine **P**rototype-**T**imestamps (APT), a lightweight and model-agnostic plug-in. APT replaces static affine transformations in standard normalization with dynamically generated parameters conditioned on timestamps, inspired by vision style transfer (Huang and Belongie 2017). This enables the forecasting pipeline to access global temporal semantics, which serve as an underlying schedule guiding the behavior of the time series system, thereby enhancing robustness and adaptability under distribution shift.

\*Corresponding authors.

Methods	CATS		Informer+RevIN		iTransformer		SparseTSF		Avg. $\Delta$									
Affine (RevIN)	True	False	True	False	True	False	True	False	True - False									
Metric	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE								
ECL	<b>0.260</b>	<b>0.165</b>	0.261	<b>0.165</b>	<b>0.307</b>	<b>0.209</b>	0.308	0.212	<b>0.264</b>	<b>0.169</b>	0.264	0.170	0.266	0.173	<b>0.265</b>	<b>0.171</b>	0.000	↓ 0.001
Weather	<b>0.274</b>	<b>0.236</b>	<b>0.274</b>	0.241	0.311	0.307	<b>0.310</b>	<b>0.306</b>	0.284	0.251	<b>0.282</b>	<b>0.248</b>	<b>0.294</b>	<b>0.267</b>	0.297	0.269	0.000	↓ 0.001
ETTh1	0.447	0.456	<b>0.441</b>	<b>0.451</b>	<b>0.582</b>	0.681	0.584	<b>0.679</b>	<b>0.462</b>	<b>0.471</b>	0.464	0.475	0.440	0.445	<b>0.430</b>	<b>0.437</b>	↑ 0.003	↑ 0.003
Exchange	0.437	0.352	<b>0.436</b>	<b>0.349</b>	<b>0.677</b>	0.708	<b>0.644</b>	0.727	0.493	0.423	<b>0.475</b>	<b>0.398</b>	0.499	0.434	<b>0.492</b>	<b>0.428</b>	↑ 0.015	↑ 0.004

Table 1: The affine transformations in RevIN do not provide performance gains, Metric  $\uparrow$  means worse performance.

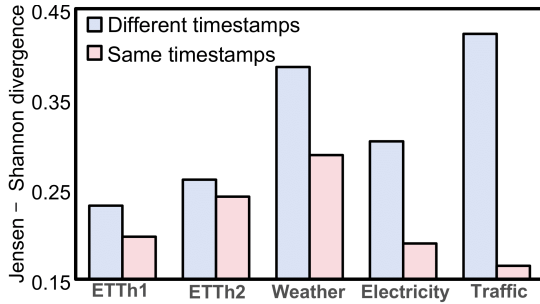


Figure 2: JS divergence of subseries with the same timestamp label in benchmark datasets

Easily accessible timestamps are widely regarded as containing global temporal semantics (Li et al. 2025), and our empirical analysis in Figure 2 confirms subseries sharing similar timestamp labels such as "Time in Day" and "Day in Week" tend to exhibit distributional similarity. APT exploits this regularity by embedding timestamps into a shared latent space to enable conditioned affine transformations.

However, the sparsity of fine-grained timestamps combinations can cause insufficient training or poor generalization to unseen cases. Inspired by few-shot learning (Li et al. 2022), APT employs prototype learning, replacing raw timestamp embeddings with nearest-neighbor prototypes to yield more robust representations. To further promote diversity and balanced prototype usage, we introduce orthogonality and load-balancing losses.

Rather than learning static affine parameters, APT utilizes global temporal semantics in timestamps to adaptively assign distinct affine parameters to different subseries. To achieve this, it first encodes timestamps and obtains robust representations through prototype matching, then converts prototype embeddings into low-dimensional affine parameters via MLPs. Crucially, this process employs self-supervised learning with the backbone and normalization frozen, leveraging orthogonal loss for embedding diversity, load-balancing loss for prototype usage, and extra affine regularization loss to ensure affine parameter convergence.

In summary, the central function of APT is to supply forecasting backbones with learnable and global distribution features suppressed by normalization or disrupted by distribution shift; all components and optimization of APT are explicitly designed to serve this purpose.

Our contributions are as follow:

- We propose APT, a lightweight and model-agnostic plugin for overcoming limitations of local statistical normalization and mitigating distribution shift.
- APT generates dynamic affine parameters via timestamp-conditioned prototype matching and self-supervised learning, delivering learnable distribution features compatible with any forecasting backbone and normalization.
- Extensive experiments with diverse backbones and normalization methods confirm APT’s effectiveness in improving forecasting accuracy under distribution shift.

## Related Work

### Deep Time Series Forecasting

As the backbone of forecasting, time series models are primarily designed to learn temporal patterns such as seasonality and trends (Wang et al. 2025; Shao et al. 2025). Deep models like Informer (Zhou et al. 2021) improve long-term forecasting, while DLinear (Zeng et al. 2023) shows that simple MLP can also perform well. Recent developments in forecasting focus on complex temporal dependencies, including intra-channel long-term dependencies, as explored by PatchTST (Nie et al. 2023) and SparseTSF (Lin et al. 2024), as well as inter-channel interactions, as addressed by iTransformer (Liu et al. 2024) and CATS (Kim et al. 2024).

Distribution shift is not typically the focus of forecasting backbone, but it is particularly important in the forecasting pipeline because it significantly increases the difficulty of pattern learning and reduces forecasting performance.

### Normalization for Distribution Shift

Normalization has become a standard component in time series forecasting due to its efficiency in handling distribution shift. RevIN (Kim et al. 2021) proposes a distribution-pattern decoupling scheme by forcing the normalization of histories to a unified domain and de-normalizing predictions, which greatly reduces the complexity of non-stationary forecasting (Fu et al. 2025).

Subsequent works aim to recover suppressed distributional features because the distribution of history and future within a single instance is inconsistent. DishTS (Fan et al. 2023) learns intra- and inter-segment drift. SAN (Liu et al. 2023) focuses on patch-level shift. SIN (Han, Ye, and Zhan 2024) relearns statistics based on local invariance and

global variability criteria, and FAN (Ye et al. 2024) replaces normalization with main frequency extraction and residual forecasting.

### Affine Transformation

Affine transformations were initially used to align parameter spaces between layers after normalization, improving training stability of deep models. In conditional generative networks, they serve to reintroduce task-specific information by modulating distributional features (Karras, Laine, and Aila 2019). CIN (Dumoulin, Shlens, and Kudlur 2017) assigns affine parameters to each style via instance normalization. AdaIN (Huang and Belongie 2017) extracts them from target images for arbitrary style transfer, and FiLM (Perez et al. 2018) extends this idea to feature-level conditioning. Such conditional affine modulation has become standard across vision and language tasks (Ziegler et al. 2019).

However, such conditional normalization techniques are rarely explored in time series forecasting, where normalization still mainly depends on local statistics. Inspired by these advances, APT introduces timestamp-based conditioning to generate dynamic affine parameters, enabling better adaptation to distribution shift.

## Preliminaries

### Time Series Forecasting

Given the historical multivariate time series  $\{x_{t-L:t}^{(i)}\}_{i=1}^C = \{x_{t-L:t}^{(1)}, \dots, x_{t-L:t}^{(C)}\} \in \mathbb{R}^{L \times C}$ , where  $x_t$  refer to the values at timestamp  $t$ ,  $L$  is the length of historical window and  $C$  is the number of channels. The objective of time series forecasting is to predict future series  $\{y_{t+1:t+H}^{(i)}\}_{i=1}^C = \{y_{t+1:t+H}^{(1)}, \dots, y_{t+1:t+H}^{(C)}\} \in \mathbb{R}^{H \times C}$  by leveraging the forecasting model  $\mathcal{M}$  to learn patterns from historical data, where  $H$  is the future horizon:

$$\{y_{t+1:t+H}^{(i)}\}_{i=1}^C = \mathcal{M}(\{x_{t-L:t}^{(i)}\}_{i=1}^C) \quad (1)$$

### Time Series Normalization

Given a normalization method  $\mathcal{N}$ , it normalizes historical time series before inputting them into  $\mathcal{M}$  to eliminate distribution differences, and then de-normalizes the time series output by  $\mathcal{M}$  to restore distribution information:

$$\{y_{t+1:t+H}^{(i)}\}_{i=1}^C = \mathcal{N}^{-1}(\mathcal{M}(\mathcal{N}(\{x_{t-L:t}^{(i)}\}_{i=1}^C))) \quad (2)$$

Since statistics often have additional networks in  $\mathcal{N}$  for adaptive learning, the statistics adopted by  $\mathcal{N}$  and  $\mathcal{N}^{-1}$  often differ from the originals and are not identical to each other. Take mean-variance normalization as an example.  $\mu_t$  is the mean of historical time series  $\{x_{t-L:t}^{(i)}\}_{i=1}^C$ , and  $\sigma_t$  is its variance. The subscripts  $l$  and  $h$  donate the normalization and de-normalization stages respectively, where the statistics are either preserved (RevIN) or relearned (DishTS, SAN), depending on the normalization strategy. Equation 2 can be expressed as:

$$y_{t+1:t+H}^{(i)} = \sigma_{h,t}^{(i)} \mathcal{M}\left(\frac{x_{t-L:t}^{(i)} - \mu_{l,t}^{(i)}}{\sigma_{l,t}^{(i)}}\right) + \mu_{h,t}^{(i)} \quad (3)$$

## Methodology

### Overview

In Figure 3, the mainstream paradigm in time series forecasting follows a normalization–forecasting–denormalization pipeline. In this framework, normalization modules (e.g., RevIN, SAN) serve as plug-in components to either remove or learn distributional information, while forecasting models (e.g., SparseTSF, iTransformer) constitute the backbone for capturing temporal patterns.

Affine Prototype-Timestamp (APT) is a lightweight plug-in module designed to break through the limitations of local statistical normalization strategies. Within the forecasting pipeline, the forward transformation of APT is applied after normalization and before model input, and the inverse transformation (de-APT) is applied after model output and before de-normalization.

Instead of learning a single set of affine parameters  $\gamma$  and  $\beta$ , APT leverages global timestamps to adaptively assign distinct affine parameters to different subseries. The goal is to reintroduce learnable distributional features, conditioned on global constraints, after distributional features have been suppressed by normalization or disrupted due to distribution shift. The following sections detail the implementation and training procedure of APT.

### Affine Prototype-Timestamp

To address global distribution shift and reduce reliance on local statistical normalization, we propose APT, a module that learns adaptive affine parameters conditioned on timestamp representations. Timestamps, commonly formatted as %Y-%m-%d %H:%M, are widely available and encode rich temporal semantics, making them a natural proxy for global distributional regularities.

Given the input historical sequence  $\{x_{t-L:t}^{(i)}\}_{i=1}^C$  and forecast horizon  $H$ , we extract the start timestamp  $ts_{t-L}$  of the history and the end timestamp  $ts_{t+H}$  of the forecasting horizon for each sample. To obtain semantically meaningful representations, we discretize timestamps into categorical attributes, such as “Time in Day” (TiD) and “Day in Week” (DiW), which are broadly applicable and empirically correlated with global temporal variation.

Each timestamp  $t$  is represented by the sum of its corresponding attribute embeddings  $\mathbf{T}_t$ :

$$\mathbf{T}_t = \mathbf{T}_t^{\text{TiD}} + \mathbf{T}_t^{\text{DiW}} \quad (4)$$

However, discrete timestamp combinations can be sparse and fail to generalize across unseen scenarios. Inspired by few-shot learning, we replace raw timestamp embeddings with learnable prototypes that capture shared temporal semantics. As a result, we match each timestamp embedding  $\mathbf{T}_t \in \mathbb{R}^D$  against a shared learnable prototype library  $\mathbf{P} = \{\mathbf{p}_j \in \mathbb{R}^D\}_{j=1}^N$  via inner product similarity:

$$\mathbf{S}_t = \mathbf{T}_t \mathbf{P}^\top \quad (5)$$

We then select the top- $k$  most similar prototypes per timestamps and compute their softmax similarity scores:

$$\mathbf{w}_{t,[1:k]} = \text{Softmax}(\mathbf{S}_{t,[1:k]}) \quad (6)$$

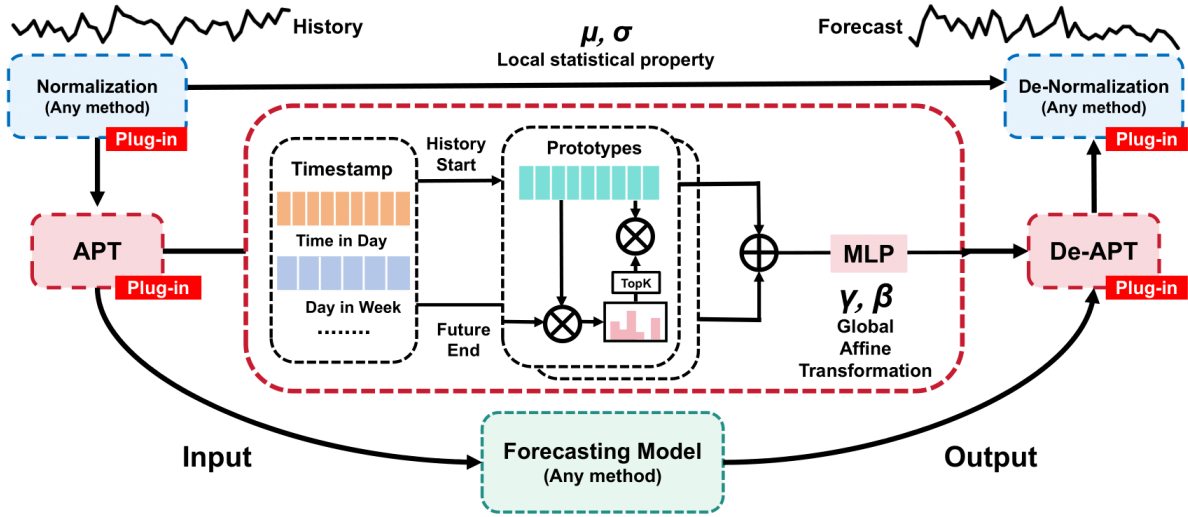


Figure 3: The pipeline of time series forecasting and the schematic of APT.

This yields a sparse weight matrix  $\mathbf{W}_t \in \mathbb{R}^N$ , where only the top- $k$  entries per row are non-zero. The final timestamp representation is formed via a weighted aggregation over the prototype embeddings:

$$\tilde{\mathbf{T}}_t = \sum_{j=1}^N \mathbf{W}_{t,j} \mathbf{p}_j \quad (7)$$

To encode both past and future temporal context, we aggregate their embeddings:  $\tilde{\mathbf{T}}_t = \tilde{\mathbf{T}}_t^{\text{history}} + \tilde{\mathbf{T}}_t^{\text{future}}$ . Optionally, channel identity embeddings  $\mathbf{ID} \in \mathbb{R}^{C \times D}$  can be added to incorporate per-channel variation when necessary.

Next, we employ two multi-layer perceptrons to map the aggregated embedding to channel-wise affine parameters  $\gamma, \beta \in \mathbb{R}^{\{C\} \times 1}$ , where  $\{C\}$  denotes the option of whether to adopt identity embeddings:

$$\gamma_t, \beta_t = \text{MLP}_\gamma(\tilde{\mathbf{T}}_t), \text{MLP}_\beta(\tilde{\mathbf{T}}_t) \quad (8)$$

These learned parameters modulate the normalized time series before and after the forecasting model. Substituting them into the normalization pipeline in Equation 3, the revised inference process becomes:

$$y_{t+1:t+H}^{(i)} = \frac{\sigma_{h,t}^{(i)}}{\gamma_t^{(i)}} \left( \mathcal{M}(\gamma_t^{(i)} \frac{x_{t-L:t}^{(i)} - \mu_{l,t}^{(i)}}{\sigma_{l,t}^{(i)}} + \beta_t^{(i)}) - \beta_t^{(i)} \right) + \mu_{h,t}^{(i)} \quad (9)$$

This formulation enables globally informed affine transformations that better adapt to distributional variations across time.

### Training Strategy

Although APT is conceptually simple, learning timestamp-conditioned affine parameters resembles self-supervised clustering, potentially forming a *bi-level* optimization problem (Gould et al. 2016), similar to SAN (Liu et al. 2023).

To stabilize training and ensure effective learning of timestamp semantics, we first freeze the normalization and backbone, and train APT with additional epochs using the following self-supervised losses:

**Orthogonal Loss.** To ensure diversity among timestamp and prototype embeddings, we encourage orthogonality within  $\mathbf{E} = \{\mathbf{T}^{\text{TiD}}, \mathbf{T}^{\text{DiW}}, \mathbf{P}\}$ . The loss is defined as:

$$\text{Loss}_{orth} = \|\mathbf{E}^\top \mathbf{E} \odot (\mathbf{I} - \mathbf{I})\|_2^2 + \|\mathbf{I} - \mathbf{E}^\top \mathbf{E} \odot \mathbf{I}\|_2^2 \quad (10)$$

Where  $\odot$  is the Hadamard product,  $\mathbf{I}$  is the identity matrix, and  $\|\cdot\|_2$  is the L2 norm.  $\text{Loss}_{orth}$  penalizes correlation between embeddings and enforces unit-norm behavior, following Barlow Twins (Zbontar et al. 2021).

**Load Balancing Loss.** Prototype matching may lead to imbalanced usage, where few prototypes dominate. Inspired by MoE imbalance issues (Wang et al. 2024). To mitigate this, we introduce a load balancing loss that encourages a uniform distribution over prototype usage. This method belongs to batch-wise loss, where we penalize the prototype weights  $\mathbf{W} \in \mathbb{R}^{B \times N}$  in each batch in Equation 7:

$$\text{Loss}_{balance} = \sum_{j=1}^N \left( \frac{\sum_{i=1}^B \mathbf{W}_{i,j}}{\sum_{i=1}^B \sum_{j=1}^N \mathbf{W}_{i,j}} - \frac{1}{N} \right)^2 \quad (11)$$

**Affine Regularization Loss.** To prevent trivial or unstable affine parameters, we follow self-supervised representation regularization strategies (Ermolov et al. 2021; Bardes, Ponce, and Lecun 2022). Taking  $\gamma$  as an example, the loss is:

$$\text{Loss}_R = \frac{1}{B} \left( 1 - \sqrt{\sum_{i=1}^B \gamma_i^2} \right)^2 + \frac{1}{B} \left( \sum_{i=1}^B \gamma_i \right)^2 \quad (12)$$

**Training Processing.** In the additional epochs of freezing the backbone and normalization strategy, the overall pre-training objective for APT is:

$$\text{Loss}_{APT} = \text{Loss}_{orth} + \text{Loss}_{balance} + \text{Loss}_R \quad (13)$$

Once pretrained, APT is jointly optimized with the normalization and forecasting modules under standard regression loss (e.g., MSE). If required by the normalization strategy (e.g., SAN), its loss can be included:

$$\text{Loss}_{main} = \text{Loss}_{normal} + \text{Loss}_{MSE} \quad (14)$$

Methods		CATS				Informer				iTransformer				SparseTSF			
Affine		+APT				+APT				+APT				+APT			
Metrics		MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
ECL	None	0.264	0.168	<b>0.262</b>	<b>0.166</b>	0.426	0.369	<b>0.408</b>	<b>0.332</b>	0.270	0.169	<b>0.265</b>	<b>0.165</b>	0.269	0.173	<b>0.267</b>	<b>0.171</b>
	+RevIN	<b>0.258</b>	<b>0.165</b>	0.259	0.166	0.337	0.252	<b>0.326</b>	<b>0.233</b>	<b>0.259</b>	0.164	<b>0.259</b>	<b>0.162</b>	0.265	<b>0.173</b>	<b>0.264</b>	<b>0.173</b>
	+Dish-TS	0.272	0.173	<b>0.270</b>	<b>0.172</b>	0.379	0.305	<b>0.371</b>	<b>0.297</b>	0.266	0.166	<b>0.261</b>	<b>0.161</b>	0.272	0.173	<b>0.267</b>	<b>0.171</b>
	+SAN	0.277	<b>0.175</b>	<b>0.276</b>	<b>0.175</b>	0.349	0.266	<b>0.324</b>	<b>0.231</b>	<b>0.258</b>	<b>0.159</b>	0.261	0.160	0.271	0.172	<b>0.266</b>	<b>0.165</b>
	+FAN	0.265	0.167	<b>0.263</b>	<b>0.165</b>	0.266	<b>0.165</b>	<b>0.264</b>	0.166	0.274	0.172	<b>0.268</b>	<b>0.169</b>	0.263	0.166	<b>0.261</b>	<b>0.164</b>
ETTh1	None	0.469	0.486	<b>0.436</b>	<b>0.431</b>	0.891	1.255	<b>0.783</b>	<b>1.134</b>	0.499	0.503	<b>0.463</b>	<b>0.457</b>	0.439	0.436	<b>0.437</b>	<b>0.435</b>
	+RevIN	0.443	0.448	<b>0.431</b>	<b>0.427</b>	0.597	0.715	<b>0.562</b>	<b>0.657</b>	0.473	0.485	<b>0.441</b>	<b>0.441</b>	0.429	0.427	<b>0.427</b>	<b>0.424</b>
	+Dish-TS	0.475	0.485	<b>0.452</b>	<b>0.456</b>	0.797	1.062	<b>0.721</b>	<b>0.947</b>	0.498	0.513	<b>0.492</b>	<b>0.495</b>	0.467	0.477	<b>0.456</b>	<b>0.461</b>
	+SAN	0.451	0.472	<b>0.447</b>	<b>0.463</b>	0.542	0.614	<b>0.528</b>	<b>0.599</b>	0.466	0.477	<b>0.453</b>	<b>0.470</b>	0.504	0.564	<b>0.464</b>	<b>0.488</b>
	+FAN	0.478	0.482	<b>0.473</b>	<b>0.478</b>	0.519	0.533	<b>0.506</b>	<b>0.529</b>	0.487	0.497	<b>0.477</b>	<b>0.480</b>	0.484	0.497	<b>0.480</b>	<b>0.488</b>
ETTh2	None	0.504	0.533	<b>0.471</b>	<b>0.486</b>	1.433	2.883	<b>1.194</b>	<b>2.269</b>	0.686	0.863	<b>0.493</b>	<b>0.507</b>	0.492	0.515	<b>0.484</b>	<b>0.498</b>
	+RevIN	0.425	0.392	<b>0.417</b>	<b>0.390</b>	0.565	0.658	<b>0.504</b>	<b>0.531</b>	0.439	0.421	<b>0.422</b>	<b>0.393</b>	<b>0.427</b>	<b>0.399</b>	0.428	0.470
	+Dish-TS	0.476	0.466	<b>0.444</b>	<b>0.419</b>	1.321	3.261	<b>0.795</b>	<b>1.318</b>	0.537	0.587	<b>0.485</b>	<b>0.497</b>	0.631	0.933	<b>0.533</b>	<b>0.612</b>
	+SAN	0.438	0.406	<b>0.433</b>	<b>0.404</b>	0.718	0.939	<b>0.447</b>	<b>0.438</b>	0.438	0.418	<b>0.414</b>	<b>0.410</b>	0.559	0.619	<b>0.485</b>	<b>0.478</b>
	+FAN	0.477	0.487	<b>0.470</b>	<b>0.470</b>	0.565	0.624	<b>0.489</b>	<b>0.499</b>	0.492	0.512	<b>0.483</b>	<b>0.497</b>	0.475	<b>0.467</b>	<b>0.472</b>	0.473
Exchange	None	0.617	0.888	<b>0.415</b>	<b>0.316</b>	0.996	1.640	<b>0.920</b>	<b>1.337</b>	0.655	0.783	<b>0.551</b>	<b>0.599</b>	0.426	0.352	<b>0.418</b>	<b>0.342</b>
	+RevIN	0.435	0.401	<b>0.424</b>	<b>0.381</b>	0.592	0.587	<b>0.519</b>	<b>0.460</b>	0.482	0.468	<b>0.447</b>	<b>0.416</b>	0.481	0.478	<b>0.438</b>	<b>0.386</b>
	+Dish-TS	0.562	0.774	<b>0.493</b>	<b>0.468</b>	0.930	2.105	<b>0.704</b>	<b>1.003</b>	0.509	0.496	<b>0.476</b>	<b>0.395</b>	0.540	0.531	<b>0.502</b>	<b>0.476</b>
	+SAN	0.483	0.521	<b>0.415</b>	<b>0.372</b>	0.457	0.410	<b>0.404</b>	<b>0.338</b>	0.496	0.499	<b>0.457</b>	<b>0.448</b>	0.418	<b>0.350</b>	<b>0.407</b>	0.356
	+FAN	0.492	0.458	<b>0.482</b>	<b>0.457</b>	0.562	0.591	<b>0.533</b>	<b>0.528</b>	0.513	0.502	<b>0.486</b>	<b>0.448</b>	0.473	0.443	<b>0.462</b>	<b>0.419</b>
Traffic	None	0.288	0.554	<b>0.286</b>	<b>0.530</b>	0.425	0.830	<b>0.402</b>	<b>0.780</b>	0.594	0.986	<b>0.428</b>	<b>0.830</b>	0.298	<b>0.443</b>	<b>0.296</b>	<b>0.443</b>
	+RevIN	0.284	0.421	<b>0.280</b>	<b>0.416</b>	0.457	0.874	<b>0.394</b>	<b>0.753</b>	<b>0.289</b>	<b>0.412</b>	0.293	0.415	<b>0.295</b>	<b>0.445</b>	<b>0.295</b>	<b>0.445</b>
	+Dish-TS	0.293	0.438	<b>0.278</b>	<b>0.417</b>	0.445	0.861	<b>0.403</b>	<b>0.761</b>	<b>0.305</b>	<b>0.429</b>	0.309	0.430	0.314	0.461	<b>0.310</b>	<b>0.457</b>
	+SAN	0.292	<b>0.432</b>	<b>0.291</b>	0.437	0.419	0.753	<b>0.394</b>	<b>0.716</b>	<b>0.294</b>	<b>0.427</b>	0.297	0.431	0.407	0.651	<b>0.363</b>	<b>0.569</b>
	+FAN	0.316	0.454	<b>0.301</b>	<b>0.438</b>	0.315	0.457	<b>0.301</b>	<b>0.445</b>	0.340	0.468	<b>0.324</b>	<b>0.454</b>	0.302	0.432	<b>0.298</b>	<b>0.431</b>
Weather	None	0.280	0.229	<b>0.268</b>	<b>0.223</b>	0.396	0.401	<b>0.311</b>	<b>0.267</b>	0.302	0.256	<b>0.285</b>	<b>0.239</b>	0.307	<b>0.254</b>	<b>0.299</b>	0.260
	+RevIN	<b>0.258</b>	<b>0.221</b>	0.259	0.224	0.294	0.279	<b>0.274</b>	<b>0.243</b>	<b>0.267</b>	<b>0.233</b>	0.269	<b>0.233</b>	0.283	<b>0.252</b>	<b>0.282</b>	<b>0.252</b>
	+Dish-TS	0.277	<b>0.225</b>	<b>0.276</b>	0.226	0.316	0.296	<b>0.295</b>	<b>0.276</b>	0.296	0.258	<b>0.285</b>	<b>0.240</b>	0.301	<b>0.240</b>	<b>0.299</b>	<b>0.240</b>
	+SAN	<b>0.277</b>	<b>0.225</b>	0.279	0.227	0.286	0.261	<b>0.276</b>	<b>0.245</b>	<b>0.271</b>	<b>0.228</b>	0.276	0.235	0.279	<b>0.226</b>	<b>0.277</b>	0.227
	+FAN	<b>0.281</b>	<b>0.231</b>	0.283	0.234	0.285	0.246	<b>0.280</b>	<b>0.245</b>	0.288	0.236	<b>0.276</b>	<b>0.230</b>	<b>0.276</b>	<b>0.228</b>	0.278	0.229
1 <sup>st</sup> count		4	6	<b>26</b>	<b>24</b>	0	1	<b>30</b>	<b>29</b>	7	6	<b>24</b>	<b>25</b>	3	11	<b>28</b>	<b>23</b>

Table 2: Results of the main experiment

## Experiments

### Experimental Setup

**Baselines.** As a plug-in designed to mitigate distribution shift, APT is compatible with any time series forecasting backbone and normalization strategy. To validate its effectiveness, we select Informer (Zhou et al. 2021), iTransformer (Liu et al. 2024), SparseTSF (Lin et al. 2024), and CATS (Kim et al. 2024) as the backbone, and RevIN (Kim et al. 2021), Dish-TS (Fan et al. 2023), SAN (Liu et al. 2023), and FAN (Ye et al. 2024) as normalization strategies.

**Datasets & Metrics.** We select six datasets for main experiments: ECL, ETTh1, ETTh2, Exchange, Traffic, Weather. The historical length  $L$  is 336, the main experiments report average results over forecasting lengths  $H = \{96, 192, 336, 720\}$ , and the average value is reported, while extra experiments use  $H = 336$ . Metrics include mean squared error (MSE) and mean absolute error (MAE).

**Details.** We follow the BasicTS benchmark setup (Shao et al. 2024). APT uses “Day in Week” and “Time in Day” as timestamp features, except on Exchange (only “Day of Week”). The number of prototypes depends on sampling

rate: 5 for daily, 30 for hourly, and 40 for 10-minute frequency. All embeddings are 20-dimensional with MLP hidden size 32, yielding 1.5K–5K parameters. APT trains with a learning rate of  $5e-5$  for 1–5 extra epochs until convergence. The optimizer is shared with other modules, with gradient updates controlled by loss weight  $\lambda$ . More details are in Appendix B and C.

### Main Results

On datasets with strong distribution shifts like ETTh1/2 and Exchange, APT yields 4%–40% performance gains across diverse backbone–normalization combinations. Combined with robust normalization, even the weaker Informer matches or surpasses state-of-the-art backbones.

CATS, iTransformer, and SparseTSF are among the most advanced backbones, while RevIN, FAN and SAN are the most effective normalization strategies. Even on stable datasets like ECL, where performance typically plateaus, APT still yields consistent gains and often achieves the best overall results with negligible risk of degradation. More and combined with data analysis are provided in Appendix C.4.

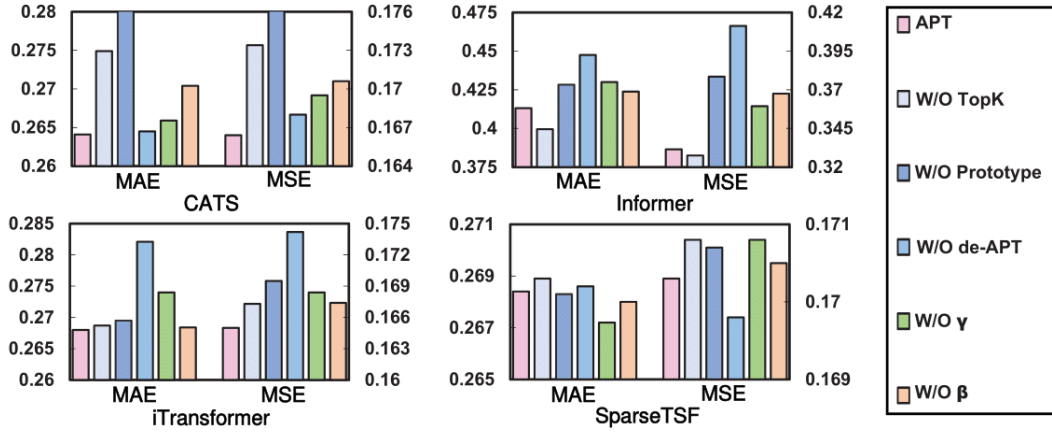


Figure 4: The ablation study results of APT components on the ECL dataset,  $L = 336, H = 336$ .

### Ablation Study

In Figure 4 and Appendix C.7, we conduct ablations from two perspectives: component-wise and training strategy.

**Component Ablation:** We evaluate APT by removing key components:

- **W/O Top- $k$ :** Remove Top- $k$  from Equation 6;
- **W/O Prototype:** Use raw timestamp embeddings without prototype matching;
- **W/O de-APT:** Remove inverse transformation of APT;
- **W/O  $\gamma$  or  $\beta$ :** Remove scaling or bias in affine transformation;

Removing **de-APT** degrades performance, as it forces the model to reconstruct distributional properties. Top- $k$  is similar to expert activation and mitigates feature over-smoothing by enforcing sparsity during prototype aggregation. In contrast, prototype learning alleviates the challenges of few-shot learning by compressing timestamp semantics, improving robustness and preventing overfitting. Between affine components,  $\gamma$  is more crucial than the bias  $\beta$  under distribution shift, because mean drift will generate abundant supervisory signals to guide the backbone to adapt.

**Training Strategy Ablation:** APT learns timestamp-aware affine parameters via self-supervised pretraining, while freezing backbone and normalization. Table 3 and Appendix C show the effect of ablating components:

- **W/O  $\lambda$ :** Remove the loss weighting coefficient from APT pre-training, sharing learning rate with the backbone;
- **W/O  $Loss_{orth}$  or  $Loss_{balance}$  or  $Loss_R$ :** Remove orthogonality, load balancing, or affine regularization loss;
- **W/O  $Loss_{APT}$ :** Skip pretraining entirely, optimizing APT only with the main loss.

Recommended learning rates for {CATS, Informer, iTransformer, SparseTSF} are  $\{5e-3, 2e-4, 5e-4, 2e-3\}$ , with APT pretraining fixed at  $5e-5$ . Accordingly,  $\lambda$  is set to  $\{5e-2, 2.5e-1, 1e-1, 2.5e-2\}$ . Without  $\lambda$ , APT shares the backbone’s rate, causing overfitting and unstable timestamp semantics.

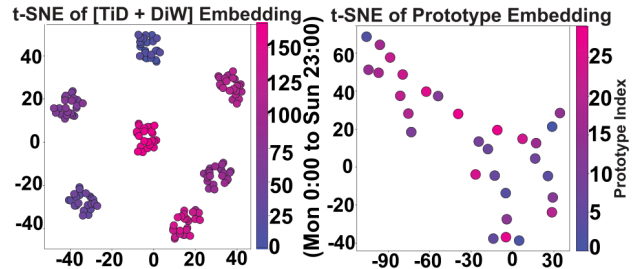


Figure 5: Visualization of APT timestamps and prototype embeddings on ECL dataset and iTransformer

Each pretraining loss targets a distinct goal:  $Loss_{orth}$  encourages embedding diversity,  $Loss_{balance}$  prevents feature collapse (Hua et al. 2021), and  $Loss_R$  stabilizes convergence for regression compatibility. As they govern embeddings, prototype matching, and parameter scaling respectively, these losses are complementary and essential for APT’s stability. Omitting pretraining altogether (W/O  $Loss_{APT}$ ) causes APT to fail, with performance reverting to backbone levels.

### Visualization

**t-SNE of Embedding:** Figure 5 shows the timestamp and prototype embeddings learned by APT with iTransformer on the ECL dataset. Timestamp embeddings form distinct clusters aligned with “Day in Week,” while “Time in Day” varies orthogonally in Appendix C.9. Prototypes are also well-separated, but their embeddings still preserve certain timestamp-related structure as orthogonal losses serve as soft supervision.

**Forecasting cases:** Figure 6 presents forecasting cases on the 6th channel of ETTh2 without normalization. While most backbones, except Informer, capture periodic patterns well, they fail to handle distribution shifts over time. APT mitigates this issue by applying timestamp-conditioned affine transformations, independently of the backbone, to better align predictions with data distributions.

Methods		APT		W/O $\lambda$		W/O $Loss_{orth}$		W/O $Loss_{balance}$		W/O $Loss_R$		W/O $Loss_{APT}$	
Metrics		MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
ECL	CATS	<b>0.264</b>	<b>0.166</b>	0.269	0.169	<u>0.265</u>	<u>0.167</u>	0.266	0.169	<u>0.265</u>	0.168	0.267	0.170
	Informer	0.413	0.332	0.441	0.379	0.429	0.350	<b>0.412</b>	<b>0.329</b>	0.436	0.377	0.453	0.395
	iTransformer	<b>0.268</b>	<b>0.165</b>	0.269	0.169	<b>0.268</b>	<b>0.165</b>	0.269	0.167	0.269	0.168	0.270	0.168
	SparseTSF	<b>0.268</b>	<b>0.170</b>	<u>0.270</u>	0.171	<u>0.269</u>	0.171	<u>0.269</u>	0.171	0.269	0.171	0.268	<u>0.171</u>
ETTh1	CATS	<b>0.437</b>	<b>0.443</b>	0.451	0.458	<u>0.438</u>	<u>0.446</u>	0.448	0.460	0.439	<u>0.446</u>	0.448	0.459
	Informer	0.797	<b>1.121</b>	0.856	1.249	0.839	1.215	1.009	1.495	<b>0.787</b>	<b>1.121</b>	0.818	<u>1.198</u>
	iTransformer	0.474	0.478	0.472	0.481	0.472	0.479	<b>0.469</b>	<b>0.474</b>	0.487	0.498	0.483	0.491
	SparseTSF	<b>0.438</b>	<b>0.451</b>	0.447	0.458	<u>0.441</u>	<u>0.453</u>	0.448	0.460	0.441	0.454	0.448	0.459
Exchange	CATS	<b>0.402</b>	<b>0.276</b>	0.411	0.286	0.407	0.290	<u>0.406</u>	<u>0.278</u>	0.402	<b>0.276</b>	0.560	0.500
	Informer	<b>0.976</b>	1.490	1.144	2.070	1.009	1.789	<u>1.009</u>	<b>1.506</b>	0.989	1.490	1.082	1.734
	iTransformer	<b>0.502</b>	<b>0.445</b>	0.634	0.6549	<u>0.629</u>	0.669	<u>0.559</u>	0.531	<u>0.544</u>	<u>0.508</u>	0.655	0.785
	SparseTSF	<b>0.431</b>	<b>0.316</b>	<u>0.432</u>	<u>0.321</u>	0.443	0.333	<b>0.431</b>	<b>0.316</b>	0.439	0.330	0.448	0.343
Weather	CATS	<b>0.283</b>	<b>0.242</b>	<u>0.287</u>	<u>0.243</u>	0.300	0.246	0.300	0.247	0.295	<u>0.243</u>	0.301	0.246
	Informer	<b>0.315</b>	<b>0.270</b>	<b>0.315</b>	0.271	0.326	0.284	0.328	0.291	0.410	0.389	0.554	0.816
	iTransformer	<b>0.293</b>	<b>0.251</b>	0.296	<u>0.254</u>	<b>0.293</b>	<b>0.251</b>	<u>0.295</u>	<u>0.254</u>	0.296	<u>0.254</u>	0.300	0.257
	SparseTSF	<b>0.312</b>	<b>0.267</b>	<u>0.314</u>	<b>0.267</b>	0.320	0.270	<u>0.322</u>	0.271	0.320	0.270	0.319	<u>0.269</u>

Table 3: Results of ablation studies on training strategies across four datasets.  $L = 336, H = 336$

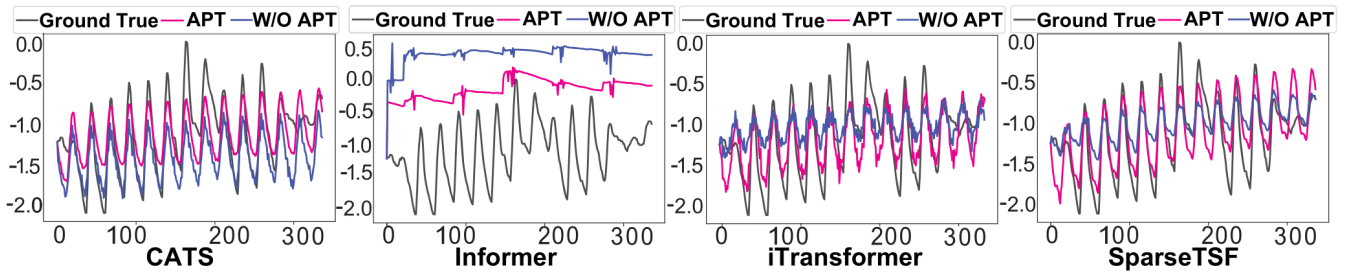


Figure 6: Visualization of forecasting results for different models on the ETTh2 dataset without normalization strategy

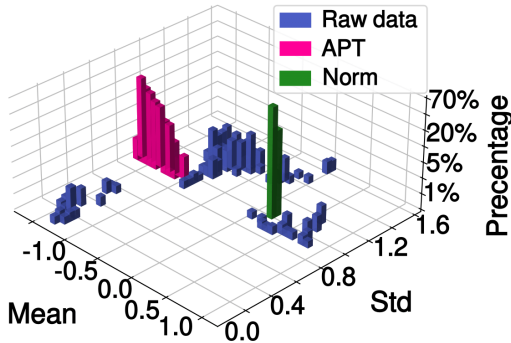


Figure 7: 3D visualization of temporal distribution and ratio of forecasting pipeline at different stages on ETTh1.

**Temporal distribution:** As shown in Figure 7, raw data exhibit severe distribution shift and its distributional features are disrupted, which is the phenomenon we emphasized earlier and it can cause great damage to backbone’s learning of temporal patterns. Normalization maps series to zero mean and unit variance; while RevIN decouples distribution from patterns to improve robustness, it also suppresses informa-

tive distribution. In contrast, APT uses timestamps to encode global distribution features, forming compact yet discriminative representations adaptable to temporal shift.

Due to space limitations, the complete results and more experiments are provided in Appendix C to better understand APT: data motivation (C.2), other plug-ins (C.3), extra main results (C.4), hyper-Parameter (C.5), parameter count (C.6), extra ablation study (C.7), cross-setting fine-tuning (C.8), visualization (C.9) and discussion (D).

## Conclusion

In this work, we propose Affine Prototype-Timestamp (APT), a lightweight and model-agnostic plug-in to enhance time series forecasting under distribution shifts. APT employs discretized timestamps and prototype learning to introduce timestamp-conditioned affine transformations, enabling forecasting models to recover global distributional features that may be suppressed by normalization or distorted by distribution shifts.

In future work, we will explore online adaptation for streaming data and incorporate external modalities such as text or images to enhance shift awareness, further strengthening the real-world applicability of deep forecasting models in domains like weather, finance, and energy.

## Acknowledgments

This work is supported by the NSFC under Grant Nos. 62372430 and 62502505, the Youth Innovation Promotion Association CAS No.2023112, the Postdoctoral Fellowship Program of CPSF under Grant Number GZC20251078, the China Postdoctoral Science Foundation No.2025M771542 and HUA Innovation fundings.

## References

- Bardes, A.; Ponce, J.; and Lecun, Y. 2022. VICReg: Variance-Invariance-Covariance Regularization For Self-Supervised Learning. In *ICLR*.
- Bebis, G.; Georgiopoulos, M.; da Vitoria Lobo, N.; and Shah, M. 1999. Learning affine transformations. *Pattern recognition*, 32(10): 1783–1799.
- Dumoulin, V.; Shlens, J.; and Kudlur, M. 2017. A Learned Representation For Artistic Style. In *ICLR*.
- Ermolov, A.; Siarohin, A.; Sangineto, E.; and Sebe, N. 2021. Whitening for self-supervised representation learning. In *ICML*, 3015–3024. PMLR.
- Fan, W.; Wang, P.; Wang, D.; Wang, D.; Zhou, Y.; and Fu, Y. 2023. Dish-ts: a general paradigm for alleviating distribution shift in time series forecasting. In *AAAI*, volume 37, 7522–7529.
- Ferté, T.; Dutartre, D.; Hejblum, B. P.; Griffier, R.; Jouhet, V.; Thiébaud, R.; Legrand, P.; and Hinaut, X. 2024. Reservoir Computing for Short High-Dimensional Time Series: an Application to SARS-CoV-2 Hospitalization Forecast. In *ICML*, 13570–13591. PMLR.
- Fu, Y.; Shao, Z.; Yu, C.; Li, Y.; An, Z.; Wang, Q.; Xu, Y.; and Wang, F. 2025. Selective Learning for Deep Time Series Forecasting. *arXiv preprint arXiv:2510.25207*.
- Gould, S.; Fernando, B.; Cherian, A.; Anderson, P.; Cruz, R. S.; and Guo, E. 2016. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *arXiv preprint arXiv:1607.05447*.
- Han, L.; Ye, H.-J.; and Zhan, D.-C. 2024. SIN: selective and interpretable normalization for long-term time series forecasting. In *ICML*.
- Hua, T.; Wang, W.; Xue, Z.; Ren, S.; Wang, Y.; and Zhao, H. 2021. On feature decorrelation in self-supervised learning. In *CVPR*, 9598–9608.
- Huang, X.; and Belongie, S. 2017. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 1501–1510.
- Karras, T.; Laine, S.; and Aila, T. 2019. A style-based generator architecture for generative adversarial networks. In *CVPR*, 4401–4410.
- Kim, D.; Park, J.; Lee, J.; and Kim, H. 2024. Are Self-Attentions Effective for Time Series Forecasting? In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Kim, T.; Kim, J.; Tae, Y.; Park, C.; Choi, J.-H.; and Choo, J. 2021. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *ICLR*.
- Li, T.; Li, Z.; Rockwell, H.; Farimani, A.; and Lee, T. S. 2022. Prototype memory and attention mechanisms for few shot image generation. In *ICLR*, volume 18.
- Li, Y.; Shao, Z.; Xu, Y.; Qiu, Q.; Cao, Z.; and Wang, F. 2024. Dynamic frequency domain graph convolutional network for traffic forecasting. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 5245–5249. IEEE.
- Li, Y.; Zezhi, S.; Yu, C.; Qian, T.; Zhang, Z.; Du, Y.; He, S.; Wang, F.; and Xu, Y. 2025. STA-GANN: A Valid and Generalizable Spatio-Temporal Kriging Approach. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management*, 1726–1736.
- Lin, S.; Lin, W.; Wu, W.; Chen, H.; and Yang, J. 2024. SparseTSF: Modeling Long-term Time Series Forecasting with\* 1k\* Parameters. In *ICML*.
- Liu, Y.; Hu, T.; Zhang, H.; Wu, H.; Wang, S.; Ma, L.; and Long, M. 2024. iTransformer: Inverted Transformers Are Effective for Time Series Forecasting. In *ICLR*.
- Liu, Z.; Cheng, M.; Li, Z.; Huang, Z.; Liu, Q.; Xie, Y.; and Chen, E. 2023. Adaptive normalization for non-stationary time series forecasting: A temporal slice perspective. *NeurIPS*, 36: 14273–14292.
- Nie, Y.; Nguyen, N. H.; Sinthong, P.; and Kalagnanam, J. 2023. A Time Series is Worth 64 Words: Long-term Forecasting with Transformers. In *ICLR*.
- Perez, E.; Strub, F.; De Vries, H.; Dumoulin, V.; and Courville, A. 2018. Film: Visual reasoning with a general conditioning layer. In *AAAI*, volume 32.
- Shao, Z.; Li, Y.; Wang, F.; Yu, C.; Fu, Y.; Qian, T.; Xu, B.; Diao, B.; Xu, Y.; and Cheng, X. 2025. Blast: Balanced sampling time series corpus for universal forecasting models. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, 2502–2513.
- Shao, Z.; Wang, F.; Xu, Y.; Wei, W.; Yu, C.; Zhang, Z.; Yao, D.; Sun, T.; Jin, G.; Cao, X.; et al. 2024. Exploring progress in multivariate time series forecasting: Comprehensive benchmarking and heterogeneity analysis. *TKDE*, 37(1): 291–305.
- Shao, Z.; Zhang, Z.; Wang, F.; Wei, W.; and Xu, Y. 2022. Spatial-temporal identity: A simple yet effective baseline for multivariate time series forecasting. In *CIKM*, 4454–4458.
- Wang, F.; Li, Y.; Shao, Z.; Yu, C.; Fu, Y.; An, Z.; Xu, Y.; and Cheng, X. 2025. ARIES: Relation Assessment and Model Recommendation for Deep Time Series Forecasting. *arXiv preprint arXiv:2509.06060*.
- Wang, L.; Gao, H.; Zhao, C.; Sun, X.; and Dai, D. 2024. Auxiliary-loss-free load balancing strategy for mixture-of-experts. *arXiv preprint arXiv:2408.15664*.
- Ye, W.; Deng, S.; Zou, Q.; and Gui, N. 2024. Frequency Adaptive Normalization For Non-stationary Time Series Forecasting. In *NeurIPS*.

Yu, C.; Wang, F.; Shao, Z.; Qian, T.; Zhang, Z.; Wei, W.; An, Z.; Wang, Q.; and Xu, Y. 2025. GinAR+: A Robust End-To-End Framework for Multivariate Time Series Forecasting with Missing Values. *TKDE*.

Zbontar, J.; Jing, L.; Misra, I.; LeCun, Y.; and Deny, S. 2021. Barlow twins: Self-supervised learning via redundancy reduction. In *ICML*, 12310–12320. PMLR.

Zeng, A.; Chen, M.; Zhang, L.; and Xu, Q. 2023. Are transformers effective for time series forecasting? In *AAAI*, volume 37, 11121–11128.

Zhou, H.; Zhang, S.; Peng, J.; Zhang, S.; Li, J.; Xiong, H.; and Zhang, W. 2021. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *AAAI*, volume 35, 11106–11115.

Ziegler, Z. M.; Melas-Kyriazi, L.; Gehrmann, S.; and Rush, A. M. 2019. Encoder-agnostic adaptation for conditional language generation. *arXiv preprint arXiv:1908.06938*.