# LEARNED DATA TRANSFORMATION: A DATA-CENTRIC PLUGIN FOR ENHANCING TIME SE RIES FORECASTING

Anonymous authors

006

007

012 013

014

015

016

017

018

019

021

025 026

027

Paper under double-blind review

## Abstract

Data-centric approaches in Time Series Forecasting (TSF) often involve heuristicbased operations on data. This paper proposes to find a general end-to-end data transformation that serves as a plugin to enhance any arbitrary TSF model's performance. Our idea is to generate transformed data during an approximating process and to co-train a predictor for evaluating data with the transformation. To achieve this, we propose the Proximal Transformation Network (PTN), which learns effective transformations while maintaining proximity to the raw data to ensure fidelity. When orthogonally integrated with popular TSF models, our method helps achieve state-of-the-art performance on seven real-world datasets. Additionally, we show that the proximal transformation process can be interpreted in terms of predictability and distribution alignment among channels, highlighting the potential of data-centric methods for future research. Our code is available at https://anonymous.4open.science/r/PTN-2FC6/.

## 1 INTRODUCTION

Time Series Forecasting (TSF) is a significant task in various domains including energy, finance, engineering, and industrial applications. Recently, along with the tremendous achievements of neural networks in deep learning, TSF has seen many learning-based methods that consistently outperform traditional methods in terms of precision, robustness, and capacity (Zhou et al., 2021; Wu et al., 2021; Salinas et al., 2020; Taylor & Letham, 2018).

034 Meanwhile, modeling specific time series data properties such as distribution shift (Kim et al., 2021), non-stationarity (Liu et al., 2022), frequency bias (Piao et al., 2024), and variate covariance (Wang 035 et al., 2024c) has achieved outstanding successes in TSF. The successes can be attributed to sharp 036 observations of commonalities in time series and sound designs from a model-centric perspective. 037 With the angle shifted, data-centric methods for time series including normalizations (Fan et al., 2023; Han et al., 2024; Deng et al., 2021; Liu et al., 2024c; September et al., 2024), augmentations (Zheng et al., 2024), and simple preprocessing methods (e.g., down-sampling (Yu et al., 2023; 040 Wang et al., 2024a) and patching (Nie et al., 2023; Wang et al., 2024b)) can improve the overall 041 performance of a given model by adopting data-level operations. These methods either explicitly 042 or implicitly entail a specific transformation applied to the raw time series data. This paper further 043 pushes the boundary of data-centric methods by answering the following: Is there a general data 044 transformation that can improve a model's generalization capability? If so, can we discover it in an 045 end-to-end manner?

We first verify the insufficiency of specific designs to generalize on data with different complexities. The evidence is provided by an example that involves learning on noisy input-target pairs generated by a linear mapping. When smoothing the noised data is the only transformation to consider, linear mapping with different complexities has different best-denoised data selected from a range of common denoising/smoothing methods. We then reformulate our problem to identify a general process, which finds an effective transformed data given any TSF model to achieve better forecasting results than if using the raw time series data directly. We refer to this problem as finding the optimal transformed data for forecasting, which extends the original optimization space by adding another dimension w.r.t. data fidelity. However, finding such a transformation is challenging. Assuming the proximal existence of our objective, we constrain the problem as a simplified
co-optimization on the proximity of the transformed data and forecasting accuracy for each level
of proximity. The transformation is provided by a Convolution Encoder, and an approaching process is learned by an attention-based Decoder. Combining these two components, we propose the
Proximal Transformation Network (PTN) that approximates the raw data at different distances to
find the transformed data that results in better performance for a given model.

060 In summary, our contributions lie in the following aspects:

- A Reformulated Problem. We consider the conventional TSF problem as a two-step problem: finding a data transformation that enables better generalization and performance, followed by learning to forecast on the transformed data.
  - A *Plug-and-Play Module*. We propose the Proximal Transformation Network to find an effective data transformation. The model includes a convolution-based Encoder and an attention-based Decoder that provide transformation on different levels of proximity. The training involves a co-optimization of the proximity of the transformed data and forecasting accuracy.
- *Extensive Experiments on Real-World Datasets*. By orthogonally adding our module to several representative models, we managed to achieve state-of-the-art (SOTA) performance on seven real-world datasets. The Proximal Transformation can be further interpreted by predictability assessments and distribution alignments.
- 072 073

074

065

066

067

## 2 RELATED WORK

We focus on data-centric aspects within TSF that transform data explicitly or implicitly. A discussion of TSF models based on their designs and advances is presented in Appendix A.2.

Normalization Normalization is crucial for neural networks on time series data (Bhanja & Das, 2019). Techniques like Reversible Instance Normalization (Kim et al., 2021) address instance-level distribution shifts, influencing models like the non-stationary Transformer (Liu et al., 2022). Extensions include local and global normalizations (Fan et al., 2023; Han et al., 2024), and channel-wise normalization that aligns multivariate series via adaptive modules (Deng et al., 2021; Zhao & Shen, 2024). Such adaptive normalizations enhance expressiveness with additional parameters (Liu et al., 2024c; September et al., 2024).

Data Augmentation Traditional data augmentations (e.g., cropping, jittering, and scaling (Wen et al., 2020)) improve time series representation learning (Yue et al., 2022; Fan et al., 2020). In-foTS (Luo et al., 2023) is a novel proposed meta-learning-based method that learns an end-to-end augmentation from a predefined set of transformations for Contrastive Learning. AutoTCL (Zheng et al., 2024) further offers theoretical insights based on Information Bottleneck Theory (Tishby et al., 2000) and demonstrates effectiveness on CoST (Woo et al., 2022). These adaptive data augmentations still explicitly rely on specific operations, unable to provide a general data transformation.



Figure 1: Four typical types of **Moving Kernel Smoothing** methods in TSF models (yellow grids indicate the elements involved in kernel computation, shadowed grids indicate the position to place the result): (a) *Moving Average* that computes average value within a kernel window (Zeng et al., 2023); (b) *Patching* that computes a learnable weighted sum within the kernel (Nie et al., 2023); (c) *Instance Normalization* that substracts the mean value of the given series (Kim et al., 2021; Li et al., 2023),  $\frac{n-1}{n}$  for the centered weight and  $-\frac{1}{n}$  for the rest; and (d) *Sparse Techinque* that computes weighted average at a defined stride (Lin et al., 2024).

102 103 104

091

096

098

100

101

Moving Kernel Smoothing Several methods function as smoothing techniques that inherently transform data (see Figure 1). DLINEAR (Zeng et al., 2023) uses moving average decomposition with linear models to outperform Transformers. RLINEAR (Li et al., 2023) simplifies this with instance normalization, achieving nearly SOTA performance. The patching operation (Nie et al.,

2023) effectively aggregates features, inspiring works on patch-level features (Wang et al., 2024b;
Tang & Zhang, 2024). Recently, SPARSETSF (Lin et al., 2024) introduces a sparse downsampling
technique (Yu et al., 2023; Wang et al., 2024a), matching SOTA models with only 1k parameters.

Most of these works have not been applied to recent models, and designs like channel-wise and multi-level normalizations are unorthogonal to some backbones. The most comparable work is AutoTCL (Zheng et al., 2024). However, unlike AutoTCL, our method does not heavily emphasize strong theoretical guarantees and focuses on supervised learning (instead of contrastive learning) that typically delivers superior performance in TSF. Moreover, our method is designed as an end-toend data transformation framework, encompassing all three key aspects mentioned above.

117 118 119

120

121

128 129

156

157 158

159

160

161

## **3** ANALYSIS ON LEARNED DATA TRANSFORMATION

## 3.1 MOTIVATION EXAMPLE

**Data Synthesis** Learning a linear mapping with a linear model from (X, Y) data pairs generated by the same linear mapping would be a good example to gain insights into effective data transformations. To make a meaningful example, we introduce independent and identically distributed (IID) **noise** and apply typical denoising methods, as depicted in Figure 1. Generally, these practices offset IID noise by computing weighted averages of data points within moving kernels. Back to our example, we generate a synthetic toy dataset as follows.

$$Y = W(X + \epsilon_X) + \epsilon_Y \tag{1}$$

s.t. 
$$|W|_1 \le k$$
 (2)

where  $X \in \mathbb{R}^{H \times N}$ ,  $Y \in \mathbb{R}^{L \times N}$  where L and H are the series lengths and N is the data point number; noises  $\epsilon_X \in \mathbb{R}^{L \times N}$ ,  $\epsilon_Y \in \mathbb{R}^{H \times N}$  are sampled from the Normal distribution  $\mathcal{N}(\mu, \sigma)$ ; and W is the weight matrix subject to an  $\ell_1$ -norm constraint.

In our synthesis, we represent the **complexity**<sup>1</sup> of the X-to-Y mapping using the  $\ell_1$ -norm of  $W^2$ and constrain it with a parameter k (see Equation 2). To ensure that X and Y have equal means thereby eliminating concerns about distribution shift — we apply Softmax to the last dimension of W. This allows us to focus on the impact of smoothing/denoising methods on the performance of linear models. We also introduce an ORACLE linear mapping method, using a hypothetically known W to directly obtain Y.

140Table 1: Typical linear models on the toy examples: (1) MALINEAR, a linear model applying141Moving Average (Figure 1(a)) on X only; (2) BIMALINEAR, a linear model applying Moving142Average on both X and Y; (3) PATCHLINEAR, using Patching (Figure 1(b)), PATCHTST with143attention removed; (4) RLINEAR, using Instance Norm (Figure 1(c)), the exact copy of (Li et al.,1442023); (5) SPARSETSF, using Sparse Technique (Figure 1(d)), the exact copy of (Lin et al., 2024);145(6) ORACLE forecasting by Y = WX, presented as a relatively strong baseline to compare with.

Models Smoothing	MALINEAR <sup>*</sup> Moving Average		B1MAI Moving	LINEAR <sup>*</sup> Average	RLI Instanc	NEAR e Norm	PATCH1 Patc	LINEAR hing	SPARS Sparse T	seTSF Technique	ORA	CLE
$ W _1$	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
0.0030	0.0290	0.1307	0.0275	0.1258	0.0119	0.0866	0.0103	0.0809	0.0101	0.0801	0.0101	0.0801
0.0052	0.0150	0.0970	0.0141	0.0943	0.0111	0.084	0.0108	0.0830	0.0101	0.0803	0.0101	0.0803
0.0104	0.0111	0.0841	0.0111	0.084	0.0106	0.0822	0.0108	0.0829	0.0103	0.0808	0.0103	0.0808
0.0143	0.0168	0.1035	0.0112	0.0844	0.0107	0.0826	0.0105	0.0818	0.0106	0.0821	0.0107	0.0826
0.0214	0.0169	0.1037	0.0114	0.0850	0.0108	0.0830	0.0106	0.0822	0.0112	0.0846	0.0117	0.0861
0.0437	0.0175	0.1056	0.0120	0.0872	0.0113	0.0847	0.0113	0.0850	0.0139	0.0941	0.0157	0.0997
0.0761	0.0148	0.0950	0.0142	0.0946	0.0127	0.0897	0.0132	0.0917	0.0229	0.1195	0.0265	0.1291

Decomposition, from another angle, interprets Moving Average results as Trend term and the residual as Seasonal term. In our toy example with synthesized noise, the DLINEAR model, which also models the Seasonal term, tends to overfit the noise. We thus use MALINEAR and BIMALINEAR, modified versions of DLINEAR that focus solely on the Trend term for denoising. In the main experiment (Section 5), we use DLINEAR instead of the other two for brevity, since both terms are meaningful in real scenarios.

**Observations** According to Table 1, our experiments with toy datasets using various linear models and smoothing methods reveal that different smoothing methods achieve the best performance de-

<sup>1</sup>A relevant notion is sparsity (Lin et al., 2024); we provide more discussions in Appendix A.7.1.

<sup>&</sup>lt;sup>2</sup>We use the convex and computationally more tractable  $\ell_1$ -norm instead of the  $\ell_0$ -norm (Tibshirani, 1996).

pending on the mapping complexities of the datasets. Since all models use the same one-layer linear
 backbone, performance differences are solely due to smoothing methods.

The performance variation related to the X-to-Y mapping complexity indicates that choosing the 165 right data transformation is crucial. However, this complexity information becomes available only 166 after an approximate training on (X, Y) is completed. For problems with intricacy beyond this 167 toy example, it is less likely to acquire the priors on the underlying data patterns. Additionally, 168 for Patching and Sparse Technique (5th and 6th columns), hyperparameters regarding their kernel 169 designs — patch length (Nie et al., 2023) and period length (Lin et al., 2024), respectively — are 170 often sensitive to different datasets. Therefore, it is essential to *develop transformations that are* 171 *learnable and general*, enabling them to adapt to any data and model without relying on heuris-172 tics. Another observation is that BIMALINEAR consistently outperforms MALINEAR (2nd and 3rd columns), suggesting that using a uniform transformation for both X and Y is more effective. 173

174 175

191

197

199

200

210

## 3.2 TSF PROBLEM REFORMULATION

In many previous studies, TSF follows a paradigm that aims to find the direct mapping from the most recent look-back window to the target series. Let the context series (*resp.* target series) be  $X \in \mathbb{R}^{L \times N}$  (*resp.*  $Y \in \mathbb{R}^{H \times N}$ ) with time series length L (*resp.* H) and the data point number N. Without loss of generality, we only consider the univariate scenario in our intuitive reformulation.

**Original Problem.** Find the optimal  $\theta$  minimizing training error:  $\theta^* = \arg \min_{\theta} ||f(X;\theta) - Y||$ ,

using a specified error metric  $|| \cdot ||$ , typically the  $\ell_1$ - or  $\ell_2$ -norm.

As noted in Section 3.1, transformations like smoothing vary in optimality across datasets. We thus propose a more general, learnable data transformation that exists for any data and predictor model.

Proposition 1. Given any time series segmented into (X, Y), and a predictor model  $f(\cdot; \theta)$ where  $||f(X_{test}; \theta^*) - Y_{test}|| > 0$ , with  $\theta^* = \arg\min_{\theta} ||f(X_{train}; \theta) - Y_{train}||$ , there exists a transformation  $(\tilde{X}, \tilde{Y})$  such that  $||f(X_{test}; \tilde{\theta}^*) - Y_{test}|| < ||f(X_{test}; \theta^*) - Y_{test}||$ , where  $\tilde{\theta}^* = \arg\min_{\theta} ||f(\tilde{X}_{train}; \theta) - \tilde{Y}_{train}||$ .

In real scenarios, this proposition holds as long as there exists a distribution difference/shift in raw series, which is common in well-recognized TSF datasets (Kim et al., 2021). We showcase that the better  $\tilde{X}$  and  $\tilde{Y}$  exist in a simple scenario for linear models in Appendix A.1. When Proposition 1 holds, we consider the following reformulation of the problem:

**Reformulated Problem.** Find a transformation  $\mathcal{G}(X; \phi^*)$ , such that

$$\phi^* = \arg\min_{\phi} \|f(X;\theta^*) - Y\|, \text{ where } \theta^* = \arg\min_{\theta} \|f(\mathcal{G}(X;\phi);\theta) - \mathcal{G}(Y;\phi)\|.$$

3.3 Loss Surface Transfer in the Reformulated Problem

Our study seeks a validated transformation  $\tilde{X} = \mathcal{G}(X; \phi)$  that enhances the predictor  $f(\cdot; \theta)$ 's generalization and robustness to varying mapping complexity in (X, Y). Searching for arbitrary transformations can be tedious, and we have no clue about what loss to use to guide such transformation. To address this, we constrain the transformation to be "somewhere near" the raw data, leading to a co-optimization process: learn prediction on the transformed data and align it with the raw series.

We use stochastic gradient descent on a combination of proximity and prediction losses. The *proximity loss* measures the difference between the transformed and raw data, while the *prediction loss* assesses prediction accuracy:

$$\mathcal{L}_{pareto} = \mathcal{L}_{prox} + \mathcal{L}_{pred} = \left( \left\| \tilde{X} - X \right\| + \left\| \tilde{Y} - Y \right\| \right) + \left\| f(\tilde{X}; \theta) - \tilde{Y} \right\| \quad \text{s.t.} \quad \nabla f \le \alpha, \quad (3)$$

where  $\tilde{X} = \mathcal{G}(X;\phi)$  is transformed from the raw data via a network parameterized by  $\phi$ , and  $\nabla f \leq \alpha$  is a convergence constraint on the predictor. When applying to a linear model, since the optimization of that model is convex, our constraint guarantees reaching the Pareto frontier.

To provide an intuitive demonstration, we employ the Proximal Transformation Network (PTN) (detailed soon) for the transformation and DLINEAR (Zeng et al., 2023) as the predictor that is



(a) Overview of surface transfer (b) Improved prediction on raw data (c) Pareto frontier w.r.t. surface

228 Figure 2: PTN with DLINEAR (colored curves) vs vanilla DLINEAR (black curves). The z-axis 229  $(l_raw)$  refers to the forecasting accuracy measured by MSE on the raw test dataset; the x-axis 230  $(l_{prox})$  and y-axis  $(l_{pred})$  are  $\mathcal{L}_{prox}$  and  $\mathcal{L}_{pred}$ , respectively, for the loss curve on our PTN. For 231 the vanilla DLINEAR,  $l_pred$  refers to the MSE on the training set, and  $l_prox$  equals zero as it 232 is trained on the raw data. (a) Loss surface transfer achieved by our Proximal Transformation; (b) 233 MSE loss curve measured on the raw test set indicating that the best parameters achieved by PTN outperform the same parameters trained on the raw data; (c) Pareto frontier formed during training, 234 illustrating the trade-off between proximity and prediction losses. 235

236

227

trained on the transformed data. The training process is showcased on the ETTh1 dataset with lookback length and prediction horizon both set to 96. Figure 2(a) illustrates that the training on the raw data occurs on the plane where the proximity loss  $l_{prox} = 0$  (gray plane), while training on transformed data happens on a curved surface (cyan surface). Intuitively, what our method does is to transfer the original plane to this curved surface (indicated by the red arrow).

242 In the xOz projection (Figure 2(b)), the minimum loss on the curved surface is clearly lower than on 243 the original plane. This surface, achieved by the co-optimization of proximity and prediction losses, 244 balances both objectives by ideally attaining the optimal prediction loss at each point. Hence, it 245 can be regarded as representing the Pareto frontier, as shown in the xOy projection (Figure 2(c)). Assuming the intersection of the two loss curves (shown as the dotted lines in Figures 2(b) and 246 (c)), we find the endpoint of this frontier an equivalent point to the raw-data-guided training. By 247 expanding the optimization from the  $l_{prox} = 0$  plane to the Pareto frontier's curved surface, we 248 discover a superior transformation. Once the global minimum is found on the validation set, we 249 fix the optimized parameters and apply them to the test set, following standard machine learning 250 practice. By analyzing the curves and surfaces, we demonstrate that our method not only maintains 251 fidelity to the raw data but also leverages the transformation to achieve enhanced performance.

252 253 254

255 256

257

258

259

## 4 PROXIMAL TRANSFORMATION NETWORK

Given the preliminaries in Section 3.2, we decompose our design into two parts: the Proximal Transformation Network (PTN) that gradually approaches the raw series via a learnable Encoder-Decoder transformation, and an arbitrary predictor model that is orthogonal to PTN. We further introduce a Mixture-of-Experts design within PTN, enabling the cooperative transformation of data patches and providing a new approach to scaling. The full architecture is shown in Figure 3.

260 261 262

263

## 4.1 ENCODER-DECODER ARCHITECTURE

**Encoder** Typical smoothing methods mentioned in Section 3.1 all compute the weighted sum within a given moving kernel, similar to a convolution layer. We then consider an Encoder (Figure 3(b)) consisting of a convolutional network that is deeper than the smoothing methods with at most two layers. Inspired by MODERNTCN (Luo & Wang, 2024), we consider the Effective Receptive Field (ERF) an important factor in designing our module. Instead of using an extra-large kernel, we enlarge the number of kernels in the convolution layer and unfold the  $(K^i \times L^i)$ -sized output to a sequence of embeddings of length  $L = L^i \times K^i$ , where  $K^i$  is the number of kernels of the *i*-th

284

285

286

287

288 289

302

303 304 305

306 307



Figure 3: Architecture of PTN. The raw data X and Y are first fed into our proposed module PTN to generate transformed data  $\tilde{X}$  and  $\tilde{Y}$  (shown in blue arrows), then the Predictor is trained on the transformed data (shown in the green arrow). (a) Overview design with MoE; (b) Convolution Encoder generates embedding by concatenating outputs from various depths; (c) Intra-patch Attention computes attention on embedding with subseries; (d) Channel-wise Attention computes attention on specific channels; (e) Point-wise Linear Head shares parameters for all embeddings.

convolution layer,  $L^i$  is the length of output from the *i*-th layer, and L is the original input length. This includes a transpose-unfold operation, shown in detail in Figure 12 in Appendix A.4.5.

By setting the 1D convolution with kernel size = 3, stride = 1, and padding = 1, we can halve the input length at each layer. Doubling the number of kernels for each subsequent layer ensures that the receptive field width at the *i*-th layer is  $(2^{i+1} - 1)$ . This design can help generate a larger ERF as discussed in Appendix A.4.5. We believe that our design is more expressive and can serve as a uniform transformation initializer. The case study in Section 5.2 demonstrates that this module adaptively aligns the data distributions of variates in time series.

Decoder We introduce two attention mechanisms and a linear projection head (Figures 3(c)-(e)).
 *Intra-patch attention* computes attention scores within patches, while *channel-wise attention* computes attention scores across channels (here "channels" refer to variates instead of image channels).

Intra-patch Attention: 
$$\mathbf{e}_{\mathrm{IP}} = \frac{\operatorname{Softmax}\left(\left(\mathbf{e}Q\right)\left(\mathbf{e}K\right)^{\top}\right)\left(\mathbf{e}V\right)}{\sqrt{d}},$$
 (4)

Channel-wise Attention: 
$$\mathbf{e}_{Ch} = \frac{\text{Softmax}\left(\left(\mathbf{e}Q\right)\left(\mathbf{e}K\right)^{\top}\right)\left(\mathbf{e}V\right)}{\sqrt{d}},$$
 (5)

where e is the output embedding from Convolution Encoder,  $\mathbf{e}_{IP} \in \mathbb{R}^{C \times L \times d}$  and  $\mathbf{e}_{Ch} \in \mathbb{R}^{L \times C \times d}$ are the outputs of intra-patch attention and channel-wise attention, respectively, and *C* denotes the number of variates. We add the two embeddings together as input of the next layer. The parameters of both attention mechanisms are shared for each layer to improve efficiency.

312 A *linear head* with a weight matrix  $W_{\theta}$  of size  $d \times 1$  (see Figure 3(e)) is shared for all data points for decoding, which is necessary for *point-wise decoding* from the latent embedding. Note that our 313 approach differs from the common approach of concatenating all latent embeddings into a single 314 long vector of length  $L \times d$ , which would require a large head of size  $Ld \times L$ . This is reasonable 315 as the whole PTN module is designed to transform data and is not required to capture temporal 316 dependencies for forecasting. *Point-wise decoding* can be also reinforced by our special design of 317 Convolution Encoder, which ensures each embedding piece has a wide enough Effective Receptive 318 Field (discussions extended in Appendices A.4.5 and A.4.6). Even with limited parameters, this 319 module can explicitly decode the embedding, serving as the initial transformation for interpretability 320 purposes. E.g., we can apply the same linear head to the latent embedding from Convolution Encoder 321 to visualize what series is taken as the initialized transformation in PTN, as shown in Section 5.2. 322

**Training Loss** Using the Encoder-Decoder architecture, we obtain the transformed data  $\hat{X}$  and  $\hat{Y}$  in Equation 3. When paired with an arbitrary predictor, the model can be trained using this equation.

6

For more complex architectures, the gradient constraint in Equation 3 may hinder convergence, so we adopt a simplified Pareto Loss without it. Experimental results in Table 3 demonstrate that PTN is also effective when applied to Transformer architectures (Nie et al., 2023; Liu et al., 2024b).

327 328

330

## 4.2 MIXTURE OF EXPERTS FOR PATCHED SUBSERIES

One major concern with Time Series Transformers and other deep-learning-based TSF methods is the scaling law. Nie et al. (2023); Liu et al. (2024b) have shown that increasing the number of model parameters does not consistently improve performance. More often than not, these Transformers achieve their best results with as few as two layers.

We thus shift our approach to scaling by incorporating Mixture of Experts (MoE) into TSF. Specifically, we combine MoE with the patching operation, a common practice adopted by (Wang et al., 2024b; Zhou et al., 2023; Yu et al., 2023; Lin et al., 2023). In our design, we route patches to different expert models, each being a separate PTN with independent parameters, to enhance both performance and efficiency. Patches are routed to each expert based on their sequential order, as depicted in Figure 3(a). Detailed performance boost and a decomposition view of MoE are provided in Appendix A.5.

342

## 343 344

345

350

351

352

353

354

355

## 5 EXPERIMENTS

We evaluate the performance of our proposed PTN using seven widely recognized datasets: Electricity, Weather, Traffic, and four ETT datasets (ETTh1, ETTh2, ETTm1, ETTm2). These datasets are
well-established benchmarks in the field and are publicly accessible (Wu et al., 2021).

We utilize ITRANSFORMER (Liu et al., 2024b), PATCHTST (Nie et al., 2023), TIMESNET (Wu et al., 2023), SPARSETSF (Lin et al., 2024), DLINEAR (Zeng et al., 2023), and FEDFORMER (Zhou et al., 2022) as our baselines, encompassing temporal Transformers, channel-wise Transformers, convolutional models, and linear models. Following previous settings (Zhou et al., 2021; Wu et al., 2021; 2023) for direct comparison, we set the forecasting horizon  $H \in \{96, 192, 336, 720\}$  and look-back length  $L \in \{96, 512\}$  for all datasets. We move the dataset descriptions, implementation details, and instructions for reproduction to Appendix A.3 and make the codebase accessible at https://anonymous.4open.science/r/PTN-2FC6/.

360

## 5.1 MAIN RESULTS

We address two main questions: Q1: Can PTN achieve SOTA results when paired with mainstream
 backbones without complex designs? Q2: Can PTN be applied to both linear models and non-linear
 Transformers to consistently improve performance?

364 For Q1, we use two representative forecasting models as backbones for PTN: RLINEAR (Li et al., 365 2023), a linear model only utilizing Instance Norm (Kim et al., 2021), and ITRANSFORMER (Liu 366 et al., 2024b), a full Transformer architecture that inverts the feature dimension to compute attention 367 scores across channels. We denote the combinations as PTN-RLI and PTN-ITR, respectively. Table 2 368 summarizes the average results for a look-back length of 96. Our models outperform others in 12 out of 14 evaluation metrics across seven datasets. Full results for look-back lengths of 96 and 369 512 are in Tables 11 and 12 in Appendix A.6. PTN-RLI leads on the four ETT datasets, where 370 linear models typically excel, while PTN-ITR performs best on the other three datasets, aligning 371 with ITRANSFORMER'S SOTA results on large, multivariate datasets. 372

For Q2, we illustrate the performance boost from applying PTN to three linear models, including a
basic single-layer linear model denoted as LINEAR. We also include two Transformer-based models:
ITRANSFORMER (channel-wise) and PATCHTST (channel-independent). The results are presented
in Table 3. Our PTN plugin consistently improves all five backbone models. For a simple LINEAR
model enhanced with PTN, performance can even match SOTA baselines on some datasets. We omit
the results for the other three ETT datasets and place their results in Table 13 in Appendix A.6.

Models	PTN-	RLI	PTN	ITR	ITRANS	FORMER	PATC	HTST	TIME	SNET	SPARS	SETSF	DLn	NEAR	FEDFC	ORMER
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	0.437	0.427	0.453	0.451	0.454	0.448	0.469	0.455	0.458	0.450	0.445	0.425	0.446	0.434	0.440	0.460
ETTh2	0.366	0.391	0.441	0.442	0.383	0.407	0.387	0.407	0.414	0.427	0.386	0.402	0.374	0.399	0.437	0.449
ETTm1	0.399	0.390	0.394	0.396	0.407	0.410	0.387	0.400	0.400	0.406	0.414	0.392	0.414	0.408	0.448	0.452
ETTm2	0.272	0.319	0.293	0.343	0.288	0.332	0.281	0.326	0.291	0.333	0.287	0.326	0.286	0.327	0.305	0.349
Electricty	0.215	0.289	0.164	0.253	0.178	0.270	0.205	0.290	0.193	0.295	0.221	0.291	0.219	0.298	0.214	0.327
Traffic	0.643	0.355	0.463	0.266	0.428	0.282	0.481	0.300	0.620	0.336	0.676	0.361	0.627	0.378	0.610	0.376
Weather	0.268	0.283	0.235	0.271	0.258	0.278	0.259	0.273	0.259	0.287	0.297	0.305	0.272	0.291	0.309	0.360

Table 2: Average forecasting results of PTN-RLI and (PTN with RLINEAR and ITRANSFORMER backbones). The best and the second-best measures are in **bold** and underlined, respectively.

Table 3: Performance boost by adding PTN to different backbones with better results in **bold**.

Models	LINEAR +PTN		TN	DLP	NEAR	+P	TN	RLIP	NEAR	+P	TN	ITRAN	SFORMER	+P	TN	PATC	HTST	+P	TN	
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	0.400	0.398	0.394	0.389	0.414	0.408	0.401	0.390	0.405	0.390	0.399	0.390	0.407	0.410	0.394	0.396	0.387	0.400	0.381	0.386
Electricity	0.215	0.296	0.212	0.289	0.219	0.298	0.215	0.289	0.217	0.289	0.215	0.289	0.178	0.270	0.164	0.253	0.205	0.290	0.196	0.274
Traffic	0.626	0.363	0.675	0.352	0.627	0.378	0.643	0.355	0.672	0.354	0.643	0.355	0.428	0.282	0.463	0.266	0.481	0.300	0.584	0.294
Weather	0.265	0.295	0.263	0.292	0.272	0.291	0.268	0.283	0.284	0.287	0.268	0.283	0.258	0.278	0.235	0.271	0.259	0.273	0.246	0.272

### 5.2 INTERPRETATION OF PROXIMAL TRANSFORMATION

We shed some light on how and why the Proximal Transformation works. First, we show a selfsupervised clustering pattern in the extended loss space for the transformed data. This pattern can be correlated to the predictability of time series, supporting the significant performance boosts from PTN. Second, we examine the training process of PTN, where the transformed data gradually approaches the raw data. In this process, the transformation adapts to channel-specific features related to data sparsity. Furthermore, we visualize the embeddings from our Convolution Encoder and observe that this design aligns distributions across channels, effectively normalizing the data.





A Self-supervised Learner for Assessing Predictibility When looking into the average loss dis-tribution in the same setting as in Section 3.3, we can observe how data samples (points) are adap-tively clustered into different distributions. Figure 4(a) shows that the average test MSE (" $l_{raw}$ " on the z-axis) forms a clustering pattern for different data samples. The samples within the prox-*imal area* (blue, lower in proximity losses) have lower prediction losses on average, leaving the rest less predictable samples outside the area. This is also supported by the distribution of variance per sample (a data sample's variance is the variance of the series it belongs to), as shown in Figure 4(b). As variance is also a commonly considered metric w.r.t. predictability (lower variance for more predictable and higher variance for less predictable), we suggest that our method enables self-supervised estimates on the predictability and based on the estimates to constrain a proximal area for more predictable time series. We also analyze the distribution of performance boost, mea-

sured by loss contribution (average loss × number of samples) changed by PTN (with vs. without) in
Figure 4. Performance boosts are mainly found in a limited range of proximity, supporting our assumption that effective data transformation exists within the proximal area close to the raw data. In
essence, our PTN can categorize time series into predictable (low variance) and unpredictable (high
variance) groups via a self-supervised manner, focusing on enhancing performance for the former.

Transformation as Channel-adaptive We analyze two representative channels from the ETTh2 dataset to illustrate Proximal Transformation. Figures 5(a) and (b) show the transformed series (in green) approximating the original (in blue). For Channel 1, the transformation evolves from over-smoothed to the original, while for Channel 5, it transitions from noisy to the original.

When comparing the distributions measured by the normalized values of both channels, it becomes evident that even after preprocessing with RevIN (Kim et al., 2021), the distributions of both channels are not aligned (see Figure 5(c)). We then apply our point-wise linear head to decode the output embeddings from the Convolution Encoder, resulting in the initially transformed data (in yellow). The distributions of both channels align (see Figure 5(d)), supporting the commonly adopted channel-independent strategy of sharing weights across channels. We provide more visualizations in Appendix A.7.2.



Figure 5: (a) and (b) Visualized raw and transformed series, respectively, in training PTN for Channels 1 and 5 in ETTh2 dataset; "raw" (blue) denotes the original series, "embedding" (yellow) the decoded series from the Convolution Encoder (Section 4.1); and transformed (green) the outputs of PTN. (c) The distributions of normalized raw series of Channels 1 and 5. (d) The distributions of decoded series of Channels 1 and 5.

5.3 TRANSFERABILITY OF THE TRANSFORMED DATA

Generalization capability is a commonly considered issue in machine learning. This topic is par-ticularly important in TSF, as the diversity of datasets across various domains presents significant challenges in model design. Transfer learning adapts models to new datasets to enhance performance or reduce training costs. Our reformulating of finding a transformation of data offers a novel per-spective on transfer learning. Specifically, our PTN design facilitates the combination of transformed data with various forecasting models. A simple, lightweight linear model can initially determine the transformation, which is then fixed and transferred to other forecasting models. Table 4 demonstrates the transfer learning results, showcasing the feasibility of our method. The data transfer protocol broadens conventional transfer learning by moving from generalizing a model for different datasets to its dual form of generalizing the data for different models. 

5.4 ABLATION STUDY ON THE PTN DESIGNS

**Encoder** We propose using a CNN-based encoder (see Section 4.1) to generate an initial transformation of the raw time series. To evaluate its effectiveness, we analyze the following three variants: 1)

Model	RL	NEAR -	PATCH	ГST		PATC	HTST		RLIN	$EAR \rightarrow I$	<b>F</b> RANSFO	ORMER		ITRANS	FORMER	
Horizon	96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720
ETTh1	0.410	0.427	0.456	0.505	0.419	0.445	0.466	0.488	0.409	0.430	0.453	0.513	0.405	0.436	0.458	0.491
ETTh2	0.349	0.402	0.443	0.526	0.348	0.400	0.433	0.446	0.400	0.403	0.450	0.619	0.349	0.400	0.432	0.445
Electricity	0.247	0.256	0.272	0.304	0.270	0.274	0.293	0.324	0.233	0.248	0.262	0.277	0.240	0.253	0.269	0.317
Traffic	0.274	0.277	0.283	0.302	0.290	0.290	0.300	0.320	0.242	0.253	0.264	0.281	0.268	0.276	0.283	0.302

Table 4: Comparison of generalization capabilities of our data transfer protocol in terms of MAE: 'A  $\rightarrow$  B' indicates the raw data transformation is learned by Model A, and then fed to Model B.

Table 5: Experiments results of Ablation Study on our Convolution Encoder under different settings.

	w. Trans.	w. Trans	Embedding	ET	Th1	ET	Th2	ET	Tm1	ET	Гm2
	on X	on $Y$	Туре	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
PatchTrans	√	√	Patch	0.444	0.436	0.471	0.452	0.399	0.390	0.279	0.329
ConvPred		$\checkmark$	Conv	0.490	0.463	0.529	0.485	0.389	0.393	0.343	0.372
SingleTrans	$\checkmark$		Conv	0.442	0.427	0.380	0.416	0.394	0.389	0.294	0.335
ConvTrans (ours)	√	$\checkmark$	Conv	0.437	0.427	0.366	0.391	0.386	0.390	0.272	0.319

PatchTrans: A common encoder alternative is patching. We perform experiments where the size of MoE is 1. Our design, named "ConvTrans", uses the CNN encoder, while the variant, "PatchTrans", uses the patching encoder. 2) ConvPred: To test the encoder's effectiveness in a Transformer-based model, we replace the embedding module in PATCHTST with our CNN-based encoder. In fact, using the same decoder to predict is equivalent to approximating the transformation only on Y. 3) SingleTrans: As analyzed in Section 3.1, we verify the need to transform both X and Y. This variant only transforms X. The configurations of all ablation models are summarized in Table 5. Our findings show that our setting consistently outperforms the others.

Decoder Regarding decoder design, we primarily focus on two attention modules: Intra-Patch (IP) attention and Channel-wise (Ch) attention. The results are presented in Figure 6. As it demonstrates, combining both modules does not necessarily yield the best results across all datasets and backbones. However, datasets with more timestamps (e.g., ETTm1) benefit from IP attention, while datasets with more channels (e.g., Electricity) benefit from Ch attention. Therefore, we incorporate both modules to address these varying needs. This also frees us from excessive hyperparameter searches, allowing for more general considerations on different datasets and models.



(a) PTN-RLI on ETTm1 (b) PTN-RLI on Electricity (c) PTN-ITR on ETTm1 (d) PTN-ITR on Electricity

Figure 6: Ablation Study results of the Decoder for PTN-RLI and PTN-ITR (metric: MSE).

## 

## 6 CONCLUSION AND FUTURE WORK

In this study, we approach time series forecasting (TSF) from a data-centric perspective by designing a co-optimization process that balances data transformation proximity and forecasting accuracy. Our design boosts the performance of a variety of TSF models and demonstrates the potential to make TSF more interpretable. For future work, it is relevant to align datasets from diverse sources using our data transformation method. Moreover, it remains open if PTNcan enhance performance in other time series tasks like anomaly detection, classification, and imputation. Further exploration into different backbone architectures is also interesting. We anticipate that a data-centric approach to generalizing datasets, tasks, and architectures could ideally pave the way for enhanced Time Series Foundation Models (Liang et al., 2024).

# 540 REFERENCES

547

553

574

575

576

581

582

583

- Samit Bhanja and Abhishek Das. Impact of data normalization on deep neural network for time se ries forecasting. In *Proceedings of conference on Advancement in Computation, Communication and Electronics Paradigm (ACCEP-2019)*, pp. 27, 2019.
- Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
   time-series forecasting. In *Forty-first International Conference on Machine Learning*, 2024.
- Jinliang Deng, Xiusi Chen, Renhe Jiang, Xuan Song, and Ivor W Tsang. St-norm: Spatial and
   temporal normalization for multi-variate time series forecasting. In *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*, 2021.
- Haoyi Fan, Fengbin Zhang, and Yue Gao. Self-supervised time series representation learning by inter-intra relational reasoning. *arXiv preprint arXiv:2011.13548*, 2020.
- Wei Fan, Pengyang Wang, Dongkun Wang, Dongjie Wang, Yuanchun Zhou, and Yanjie Fu. Dish-ts:
  a general paradigm for alleviating distribution shift in time series forecasting. In *Proceedings* of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence, 2023.
- Lu Han, Han-Jia Ye, and De-Chuan Zhan. Sin: Selective and interpretable normalization for long term time series forecasting. In *Forty-first International Conference on Machine Learning*, 2024.
- Romain Ilbert, Ambroise Odonnat, Vasilii Feofanov, Aladin Virmaux, Giuseppe Paolo, Themis Palpanas, and Ievgen Redko. Samformer: Unlocking the potential of transformers in time series forecasting with sharpness-aware minimization and channel-wise attention. In *Forty-first International Conference on Machine Learning*, 2024.
- Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2021.
- Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Seunghan Lee, Taeyoung Park, and Kibok Lee. Learning to embed time series patches independently. *arXiv preprint arXiv:2312.16427*, 2023.
  - Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting long-term time series forecasting: An investigation on linear mapping. *arXiv preprint arXiv:2305.10721*, 2023.
- 577 Yuxuan Liang, Haomin Wen, Yuqi Nie, Yushan Jiang, Ming Jin, Dongjin Song, Shirui Pan, and
  578 Qingsong Wen. Foundation models for time series analysis: A tutorial and survey. In *Proceedings*579 of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 6555–
  580 6565, 2024.
  - Shengsheng Lin, Weiwei Lin, Wentai Wu, Feiyu Zhao, Ruichao Mo, and Haotong Zhang. Segrnn: Segment recurrent neural network for long-term time series forecasting. *arXiv preprint arXiv:2308.11200*, 2023.
- Shengsheng Lin, Weiwei Lin, Wentai Wu, Haojun Chen, and Junjie Yang. Sparsetsf: Modeling
   long-term time series forecasting with\* 1k\* parameters. In *Forty-first International Conference on Machine Learning*, 2024.
- Hanwen Liu, Yibing Zhang, Ximeng Wang, Bin Wang, and Yanwei Yu. St-moe: Spatio-temporal mixture of experts for multivariate time series forecasting. In 2023 18th International Conference on Intelligent Systems and Knowledge Engineering (ISKE). IEEE, 2023.
- Xu Liu, Junfeng Hu, Yuan Li, Shizhe Diao, Yuxuan Liang, Bryan Hooi, and Roger Zimmermann.
   Unitime: A language-empowered unified model for cross-domain time series forecasting. In *Proceedings of the ACM on Web Conference 2024*, 2024a.

611

618

594	Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring
595	the stationarity in time series forecasting. Advances in Neural Information Processing Systems.
596	2022.
597	

- Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.
   itransformer: Inverted transformers are effective for time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024b.
- <sup>601</sup>
   <sup>602</sup>
   <sup>603</sup>
   <sup>603</sup>
   <sup>604</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>606</sup>
   <sup>606</sup>
   <sup>607</sup>
   <sup>607</sup>
   <sup>608</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>601</sup>
   <sup>601</sup>
   <sup>601</sup>
   <sup>602</sup>
   <sup>603</sup>
   <sup>603</sup>
   <sup>604</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>606</sup>
   <sup>607</sup>
   <sup>607</sup>
   <sup>608</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>601</sup>
   <sup>601</sup>
   <sup>601</sup>
   <sup>602</sup>
   <sup>602</sup>
   <sup>603</sup>
   <sup>603</sup>
   <sup>604</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>606</sup>
   <sup>607</sup>
   <sup>607</sup>
   <sup>608</sup>
   <sup>608</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>601</sup>
   <sup>602</sup>
   <sup>602</sup>
   <sup>603</sup>
   <sup>603</sup>
   <sup>603</sup>
   <sup>604</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>606</sup>
   <sup>607</sup>
   <sup>607</sup>
   <sup>608</sup>
   <sup>608</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>601</sup>
   <sup>602</sup>
   <sup>602</sup>
   <sup>603</sup>
   <sup>603</sup>
   <sup>603</sup>
   <sup>604</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>605</sup>
   <sup>606</sup>
   <sup>607</sup>
   <sup>607</sup>
   <sup>608</sup>
   <sup>608</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>609</sup>
   <sup>601</sup>
   <sup>602</sup>
- Donghao Luo and Xue Wang. Moderntcn: A modern pure convolution structure for general time series analysis. In *The Twelfth International Conference on Learning Representations*, 2024.
- Dongsheng Luo, Wei Cheng, Yingheng Wang, Dongkuan Xu, Jingchao Ni, Wenchao Yu, Xuchao Zhang, Yanchi Liu, Yuncong Chen, Haifeng Chen, et al. Time series contrastive learning with information-aware augmentations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023.
- Ronghao Ni, Zinan Lin, Shuaiqi Wang, and Giulia Fanti. Mixture-of-linear-experts for long-term
   time series forecasting. In *International Conference on Artificial Intelligence and Statistics*.
   PMLR, 2024.
- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In *The Eleventh International Conference on Learning Representations*, 2023.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
   Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Kihao Piao, Zheng Chen, Taichi Murayama, Yasuko Matsubara, and Yasushi Sakurai. Fredformer:
   Frequency debiased transformer for time series forecasting. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2024.
- David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic fore casting with autoregressive recurrent networks. *International journal of forecasting*, 36(3):1181–
   1191, 2020.
- Marcus AK September, Francesco Sanna Passino, Leonie Goldmann, and Anton Hinel. Extended deep adaptive input normalization for preprocessing time series data for neural networks. In *AISTATS*, 2024.
- Xiaoming Shi, Shiyu Wang, Yuqi Nie, Dianqi Li, Zhou Ye, Qingsong Wen, and Ming Jin. Time Billion-scale time series foundation models with mixture of experts. *arXiv preprint arXiv:2409.16040*, 2024.
- Mingtian Tan, Mike A Merrill, Vinayak Gupta, Tim Althoff, and Thomas Hartvigsen. Are language
   models actually useful for time series forecasting? *arXiv preprint arXiv:2406.16964*, 2024.
- Peiwang Tang and Weitai Zhang. Pdmlp: Patch-based decomposed mlp for long-term time series
   forecasting. *arXiv preprint arXiv:2405.13575*, 2024.
- Sean J Taylor and Benjamin Letham. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- 647 Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.

648 William Toner and Luke Nicholas Darlow. An analysis of linear time series forecasting models. In 649 Forty-first International Conference on Machine Learning, 2024. 650 651 Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. Advances in neural information processing systems, 30, 2017. 652 653 Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y Zhang, and 654 JUN ZHOU. Timemixer: Decomposable multiscale mixing for time series forecasting. In The 655 Twelfth International Conference on Learning Representations, 2024a. 656 657 Xue Wang, Tian Zhou, Qingsong Wen, Jinyang Gao, Bolin Ding, and Rong Jin. Card: Channel 658 aligned robust blend transformer for time series forecasting. In The Twelfth International Confer-659 ence on Learning Representations, 2024b. 660 Yuxuan Wang, Haixu Wu, Jiaxiang Dong, Yong Liu, Yunzhong Qiu, Haoran Zhang, Jianmin Wang, 661 and Mingsheng Long. Timexer: Empowering transformers for time series forecasting with ex-662 ogenous variables. arXiv preprint arXiv:2402.19072, 2024c. 663 664 Qingsong Wen, Liang Sun, Fan Yang, Xiaomin Song, Jingkun Gao, Xue Wang, and Huan Xu. Time 665 series data augmentation for deep learning: A survey. arXiv preprint arXiv:2002.12478, 2020. 666 667 Gerald Woo, Chenghao Liu, Doyen Sahoo, Akshat Kumar, and Steven Hoi. Cost: Contrastive learning of disentangled seasonal-trend representations for time series forecasting. In International 668 Conference on Learning Representations, 2022. 669 670 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans-671 formers with auto-correlation for long-term series forecasting. Advances in neural information 672 processing systems, 2021. 673 674 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: 675 Temporal 2d-variation modeling for general time series analysis. In The Eleventh International Conference on Learning Representations, 2023. 676 677 Zhijian Xu, Ailing Zeng, and Qiang Xu. Fits: Modeling time series with 10k parameters. In The 678 Twelfth International Conference on Learning Representations, 2024. 679 680 Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian, Long-681 bing Cao, and Zhendong Niu. Frequency-domain mlps are more effective learners in time series 682 forecasting. Advances in Neural Information Processing Systems, 36, 2024. 683 Chengqing Yu, Fei Wang, Zezhi Shao, Tao Sun, Lin Wu, and Yongjun Xu. Dsformer: A double 684 sampling transformer for multivariate time series long-term prediction. In Proceedings of the 685 32nd ACM international conference on information and knowledge management, pp. 3062–3072, 686 2023. 687 688 Zhihan Yue, Yujing Wang, Juanyong Duan, Tianmeng Yang, Congrui Huang, Yunhai Tong, and 689 Bixiong Xu. Ts2vec: Towards universal representation of time series. In Proceedings of the AAAI 690 Conference on Artificial Intelligence, 2022. 691 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series 692 forecasting? In Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence 693 and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth 694 Symposium on Educational Advances in Artificial Intelligence, pp. 11121–11128, 2023. 695 696 Yitian Zhang, Liheng Ma, Soumyasundar Pal, Yingxue Zhang, and Mark Coates. Multi-resolution 697 time-series transformer for long-term forecasting. In International Conference on Artificial Intel-698 ligence and Statistics. PMLR, 2024. 699 Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency 700 for multivariate time series forecasting. In The Eleventh International Conference on Learning 701 Representations, 2023.

702 703 704	Lifan Zhao and Yanyan Shen. Rethinking channel dependence for multivariate time series fore- casting: Learning from leading indicators. In <i>The Twelfth International Conference on Learning</i> <i>Representations</i> , 2024.
705 706 707 708	Xu Zheng, Tianchun Wang, Wei Cheng, Aitian Ma, Haifeng Chen, Mo Sha, and Dongsheng Luo. Parametric augmentation for time series contrastive learning. In <i>The Twelfth International Con-</i> <i>ference on Learning Representations</i> , 2024.
709 710 711	Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In <i>Proceedings</i> of the AAAI conference on artificial intelligence, 2021.
712 713 714 715	Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In <i>International Conference on Machine Learning</i> . PMLR, 2022.
716 717 718	Tian Zhou, Peisong Niu, Liang Sun, Rong Jin, et al. One fits all: Power general time series analysis by pretrained lm. <i>Advances in neural information processing systems</i> , 36:43322–43355, 2023.
719 720 721	
722 723 724	
725 726 727 728	
729 730 731	
732 733 734	
735 736 737	
738 739 740 741	
742 743 744	
745 746 747	
748 749 750	
751 752 753	
754 755	

#### APPENDIX А

#### A.1 DISCUSSION ON PROPOSITION 1

We are going to briefly discuss Proposition 1 when  $X_{train}$  and  $X_{test}$  are identically distributed (same for  $Y_{train}$  and  $Y_{test}$ , therefore we no longer use different subscripts to distinguish them). Specifically, we briefly discuss a very simple scenario that X and Y satisfies Y = WX, where  $X \in \mathbb{R}^{L \times N}$  and  $Y \in \mathbb{R}^{H \times N}$ . And we have noised observations  $Y = WX + \epsilon$  that parameters are trained on, where  $\epsilon \in \mathbb{R}^{H \times N}$  and each elements in  $\epsilon$  is sampled from  $\mathcal{N}(\mu, \sigma)$ 

**Lemma 1** Suppose we have a linear predictor  $f(x; W_{\theta}) = W_{\theta}x$ , the closed-form solution for  $\min_{W} ||W_{\theta}X - (Y + \epsilon)||_2 \text{ is } W_{\theta}^* = (Y + \epsilon)X^{\top}(XX^{\top})^{-1}.$ 

*Proof.*  $W^*_{\theta}$  can be solved by the following:

$$\frac{\mathrm{d}\left(W_{\theta}^{*}X - (Y + \epsilon)\right)^{2}}{W_{\theta}^{*}} = 0, \tag{6}$$

because

$$\frac{\mathrm{d}\left(W_{\theta}^{*}X - (Y+\epsilon)\right)^{2}}{W_{\theta}^{*}} = 2\left(W_{\theta}^{*}X - (Y+\epsilon)\right)X^{\top}.$$
(7)

We then solve Equation 6:

$$2 \left( W_{\theta}^* X - (Y + \epsilon) \right) X^{\top} = 0$$
  

$$\Leftrightarrow W_{\theta}^* X X^{\top} - (Y + \epsilon) X^{\top} = 0$$
  

$$\Leftrightarrow W_{\theta}^* = (Y + \epsilon) X^{\top} (X X^{\top})^{-1}$$
(8)

This follows the typical form of a linear regression problem, and this closed-form solution always exists because of its convex nature.

**Theorem 1** We have  $\tilde{W}_{\theta} = W$  that  $\mathbb{E}||\tilde{W}_{\theta}X - (Y + \epsilon)||_2 < \mathbb{E}||W_{\theta}^*X - (Y + \epsilon)||_2$ , when  $XX^{\top}$  is slightly ill-conditioned that its condition number  $\kappa(XX^{\top}) = ||XX^{\top}||_2 \cdot ||(XX^{\top})^{-1}||_2 > 2.$ 

Proof.

$$\mathbb{E}||\tilde{W}_{\theta}X - (Y+\epsilon)||_{2} = \mathbb{E}||WX - (Y+\epsilon)||_{2}$$
$$= \mathbb{E}||Y - Y - \epsilon||_{2}$$
$$= \mathbb{E}||\epsilon||_{2}$$
(9)

$$\mathbb{E} \|W_{\theta}^{*} - (Y+\epsilon)\|_{2} = \mathbb{E} \|(Y+\epsilon)X^{\top}(XX^{\top})^{-1}X - (Y+\epsilon)\|_{2}$$

$$= \mathbb{E} \|(WX+\epsilon)X^{\top}(XX^{\top})^{-1}X - (Y+\epsilon)\|_{2}$$

$$= \mathbb{E} \|W\left[XX^{\top}(XX^{\top})^{-1}\right]X + \epsilon X^{\top}(XX^{\top})^{-1}X - (Y+\epsilon)\|_{2}$$

$$= \mathbb{E} \|\epsilon\left[X^{\top}(XX^{\top})^{-1}X - I\right]\|_{2}$$

$$= \mathbb{E} \frac{\|\epsilon\left[X^{\top}(XX^{\top})^{-1}XX^{\top} - X^{\top}\right]\|_{2}}{\|X^{\top}\|_{2}}$$

$$\geq \mathbb{E} \frac{\|\epsilon\|_{2} \|X^{\top}\|_{2} \left(\left\|(XX^{\top})^{-1}XX^{\top}\|_{2} - 1\right)\right)}{\|X^{\top}\|_{2}}$$

$$= \mathbb{E} \|\epsilon\|_{2} \cdot \left(\left\|(XX^{\top})^{-1}XX^{\top}\|_{2} - 1\right)\right)$$

$$= \mathbb{E} \|\epsilon\|_{2} \cdot \left(\|(XX^{\top})^{-1}XX^{\top}\|_{2} - 1\right)$$

so we have  $\mathbb{E}\left\|\tilde{W}_{\theta}X - (Y+\epsilon)\right\|_{2} < \mathbb{E}\left\|W_{\theta}^{*}X - (Y+\epsilon)\right\|_{2}$ , when  $\left\|XX^{\top}\right\|_{2} \cdot \left\|\left(XX^{\top}\right)^{-1}\right\|_{2} > \Box$ 2.

ETTh1	ETTh2	ETTm1	ETTm2	electricity	traffic	weather
96   6084.71	24703.02	25834.22	67423.96	4609.54	1710.74	2334938.58
192   12797.42	52742.00	52918.36	147632.44	10073.12	4646.02	4461250.74
336   24097.35	97711.27	94813.64	277571.90	19957.84	12035.87	7681729.05
720   76941.31	280377.34	217174.44	646904.38	53569.14	43433.10	16193817.02

Table 6: Condition numbers for different datasets under the adopted look-back lengths.

the figures here are for reference only and to provide some insights on future research such as why or why not simple linear models can outperform Transformer-based models on certain datasets.

This simple example is meaningful for a variety of linear models covering temporal and frequency domains, as they are equivalent to the closed-formed solution in **Lemma 1** Toner & Darlow (2024). In real scenarios, the raw data is often ill-conditioned as discussed in Toner & Darlow (2024) (shown in Table 6). By the above Theorem, we demonstrate that under very trivial conditions, Proposition 1 often holds when using linear models, even if the train sets and test sets are identically distributed.

A.2 RELATED WORK REVISIT

In this part, we will introduce some popularly adopted methods for TSF, in which our baseline models are included.

Channel-independent Transformers Time Series Transformers (TST) (Zhou et al., 2021; Wu
 et al., 2021; Zhou et al., 2022) have recently led to significant progress in the TSF problem, demonstrating convincing superiority over traditional methods and convolution-based models.

Inspired by Zeng et al. (2023), Nie et al. (2023) incorporates the *Channel Independent* (CI) design
(sharing weights across all channels) to introduce PATCHTST, a new state-of-the-art (SOTA) model
that significantly benefits from CI and the patching operation. More works follow this practice and
achieve excellent performance in TSF (Das et al., 2024; Liu et al., 2024a; Zhang et al., 2024).

**Channel-wise Transformers** Building on this, more recent research has focused on designing Transformers capable of capturing channel dependencies inherent in multivariate time series data. Notable examples include CROSSFORMER (Zhang & Yan, 2023), ITRANSFORMER (Liu et al., 2024b), DSFORMER (Yu et al., 2023), and CARD (Wang et al., 2024b). These models shed light on exploiting inter-series relationships to improve forecasting accuracy. Furthermore, Ilbert et al. (2024) propose that Transformers are inherently unstable when trained on time series datasets due to sharp gradient curvature, and suggest a sharpness-aware optimizer to mitigate such issues. Nonetheless, channel-wise Transformers still suffer from overfitting on small datasets.

**Linear Models** In contrast to the quadratic complexity of Transformers, lightweight linear mod-els have emerged as competitive alternatives offering simplicity and efficiency. RLINEAR (Li et al., 2023) verifies that a vanilla single-linear model, when combined with a widely-adopted normaliza-tion method (Kim et al., 2021), can achieve near SOTA performance in TSF. Further research (Zhao & Shen, 2024) on channel dependencies within linear and MLP-based models has yielded performance improvements over previous CI approaches. Moreover, models (Xu et al., 2024; Yi et al., 2024) that directly learn linear regression or MLP-based models on complex frequency features have achieved remarkable performances. Through theoretical analysis of these linear models, Toner & Darlow (2024) conclude that normalization-enhanced (Li et al., 2023), decomposition-based (Zeng et al., 2023) and frequency-domain linear (Xu et al., 2024) models are essentially equivalent to vanilla linear regression. Most recently, SPARSETSF (Lin et al., 2024), a SOTA lightweight model, incorporates 1D convolutions as down-sampling modules and learns linear parameters on the down-sampled values of the original series.

# A.3 EXPERIMETAL DETAILS

## A.3.1 DATASETS

867

882 883

885

887 888

889

We conduct experiments on 11 real-world datasets to evaluate the performance of the proposed 868 PTN. (1) ETT (Zhou et al., 2021) contains 7 factors of electricity transformers from July 2016 to July 2018. The ETTh1 and ETTh2 are subsets that are recorded every hour, and ETTm1 and 870 ETTm2 are recorded every 15 minutes. (2) Electricity (Wu et al., 2021) records the hourly electricity 871 consumption data of 321 clients. (3) Traffic (Wu et al., 2021) collects hourly road occupancy rates 872 measured by 862 sensors of San Francisco Bay area freeways from January 2015 to December 2016. 873 (4) Weather (Wu et al., 2021) includes 21 meteorological factors collected every 10 minutes from 874 the Weather Station of the Max Planck Biogeochemistry Institute in 2020. (5) PeMS contains the 875 public traffic network data in California collected by 5-minute windows. We use the same four 876 public subsets (PeMS03, PeMS04, PeMS07, PeMS08) adopted in iTransformer (Liu et al., 2024b). 877 For ETT datasets, we divide datasets by ratio  $\{0.6, 0.2, 0.2\}$  into train set, validation set, and test set. For Electricity, Traffic, Weather and PeMS datasets, we divide by ratio  $\{0.7, 0.1, 0.2\}$  in the same 878 setting of TimesNet (Wu et al., 2023) and many previous works. All datasets are scaled by the mean 879 and variance of train sets, a common practice in TSF (Wu et al., 2023). 880

Table 7: Statistics of evaluation datasets.

Datasets	ETTh1	ETTh2	ETTm1	ETTm2	Electricity	Traffic	Weather	PeMS03	PeMS04	PeMS07	PeMS08
# of TS Variates	7	7	7	7	321	862	21	358	307	883	170
TS Length	17420	17420	69680	69680	26304	17544	52696	26209	16992	28224	17856

## A.3.2 BASELINES

890 For horizontal comparisons, we include RLINEAR, DLINEAR, SPARSETSF, ITRANSFORMER, 891 PATCHTST, FEDFORMER, and TIMESNET, covering linear, transformed-based, convolution-based 892 and frequecy-domain methods. (1) RLINEAR (Li et al., 2023), a linear model with Reversible In-893 stance Normalization (Kim et al., 2021); (2) DLINEAR (Zeng et al., 2023), a linear model on the 894 decomposed Trend term and Seasonal term; (3) SPARSETSF (Lin et al., 2024). A linear model using sparse technique; (4) ITRANSFORMER (Liu et al., 2024b), a Transformer-based model the 895 compute attention score on the inverted series (along the channel dimension); (5)PATCHTST (Nie 896 et al., 2023), a Channel-Independent Transformer-based model that uses patching to tokenize; (6) 897 FEDFORMER, a Transformer-based model utilizing frequency-domain information. The Table 7 898 shows the detailed statistics of 7 datasets. It should be pointed out that the scale of a time series 899 dataset is determined by the number of variates and dataset length combined. Therefore a more fair 900 comparison should cover both two factors. In our experiment, we observe that some baseline models 901 like (Liu et al., 2024b) can benefit greatly from a larger number of variates, while some may prefer 902 longer datasets (Nie et al., 2023) or longer look-back lengths (Lin et al., 2024).

903 904

905

A.4 IMPLEMENTATION DETAILS

# 906 A.4.1 Hyperparameters and Settings

All experiments and methods are implemented in Python and PyTorch (Paszke et al., 2019) and conducted on two Nvidia RTX A5000 ada generation GPUs (32GB) and two Nvidia RTX A6000 GPUs (48GB). We use ADAM (Kingma, 2014) and the learning rate initialized as 0.0003 for all settings. Instead of setting an exceptionally small epoch for all experiments as adopted in previous works (Liu et al., 2024b; Lin et al., 2024), we use early stopping on the MSE metric of the validation set with patience set as 10 epochs.

914 The hyperparameters of convolution encoder will be set as fixed as will be explained in Sec915 tion A.4.5. We the size of hidden dimension of attention-based decoder as 32 and use two attention
916 layers for both intra-patch and channel-wise modules. The patch length also referred as the period
917 length in SparseTSF Lin et al. (2024) is set to 4 to support an enfficient implementation. The total number of parameters are 26285, which is a relatively light amount of parameters when counting

static memory, i.e. the memory used to store parameters. For more practical scenes when it comes to time and memory complexity, we are going to analyze in the following.

## A.4.2 OVERALL PIPELINES OF PTN IN DIFFERENT WORKING MODES



Figure 7: The overall pipelines for PTN in different working modes.

As illusrated in Figure 3, the proposed Proximal Transformation Network (PTN) consists of three distinct pipelines: one for training, one for evaluation with PTN, and one for evaluation without PTN.

937 During the training phase, the pipeline involves two transformation processes applied to both X and 938 Y (see blue arrows in Figure 3), followed by the forecasting process, which is parameterized by an 939 arbitrary predictor.

For evaluation (inference), there are two options available. The first option (see the red arrow) is the default evaluation, applicable to any plugin module, where the backbone model (the predictor in this case) is utilized. The second option (see the yellow arrow) provides a simplified approach for evaluation or inference. In this case, the raw data X is fed directly into a *distilled* student model, which has the same architecture and produces outputs trained using  $\tilde{Y}$  as labels. Further details regarding the student model are provided in Appendix A.4.4.

946 947

948

921

922 923 924

925

926

927

928

929

930

931 932

933 934

## A.4.3 COMPLEXITY ANALYSIS OF PTN

The proposed PTN serves as a plugin that can enhance a backbone's performance at the cost of additional computation and memory usage. Hence, we proceed to analyze the computational complexity
of PTN, bridging some of the gaps between theoretical complexity from past studies and our empirical findings. For a brief conclusion, our complexity analysis offers guidance on implementing PTN
efficiently for practical usage. By limiting the size of parallel computations within a batch, we can
reduce the relative multiplier of time and memory costs from a number of times to a factor of less
than one during training and inference.

To be specific, we first identify the actual batch size as a crucial factor in complexity analysis and extend the analysis from previous studies to better align with the empirical results (see P1 below). Second, by reducing actual batch size, we propose a more efficient implementation for the combinational uses of two orthogonal attentions (core components in PTN). Under certain conditions, this method can achieve reasonable speed-up at the cost of memory usage (see P2). Third, we verify the effectiveness of the variate sampling strategy, which is also supported by our analysis (see P3). Both methods here do not incur performance drop with improvement in efficiency.

962 963

**P1.** Complexities of different Attentions. A previous study by Liu et al. (2024b) reveals that 964 the channel-wise and inter-patch attentions have different complexities based on dataset features. 965 Compared to inter-patch attention, channel-wise attention uses more resources on datasets with more 966 variables and less on those with fewer variables. This difference is due to the dimension along which 967 attention is calculated. The complexity of attention is  $\mathcal{O}(N^2)$ , where N represents the number of 968 tokens. For channel-wise attention, N is the number of variables, while for inter-patch attention, N969 is the number of patches. These complexities are represented as  $\mathcal{O}(C^2)$  for channel-wise and  $\mathcal{O}(n^2)$ for inter-patch, where C is the number of variables and n is the number of patches (calculated as 970  $L/l_p$ , with  $l_p$  being the patch length). Intra-patch attention has a constant complexity of  $\mathcal{O}(l_n^2)$ , 971 meaning it doesn't scale with the series length or horizon.

987

990

991

993 994

995

996

997

998

999

1000 1001

1003

1004

1012

972 However, despite not computing along the variable dimension, the times of inter-patch and 973 intra-patch attention computation increase nearly linearly with more variables. Due to channel-974 independence strategy Zeng et al. (2023), the number of variates will multiply the batch size, re-975 sulting in a larger number of parallel computations. Since GPU parallel capacity is limited, a larger 976 batch size does not necessarily speed up computation. What is more, the increased batch size does not lower the number of batches to be computed in a dataset. Therefore we should take actual batch 977 size into consideration for complexity. For large C, inter-patch attention complexity is  $\mathcal{O}(C \times n^2)$ . 978 As shown in Figure 10 and Figure 11, datasets with many variables (like the Traffic dataset with 862 979 variables) result in faster training and inference speeds for models like PATCHTST and DLINEAR 980 compared to datasets with fewer variables (like the Weather dataset with 21 variables). This leads us 981 to adjust the empirical computation complexity of intra-patch attention. We compute intra-patch at-982 tention for  $n \times C = L/l_p \times C$  times, resulting in a complexity of  $\mathcal{O}(l_p \cdot LC)$ . Similarly, channel-wise 983 attention complexity is  $\mathcal{O}(L \cdot C^2)$ . By adding the outputs of both attentions, the total complexity 984 becomes  $\mathcal{O}((l_p + C) \cdot LC)$ . 985

**P2.** Mask Fusion. As the analysis may suggest, constraining the actual batch size by computing longer sequences of attention may increase efficiency in cases. Indeed, we have developed an alter-988 native approach to implement our attention mechanism, which can reduce computational complexity in certain scenarios. As discussed, the attention components (Q, K, V) have a shape of [B, N, D], 989 where variations in sequence length N and the batch size B affect practical complexity with given GPU parallel capacity. In this sense, increasing the sequence length can enhance efficiency in some cases by having a smaller number of B. 992





Figure 9: Empirical results of fusing masks with fixed 64 variates

1005 Figure 8 demonstrates that by fusing the attention masks, we can compute intra-patch and channel-1006 wise attention simultaneously. This results in a fused attention complexity of  $\mathcal{O}(l_p^2 C^2)$ . When 1007  $l_p^2 C < (l_p + C)L$ , this fused approach is more efficient than the standard implementation. In 1008 practice, we use a small patch length of  $l_p = 4$ . For empirical analysis, we set a batch size of 128, 1009 patch length of 4, and a number of variables as 64 to evaluate the effectiveness of mask fusion. 1010 As shown in Figure 9, time consumption is reduced by 23.34% during training and 27.70% during 1011 inference, though it requires more memory.

**P3.** Variate Sampling. According to (Liu et al., 2024b), even for channel-dependent methods, 1013 sampling variables from different channels does not bring performance degrades while significantly 1014 improving efficiency. As analyzed, this technique should reduce complexity related to the number 1015 of channels and apply to channel-independent methods as well. We also compare memory and time 1016 complexity with this sampling strategy, using a sampling size of 64. The results are presented in 1017 Table 8 and Table 9. 1018

For a more comprehensive empirical study, we conducted complexity experiments on Traffic and 1019 Weather datasets, in line with previous work like SPARSETSF and ITRANSFORMER. We also in-1020 cluded DLINEAR to represent linear models, which should have similar complexities. The results 1021 are illustrated in Figure 10 and Figure 11. 1022

1023 Additionally, inference and training cost increases may differ across different backbones. For instance, training speed might decrease more for DLINEAR and ITRANSFORMER but maintain a rel-1024 atively low decrease for PATCHTST, potentially reducing time complexity due to efficient attention 1025 mask implementation. The efficiency-friendly methods here do not sacrifice performance. An al-

1050 1051 1052

1063

1064

Table 8: Training complexity increases under different settings on Traffic and Weather datasets. For the sampling strategy, we compute the training time of only the sampled variates. The time complexity is measured by "ms/iter" for the time consumed for each iteration. The memory complexity is measured by "GB" for GPU memory consumed. The additional cost is measured by "inc%" for the relative increase in time or memory compared with backbone models.

	Models		DLIN	EAR			ITRANSF	ORMER			PATCH	ITST	
	metric	ti	me	me	emory	tin	ne	me	mory	tii	me	men	nory
	metric	ms/iter	inc%	GB	inc%	ms/iter	inc%	GB	inc%	ms/iter	inc%	GB	inc%
Troffic	w./o. sampling	150.105	710.45%	6.830	563.58%	149.201	82.51%	7.928	62.57%	79.701	31.56%	3.330	33.38%
manie	w. sampling	10.775	51.92%	0.042	3.83%	9.383	36.56%	0.033	2.48%	4.779	15.74%	-0.139	-8.64%
Waathan	w./o. sampling	9.452	86.11%	0.070	8.66%	4.301	30.15%	0.039	3.48%	4.429	24.87%	-0.137	-9.84%
weather	w. sampling	2.449	28.30%	0.035	8.39%	1.941	18.40%	0.023	3.90%	-0.690	-5.28%	-0.008	-1.02%

Table 9: Inference complexity increases under different settings on Traffic and Weather datasets. For
the sampling strategy, we compute the inference time of evaluating all variates. The time complexity
is measured by "ms/iter" for the time consumed for each iteration. The memory complexity is
measured by "GB" for GPU memory consumed. The additional cost is measured by "inc%" for the
relative increase in time or memory compared with backbone models.

1045														
	1	Models		PTN-	DLi			PTN-	ITR			PTN-	Pat	
1046		metric	ti	me	me	emory	tii	me	men	nory	ti	me	mer	nory
10/7		Traffic w./o. sampling		inc%	GB	inc%	ms/iter	inc%	GB	inc%	ms/iter	inc%	GB	inc%
1047	Troffic			237.19%	1.705	301.55%	50.566	81.66%	-0.223	-4.13%	20.872	21.35%	-0.039	3.03%
1048	manne	w. sampling	29.252	133.74%	0.002	0.31%	29.260	88.13%	0.039	5.79%	-4.473	-5.94%	0.004	0.54%
1040	Waathar	w./o. sampling	2.047	23.54%	0.008	1.98%	1.963	20.31%	0.055	11.59%	1.101	11.13%	0.021	4.55%
1049	weather	w. sampling	2.449	28.30%	0.035	8.39%	1.941	18.40%	0.023	3.90%	-0.690	-5.28%	-0.008	-1.02%

0.45

0.35

0.25

MAE

0.57GB, 17.51ms









Traffic (862 variates)

2.96GB, 97.75ms

0.5GB

PTN-Pat

2.92GB, 118.62ms

PTN-iTr

5.17GB. 112.48ms

4.5GB

PTN-DI

.27GB, 59.03ms

(b) Inference complexity for Traffic dataset with all variates



(c) Training complexity for Traffic dataset with
 variate sampling

1078 1079

(d) Inference complexity for Traffic dataset with variate sampling

Figure 10: Visualization of computation and memory complexity on Traffic dataset









Weather (variate sampling)

PTN-DLi

0.45GB, 11.10

10

0.60GB, 10.55ms

0.5GB

0.7566

PTN-Pat

DTN-iTr

0.76GB, 12.39

0.62GB, 12.49

inference time (ms/iter)

1GB

0.76GB.13.08ms

14



(c) Training complexity for Weather dataset with variate sampling





ternative option, detailed in Appendix A.4.4, incurs no additional inference costs with performance
 compromises, allowing for a flexible trade-off between efficiency and effectiveness.

## 1112 A.4.4 EFFICIENT STUDENT MODEL BY TRAIN-TIME DISTILLATION

1113 The training cost for non-linear attention-based models can be reduced by a data transfer protocol 1114 as described in Section 5.3. The transfer utilizes the transformation learned on a linear backbone to 1115 train a non-linear model from scratch, which improves training efficiency at a minor compromise 1116 of performance. We have a more direct and burden-free version for inference as well, ss shown in Figure 7 and briefly explained in the previous section, The core idea is to train a student model that 1117 directly uses the raw input X to predict the transformed output  $\dot{Y}$ . This practice is an intermediate 1118 version between a standard PTN and the data transfer version. Unlike the data transfer protocol in 1119 Section 5.3, the student model learns the mapping from X to Y during training, which means it 1120 requires more samples than just the converged  $\bar{Y}$  after training PTN. When training a standard PTN, 1121 the backbone predictor uses the transformed series as inputs and labels, and updates computed by 1122 the  $l_p red$  will also back-propagate the gradients of the PTN parameters. Whereas when training 1123 with a student predictor, we use straight-through method (Van Den Oord et al., 2017) to omit the 1124 calculation of PTN gradients (i.e. use "tensor.detach()" in pytorch). Thereafter, we only added the 1125 cost of another backbone model and can achieve similar improvements without any additional cost 1126 increase during the inference phase.

1127 1128

1090

1091

1092 1093

1094

1095

1098

1099

1100

1101

1102 1103

1104

1105 1106

1107 1108

1111

## A.4.5 CONVOLUTION ENCODER

**Transpose and Unfold** The implementation details regarding Convolution Encoder involve two operations that are designed for the following two benefits. (1) The convolution outputs can be concatenated in the same length of embedding so that features from different frequencies can be ideally fused into one uniform embedding. (2) The features can be evenly arranged along the temporal dimension so that each embedding in the sequence can have the same large Receptive Field. There-

137	Models	DLIN	NEAR	+PTN	i(stu)	PATC	HTST	+PTN	v(stu)	ITRANS	FORMER	+PTN	(stu)	
107	M	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
138		96	0.386	0.395	0.384	0.391	0.414	0.419	0.396	0.411	0.386	0.405	0.401	0.409
139	hl	192	0.437	0.424	0.436	0.422	0.460	0.445	0.455	0.447	0.441	0.436	0.456	0.449
1/10	L	336	0.479	0.446	0.480	0.444	0.501	0.466	0.496	0.473	0.487	0.458	0.540	0.504
140	Ы	720	0.481	0.470	0.475	0.464	0.500	0.488	0.492	0.490	0.503	0.491	0.594	0.554
141		Avg.	0.446	0.434	0.444	0.430	0.469	0.455	0.460	0.456	0.454	0.448	0.498	0.479
142		96	0.288	0.338	0.284	0.335	0.302	0.348	0.296	0.347	0.297	0.349	0.302	0.349
1/10	h2	192	0.374	0.390	0.370	0.387	0.388	0.400	0.377	0.396	0.380	0.400	0.392	0.410
143	E	336	0.415	0.426	0.414	0.424	0.426	0.433	0.423	0.437	0.428	0.432	0.451	0.449
144	ЕÚ	720	0.420	0.440	0.416	0.436	0.431	0.446	0.447	0.459	0.427	0.445	0.603	0.550
145		Avg.	0.374	0.399	0.371	0.396	0.387	0.407	0.386	0.410	0.383	0.407	0.437	0.440
	y	96	0.201	0.281	0.206	0.278	0.181	0.270	0.174	0.258	0.148	0.240	0.143	0.234
146	icit	192	0.201	0.283	0.205	0.281	0.188	0.274	0.184	0.268	0.162	0.253	0.154	0.247
147	ctri	336	0.215	0.298	0.219	0.296	0.204	0.293	0.200	0.284	0.178	0.269	0.169	0.263
1/10	Ie	720	0.257	0.331	0.260	0.328	0.246	0.324	0.239	0.316	0.225	0.317	0.209	0.300
140	щ	Avg.	0.219	0.298	0.222	0.296	0.205	0.290	0.199	0.281	0.178	0.270	0.169	0.261
149		96	0.649	0.389	0.664	0.384	0.462	0.290	0.474	0.277	0.395	0.268	0.406	0.255
150	fic	192	0.601	0.366	0.615	0.354	0.466	0.290	0.478	0.279	0.417	0.276	0.436	0.276
124	rafi	336	0.609	0.369	0.621	0.357	0.482	0.300	0.487	0.287	0.433	0.283	0.448	0.275
101	F	720	0.647	0.387	0.657	0.380	0.514	0.320	0.525	0.303	0.467	0.302	0.476	0.291
152		Avg.	0.627	0.378	0.639	0.369	0.481	0.300	0.491	0.286	0.428	0.282	0.442	0.274

Table 10: Full results for performance comparisons between backbone models and student models
 with train-time distillation by PTN.

fore, we introduce the two-step operation of Transpose and Unfold, which brings the two benefits together. Specifically, we set kernel size = 3, stride = 2, paddig = 1, and the number of kernels doubled for each next layer. In this way, as shown in Figure 12, we can ensure the number of fea-tures is invariant for different layers, and only the shape is changed. Now we can fuse the outputs from different convolution layers together by flattening/unfolding the features to the original shape of  $L \times 1$ . Again, considering the effectiveness of point-wise linear head 4.1, we want the concate-nated features to be near-equally arranged along the temporal dimension to maintain the sequential relationship in the embedding. So we transpose the features first and unfold them afterwards. This practice can ensure a Receptive Field of  $(2^{l+1}-1)$  wide for each embedding, where l is the total number of convolution layers.

Effective Receptive Field of Convolution Encoder As what was proposed by (Luo & Wang, 2024), the Effective Receptive Field (ERF) is a reasonable consideration for designing convolution-based architecture. We also input the impulse function into our Convolution Encoder to visualize the ERF, as shown in Figure 13. It has shown that without using an extra large convolution kernel, our proposed method can utilize a near-global ERF by combining outputs from different layers, with different frequency patterns.



Figure 12: The illustration of transpose and unfold operation in the Convolution Encoder





Like any Mix-of-Experts model, increasing the size of MoE can help generate better results. We conduct such a hyperparameter study of two versions of our PTN model on the Weather dataset. The results are reported in Fig 15.

1246 A.5.2 REVISITING TIME SERIES MIXTURE-OF-EXPERTS: A DECOMPOSITION VIEW

Time series MoE has been recently proposed as an effective architecture in scaling parameters for time series models (Shi et al., 2024; Liu et al., 2023; Ni et al., 2024). Here, we conduct simple attempts to extend our MoE model to a larger number of 12, without subseries routing (i.e., each Expert transforms the whole series). This practice can be computation-intensive and beyond our capacity for extensive research. However, when we train the larger-in-scale MoE on par with its simplified counterpart (with subseries routing) on ETTh1, the visualization for each Expert shows patterns as a decomposition of time series, as shown in Figure 16.



Figure 16: Transformed series from different Expert Models under full-transform setting.

## A.6 MORE RESULTS OF EXPERIMENTS

Table 11 and Table 12 show the full results under the look-back length L = 96 and L = 512 settings respectively. Tab 13 shows full results of a comparison study between popular backbones and them integrated with PTN-RLI.

1293

1290 1291

1278 1279

1280

1281

1282

1283

- 1294
- 1295

Table 11: Multivariate long-term forecasting results of PTN-RLI (i.e. PTN with RLINEAR backbone) and PTN-ITR (i.e. PTN with ITRANSFORMER backbone). We use look-back length fixed as 96 and prediction lengths  $T \in \{96, 192, 336, 720\}$  for all datasets. The best results are in **bold** and the second best is underlined. 

1313																		
1313	M	odels	PTN	-RLi	PTN	-ITR	ITRAN	SFORMER	PATCI	HTST	TIME	SNET	SPARS	SETSF	DLI	NEAR	FEDF	ORMER
1314	Μ	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
1315		96	0.372	0.385	0.389	0.407	0.386	0.405	0.414	0.419	0.384	0.402	0.391	0.389	0.386	0.395	0.376	0.419
1010	hl	192	0.431	0.418	0.432	0.432	0.441	0.436	0.460	0.445	0.436	0.429	0.441	0.419	0.437	0.424	0.420	0.448
1316	E	336	0.463	0.438	0.460	0.443	0.487	0.458	0.501	0.466	0.491	0.469	0.481	0.439	0.479	0.446	0.459	0.465
1317	щ	720	0.464	0.460	0.532	0.523	0.503	0.491	0.500	0.488	0.521	0.500	0.467	0.453	0.481	0.470	0.506	0.507
1017		Avg.	0.437	0.427	0.453	0.451	0.454	0.448	0.469	0.455	0.458	0.450	0.445	0.425	0.446	0.434	0.440	0.460
1318		96	0.281	0.330	0.323	0.367	0.297	0.349	0.302	0.348	0.340	0.374	0.310	0.346	0.288	0.338	0.358	0.397
1319	h2	192	0.361	0.380	0.377	0.399	0.380	0.400	0.388	0.400	0.402	0.414	0.391	0.394	0.374	0.390	0.429	0.439
1010	E	336	0.409	0.418	0.453	0.450	0.428	0.432	0.426	0.433	0.452	0.452	0.424	0.428	0.415	0.426	0.496	0.487
1320	щ	720	0.413	0.434	0.610	0.552	0.427	0.445	0.431	0.446	0.462	0.468	0.421	0.437	0.420	0.440	0.463	0.474
1321		Avg.	0.366	0.391	0.441	0.442	0.383	0.407	0.387	0.407	0.414	0.427	0.386	0.402	0.374	0.399	0.437	0.449
1021		96	0.333	0.351	0.330	0.357	0.334	0.368	0.329	0.367	0.338	0.375	0.345	0.355	0.355	0.376	0.379	0.419
1322	n1	192	0.376	0.374	0.372	0.381	0.377	0.391	0.367	0.385	0.374	0.387	0.390	0.377	0.391	0.392	0.426	0.441
1323	Ê	336	0.404	0.395	0.408	0.405	0.426	0.420	0.399	0.410	0.410	0.411	0.425	0.399	0.424	0.415	0.445	0.459
	Ш	720	0.463	0.429	0.467	0.440	0.491	0.459	0.454	0.439	0.478	0.450	0.496	0.438	0.487	0.450	0.543	0.490
1324		Avg.	0.399	0.390	0.394	0.396	0.407	0.410	0.387	0.400	0.400	0.406	0.414	0.392	0.414	0.408	0.448	0.452
1325		96	0.178	0.258	0.184	0.268	0.180	0.264	0.175	0.259	0.187	0.267	0.184	0.264	0.182	0.265	0.203	0.287
	m2	192	0.234	0.298	0.257	0.323	0.250	0.309	0.241	0.302	0.249	0.309	0.248	0.304	0.246	0.304	0.269	0.328
1326	Ê	336	0.294	0.333	0.324	0.367	0.311	0.348	0.305	0.343	0.321	0.351	0.308	0.338	0.307	0.342	0.325	0.366
1327	Ш	720	0.383	0.386	0.405	0.412	0.412	0.407	0.402	0.400	0.408	0.403	0.407	0.396	0.407	0.398	0.421	0.415
		Avg.	0.272	0.319	0.293	0.343	0.288	0.332	0.281	0.326	0.291	0.333	0.287	0.326	0.286	0.327	0.305	0.349
1328	~	96	0.200	0.271	0.143	0.227	0.148	0.240	0.181	0.270	0.168	0.272	0.206	0.274	0.201	0.281	0.193	0.308
1329	cit	192	0.198	0.274	0.155	0.245	0.162	0.253	0.188	0.274	0.184	0.289	0.203	0.276	0.201	0.283	0.201	0.315
1000	ctri	336	0.212	0.289	0.167	0.259	0.178	0.269	0.204	0.293	0.198	0.300	0.216	0.291	0.215	0.298	0.214	0.329
1330	Еle	720	0.252	0.322	0.189	0.282	0.225	0.317	0.246	0.324	0.220	0.320	0.257	0.325	0.257	0.331	0.246	0.355
1331		Avg.	0.215	0.289	0.164	0.253	0.178	0.270	0.205	0.290	0.193	0.295	0.221	0.291	0.219	0.298	0.214	0.327
1000		96	0.662	0.366	0.429	0.251	0.395	0.268	0.462	0.290	0.593	0.321	0.705	0.372	0.649	0.389	0.587	0.366
1332	fic	192	0.621	0.346	0.449	0.258	0.417	0.276	0.466	0.290	0.617	0.336	0.656	0.349	0.601	0.366	0.604	0.373
1333	raf	336	0.634	0.345	0.469	0.269	0.433	0.283	0.482	0.300	0.629	0.336	0.657	0.351	0.609	0.369	0.621	0.383
1001	Е	720	0.654	0.364	0.503	0.287	0.467	0.302	0.514	0.320	0.640	0.350	0.688	0.372	0.647	0.387	0.626	0.382
1334		Avg.	0.643	0.355	0.463	0.266	0.428	0.282	0.481	0.300	0.620	0.336	0.676	0.361	0.627	0.378	0.610	0.376
1335		96	0.199	0.229	0.160	0.207	0.174	0.214	0.177	0.210	0.172	0.220	0.224	0.254	0.192	0.232	0.217	0.296
1000	her	192	0.234	0.261	0.205	0.253	0.221	0.254	0.225	0.250	0.219	0.261	0.263	0.285	0.240	0.271	0.276	0.336
1336	eat	336	0.285	0.297	0.253	0.285	0.278	0.296	0.278	0.290	0.280	0.306	0.314	0.319	0.292	0.307	0.339	0.380
1337	3	720	0.352	0.344	0.323	0.338	0.358	0.349	0.354	0.340	0.365	0.359	0.386	0.363	0.364	0.353	0.403	0.428
1000		Avg.	0.268	0.283	0.235	0.271	0.258	0.278	0.259	0.273	0.259	0.287	0.297	0.305	0.272	0.291	0.309	0.360
1 5 5 25																		

Table 12: Multivariate long-term forecasting results of PTN-RLI (i.e. PTN with RLINEAR backbone) and PTN-ITR (i.e. PTN with ITRANSFORMER backbone). We use look-back length fixed as **512** and prediction lengths  $T \in \{96, 192, 336, 720\}$  for all datasets. The best results are in **bold** and the second best is underlined.

1367	Models		Models PTN-RLI		PTN-ITR I		ITRAN	ITRANSFORMER		PATCHTST		TIMESNET		SPARSETSE		DLINEAR		FEDFORMER	
1368	M	letric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
1000		96	0.357	0.383	0.413	0.435	0.421	0.442	0.375	0.397	0 384	0.402	0.362	0.388	0.382	0.405	0.376	0.419	
1369	-	192	0.397	0.407	0.439	0.450	0.448	0.454	0.412	0.421	0.436	0.429	0.403	0.411	0.417	0.425	0.420	0.448	
1370	Ę	336	0.427	0.426	0.468	0.474	0.452	0.462	0.436	0.437	0.491	0.469	0.434	0.428	0.436	0.442	0.459	0.465	
1071	Щ	720	0.431	0.450	0.542	0.546	0.564	0.558	0.465	0.471	0.521	0.500	0.426	0.447	0.433	0.455	0.506	0.507	
1371		Avg.	0.403	0.417	0.466	0.476	0.471	0.479	0.422	0.431	0.458	0.450	0.406	0.419	0.417	0.432	0.440	0.460	
1372		96	0.268	0.326	0.322	0.375	0.353	0.387	0.287	0.342	0.340	0.374	0.294	0.346	0.272	0.336	0.358	0.397	
1070	2	192	0.334	0.371	0.379	0.410	0.444	0.449	0.349	0.383	0.402	0.414	0.339	0.377	0.333	0.375	0.429	0.439	
1373	Ē	336	0.373	0.403	0.442	0.452	0.434	0.451	0.364	0.404	0.452	0.452	0.359	0.397	0.355	0.396	0.496	0.487	
1374	Ш	720	0.417	0.441	0.489	0.487	0.545	0.520	0.411	0.437	0.462	0.468	0.383	0.424	0.378	0.423	0.463	0.474	
1275		Avg.	0.348	0.385	0.408	0.431	0.444	0.452	0.353	0.392	0.414	0.427	0.344	0.386	0.335	0.383	0.437	0.449	
1375		96	0.293	0.336	0.322	0.354	0.324	0.363	0.297	0.337	0.338	0.375	0.312	0.354	0.311	0.354	0.364	0.416	
1376	nl	192	0.330	0.358	0.371	0.394	0.377	0.396	0.340	0.366	0.374	0.387	0.347	0.376	0.340	0.369	0.426	0.441	
1377	Ē	336	0.361	0.378	0.399	0.412	0.400	0.409	0.375	0.386	0.410	0.411	0.367	0.386	0.367	0.385	0.445	0.459	
1011	Ш	720	0.420	0.411	0.452	0.441	0.455	0.443	0.424	0.419	0.478	0.450	0.419	0.413	0.416	0.412	0.543	0.490	
1378		Avg.	0.351	0.371	0.386	0.400	0.389	0.403	0.380	0.390	0.400	0.406	0.361	0.382	0.359	0.380	0.448	0.452	
1379		96	0.162	0.248	0.189	0.275	0.186	0.272	0.176	0.259	0.187	0.267	0.163	0.252	0.163	0.254	0.203	0.287	
	m2	192	0.216	0.285	0.336	0.375	0.250	0.315	0.237	0.304	0.249	0.309	0.217	0.290	0.217	0.291	0.269	0.328	
1380	Ê	336	0.269	0.321	0.320	0.360	0.303	0.347	0.301	0.346	0.321	0.351	0.270	0.327	0.268	0.326	0.325	0.366	
1381	Щ	720	0.357	0.377	0.403	0.412	0.395	0.407	0.409	0.408	0.408	0.403	0.352	0.379	0.349	0.378	0.421	0.415	
1000		Avg.	0.251	0.308	0.312	0.356	0.283	0.335	0.281	0.329	0.291	0.333	0.251	0.312	0.249	0.312	0.305	0.349	
1382	2	96	0.135	0.230	0.129	0.221	0.130	0.222	0.133	0.224	0.168	0.272	0.138	0.233	0.145	0.248	0.193	0.308	
1383	icit	192	<u>0.149</u>	0.243	0.149	0.240	0.149	0.242	0.148	0.237	0.184	0.289	0.151	0.244	0.159	0.260	0.201	0.315	
100/	ectr	336	0.165	0.259	0.167	0.259	0.170	0.263	0.165	0.256	0.198	0.300	0.166	0.260	0.175	0.275	0.214	0.329	
1304	Ē	720	0.203	0.290	0.188	0.278	0.196	0.288	0.203	0.287	0.220	0.320	0.205	0.293	0.212	0.305	0.246	0.355	
1385		Avg.	0.163	0.255	0.158	0.250	0.161	0.254	0.162	0.251	0.193	0.295	0.165	0.258	0.173	0.272	0.214	0.327	
1226		96	0.414	0.269	0.370	0.234	0.349	0.243	0.397	0.259	0.593	0.321	0.389	0.268	0.398	0.286	0.587	0.366	
1300	ffic	192	0.422	0.271	0.394	0.245	0.363	0.246	0.410	0.266	0.617	0.336	0.398	0.270	0.409	0.289	0.604	0.373	
1387	Irai	336	0.433	0.276	0.415	0.256	0.381	0.258	0.417	0.269	0.629	0.336	0.411	0.275	0.421	0.294	0.621	0.383	
1388		720	0.466	0.293	0.451	0.274	0.425	0.282	0.451	0.287	0.640	0.350	0.448	0.297	0.457	0.311	0.626	0.382	
1300		Avg.	0.434	0.277	0.408	0.252	0.380	0.257	0.419	0.270	0.620	0.336	0.412	0.278	0.421	0.295	0.610	0.376	
1389	-	96	0.162	0.209	0.163	0.208	$\frac{0.158}{0.200}$	0.213	0.154	0.190	0.172	0.220	0.169	0.223	0.170	0.225	0.217	0.296	
1300	the	192	0.205	0.246	0.216	0.254	0.209	0.262	0.204	0.239	0.219	0.261	0.214	0.262	0.212	0.260	0.276	0.336	
1000	Vea	336	0.250	0.282	0.271	0.296	0.253	0.295	0.268	0.286	0.280	0.306	0.257	0.293	0.258	0.294	0.339	0.380	
1391	2	720	0.318	0.333	0.382	0.366	0.319	0.347	0.354	0.344	0.365	0.359	0.321	0.340	0.320	0.339	0.403	0.428	
1202		Avg.	0.234	0.268	0.258	0.281	0.235	0.279	0.245	0.265	0.259	0.287	0.240	0.280	0.240	0.280	0.309	0.360	

1407	Models		Models LINEAR		+PTN		DLINEAR		+PTN		RLINEAR		+P	TN	ITRANS	FORMER	+PTN		PATCI	HTST	Γ +PTN	
1400	Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
1400		96	0.379	0.394	0.381	0.391	0.386	0.395	0.373	0.387	0.383	0.386	0.372	0.385	0.386	0.405	0.389	0.407	0.414	0.419	0.383	0.401
1/100	- 2	192	0.430	0.425	0.444	0.434	0.437	0.424	0.431	0.419	0.437	0.417	0.431	0.418	0.441	0.436	0.432	0.432	0.460	0.445	0.429	0.424
1403	E	336	0.474	0.452	0.477	0.449	0.479	0.446	0.469	0.440	0.482	0.440	0.463	0.438	0.487	0.458	0.460	0.443	0.501	0.466	0.473	0.450
1410	ш	720	0.505	0.504	0.501	0.497	0.481	0.470	0.468	0.462	0.476	0.460	0.464	0.460	0.503	0.491	0.532	0.523	0.500	0.488	0.495	0.494
1110		Avg.	0.447	0.444	0.451	0.443	0.446	0.434	0.435	0.427	0.444	0.426	0.437	0.427	0.454	0.448	0.453	0.451	0.469	0.455	0.445	0.442
1411		96	0.294	0.346	0.283	0.335	0.288	0.338	0.287	0.336	0.291	0.335	0.281	0.330	0.297	0.349	0.323	0.367	0.302	0.348	0.294	0.343
	圮	192	0.380	0.401	0.370	0.394	0.374	0.390	0.367	0.385	0.376	0.387	0.361	0.380	0.380	0.400	0.377	0.399	0.388	0.400	0.379	0.400
1412	E	336	0.449	0.453	0.450	0.451	0.415	0.426	0.417	0.423	0.420	0.423	0.409	0.418	0.428	0.432	0.453	0.450	0.426	0.433	0.450	0.443
1/10	ш	720	0.588	0.541	0.580	0.534	0.420	0.440	0.436	0.454	0.418	0.435	0.413	0.434	0.427	0.445	0.610	0.552	0.431	0.446	0.474	0.472
1413		Avg.	0.428	0.435	0.421	0.429	0.374	0.399	0.377	0.400	0.376	0.395	0.366	0.391	0.383	0.407	0.441	0.442	0.387	0.407	0.399	0.415
1414		96	0.338	0.360	0.328	0.349	0.355	0.376	0.335	0.354	0.337	0.352	0.333	0.351	0.334	0.368	0.330	0.357	0.329	0.367	0.315	0.344
1414	F	192	0.379	0.382	0.374	0.374	0.391	0.392	0.379	0.376	0.383	0.375	0.376	0.374	0.377	0.391	0.372	0.381	0.367	0.385	0.363	0.373
1415	El	336	0.410	0.406	0.405	0.396	0.424	0.415	0.413	0.398	0.417	0.397	0.404	0.395	0.426	0.420	0.408	0.405	0.399	0.410	0.389	0.392
	Ш	720	0.472	0.445	0.470	0.437	0.487	0.450	0.475	0.434	0.485	0.435	0.463	0.429	0.491	0.459	0.467	0.440	0.454	0.439	0.457	0.434
1416		Avg.	0.400	0.398	0.394	0.389	0.414	0.408	0.401	0.390	0.405	0.390	0.399	0.390	0.407	0.410	0.394	0.396	0.387	0.400	0.381	0.386
		96	0.183	0.263	0.179	0.260	0.182	0.265	0.179	0.258	0.183	0.257	0.178	0.258	0.180	0.264	0.184	0.268	0.175	0.259	0.171	0.258
1417	n2	192	0.245	0.307	0.238	0.300	0.246	0.304	0.239	0.302	0.247	0.299	0.234	0.298	0.250	0.309	0.257	0.323	0.241	0.302	0.231	0.300
1/110	Ē	336	0.306	0.349	0.297	0.342	0.307	0.342	0.296	0.336	0.307	0.337	0.294	0.333	0.311	0.348	0.324	0.367	0.305	0.343	0.344	0.343
1410	ш	720	0.403	0.407	0.392	0.397	0.407	0.398	0.391	0.396	0.408	0.394	0.383	0.386	0.412	0.407	0.405	0.412	0.402	0.400	0.424	0.421
1/110		Avg.	0.284	0.332	0.277	0.325	0.286	0.327	0.276	0.323	0.286	0.322	0.272	0.319	0.288	0.332	0.293	0.343	0.281	0.326	0.293	0.331
1413	~	96	0.200	0.277	0.197	0.271	0.201	0.281	0.199	0.271	0.201	0.270	0.200	0.271	0.148	0.240	0.143	0.227	0.181	0.270	0.169	0.248
1420	icit.	192	0.199	0.281	0.196	0.274	0.201	0.283	0.198	0.274	0.200	0.273	0.198	0.274	0.162	0.253	0.155	0.245	0.188	0.274	0.177	0.257
	ct	336	0.212	0.297	0.209	0.289	0.215	0.298	0.211	0.289	0.214	0.289	0.212	0.289	0.178	0.269	0.167	0.259	0.204	0.293	0.194	0.275
1421	Ξ	720	0.248	0.328	0.244	0.320	0.257	0.331	0.251	0.321	0.255	0.322	0.252	0.322	0.225	0.317	0.189	0.282	0.246	0.324	0.242	0.318
4.400		Avg.	0.215	0.296	0.212	0.289	0.219	0.298	0.215	0.289	0.217	0.289	0.215	0.289	0.178	0.270	0.164	0.253	0.205	0.290	0.196	0.274
1422		96	0.653	0.376	0.702	0.366	0.649	0.389	0.670	0.365	0.662	0.366	0.694	0.365	0.395	0.268	0.429	0.251	0.462	0.290	0.548	0.281
1/102	Э	192	0.603	0.351	0.657	0.343	0.601	0.366	0.625	0.343	0.621	0.346	0.651	0.342	0.417	0.276	0.449	0.258	0.466	0.290	0.554	0.277
1423	Taf	336	0.607	0.353	0.666	0.348	0.609	0.369	0.634	0.345	0.634	0.345	0.655	0.344	0.433	0.283	0.469	0.269	0.482	0.300	0.576	0.288
1494		720	0.642	0.373	0.674	0.352	0.647	0.387	0.667	0.364	0.654	0.364	0.686	0.364	0.467	0.302	0.503	0.287	0.514	0.320	0.659	0.330
1-14-1		Avg.	0.626	0.363	0.675	0.352	0.627	0.378	0.649	0.354	0.643	0.355	0.672	0.354	0.428	0.282	0.463	0.266	0.481	0.300	0.584	0.294
1425		96	0.200	0.237	0.192	0.228	0.192	0.232	0.194	0.227	0.209	0.230	0.199	0.229	0.174	0.214	0.160	0.207	0.177	0.210	0.168	0.208
-	her	192	0.238	0.275	0.238	0.274	0.240	0.271	0.234	0.261	0.253	0.266	0.234	0.261	0.221	0.254	0.205	0.253	0.225	0.250	0.214	0.252
1426	eat	336	0.281	0.310	0.283	0.315	0.292	0.307	0.278	0.296	0.302	0.302	0.285	0.297	0.278	0.296	0.253	0.285	0.278	0.290	0.261	0.289
4.407	*	720	0.341	0.357	0.339	0.353	0.364	0.353	0.352	0.344	0.374	0.350	0.352	0.344	0.358	0.349	0.323	0.338	0.354	0.340	0.343	0.341
1427		Avg.	0.265	0.295	0.263	0.292	0.272	0.291	0.265	0.282	0.284	0.287	0.268	0.283	0.258	0.278	0.235	0.271	0.259	0.273	0.246	0.272

Table 13: Full results for performance comparisons between backbone models and them integrated with our PTN on ETT, electricity, traffic and weather datasets. 

Table 14: Full results for performance comparisons between backbone models and them integrated with our PTN on PeMS datasets.

M	Models DLINEAR Metric MSE MA		NEAR	+P	TN	PATC	htst	+P	TN	ITRANS	SFORMER	+PTN		
M	letric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
~	12	0.116	0.224	0.097	0.207	0.113	0.229	0.076	0.179	0.071	0.174	0.069	0.171	
20	24	0.233	0.316	0.143	0.254	0.211	0.313	0.110	0.218	0.093	0.201	0.091	0.195	
Ň	36	0.381	0.410	0.202	0.304	0.362	0.420	0.143	0.251	0.125	0.236	0.115	0.222	
ΡE	48	0.536	0.505	0.246	0.345	0.504	0.512	0.164	0.271	0.164	0.275	0.135	0.240	
_	Avg.	0.317	0.364	0.172	0.277	0.298	0.368	0.123	0.230	0.113	0.222	0.102	0.207	
	12	0.128	0.236	0.107	0.215	0.128	0.244	0.094	0.198	0.078	0.183	0.082	0.183	
04	24	0.243	0.329	0.161	0.269	0.252	0.355	0.143	0.248	0.095	0.205	0.107	0.213	
Ň	36	0.390	0.424	0.238	0.330	0.413	0.472	0.172	0.276	0.120	0.233	0.134	0.240	
E	48	0.555	0.518	0.266	0.352	0.573	0.573	0.202	0.305	0.150	0.262	0.161	0.266	
-	Avg.	0.329	0.377	0.193	0.292	0.341	0.411	0.153	0.257	0.111	0.221	0.121	0.225	
~	12	0.109	0.219	0.096	0.203	0.099	0.212	0.075	0.167	0.067	0.165	0.066	0.158	
202	24	0.230	0.318	0.160	0.265	0.203	0.311	0.115	0.214	0.088	0.190	0.089	0.184	
Ň	36	0.385	0.419	0.209	0.308	0.333	0.406	0.135	0.236	0.110	0.215	0.113	0.209	
E	48	0.551	0.515	0.262	0.355	0.535	0.536	0.156	0.255	0.139	0.245	0.134	0.226	
_	Avg.	0.319	0.368	0.182	0.283	0.293	0.366	0.120	0.218	0.101	0.204	0.101	0.194	
	12	0.121	0.229	0.152	0.223	0.120	0.239	0.140	0.198	0.079	0.182	0.128	0.184	
808	24	0.235	0.321	0.242	0.292	0.226	0.331	0.187	0.240	0.115	0.219	0.161	0.212	
MS	36	0.380	0.419	0.318	0.351	0.345	0.414	0.230	0.276	0.186	0.235	0.188	0.233	
E	48	0.549	0.518	0.325	0.368	0.506	0.502	0.245	0.294	0.221	0.267	0.215	0.252	
4	Avg.	0.321	0.372	0.259	0.308	0.299	0.372	0.200	0.252	0.150	0.226	0.173	0.220	

## A.7 MORE VISUALIZATIONS

#### A.7.1 WEIGHT SPARSITY COMPARISONS

As proposed in SparseTSF (Lin et al., 2024), TSF datasets can require far fewer parameters involved and information provided to make near-SOTA forecasting. The explanation for this might be the sparsity. As some work (Ilbert et al., 2024; Tan et al., 2024) also suggests, complex modeling of Transformer models might suffer from overfitting non-sparse features in time series data that even have degraded Indentity matrix as attention scores. The visualization of weight matrices learned can

M	Models		нтѕт	RLINEA	R→PATCHTST	ITRANS	FORMER	RLINEAR→ITRANSFORMER		
M	etric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	
	96	0.414	0.419	0.402	0.410	0.386	0.405	0.395	0.409	
hl	192	0.460	0.445	0.425	0.427	0.441	0.436	0.432	0.430	
E	336	0.501	0.466	0.473	0.456	0.487	0.458	0.473	0.453	
田	720	0.500	0.488	0.515	0.505	0.503	0.491	0.524	0.513	
	Avg.	0.469	0.455	0.454	0.450	0.454	0.448	0.456	0.451	
	96	0.302	0.348	0.306	0.349	0.297	0.349	0.372	0.400	
h2	192	0.388	0.400	0.389	0.402	0.380	0.400	0.390	0.403	
E	336	0.426	0.433	0.460	0.443	0.428	0.432	0.452	0.450	
函	720	0.431	0.446	0.527	0.526	0.427	0.445	0.777	0.619	
	Avg.	0.387	0.407	0.420	0.430	0.383	0.407	0.498	0.468	
ty	96	0.181	0.270	0.169	0.247	0.148	0.240	0.147	0.233	
::	192	0.188	0.274	0.177	0.256	0.162	0.253	0.162	0.248	
Ľ,	336	0.204	0.293	0.192	0.272	0.178	0.269	0.172	0.262	
lec	720	0.246	0.324	0.230	0.304	0.225	0.317	0.185	0.277	
Щ	Avg.	0.205	0.290	0.192	0.270	0.178	0.270	0.166	0.255	
	96	0.462	0.290	0.549	0.274	0.395	0.268	0.486	0.242	
fic	192	0.466	0.290	0.552	0.277	0.417	0.276	0.510	0.253	
af	336	0.482	0.300	0.568	0.283	0.433	0.283	0.527	0.264	
T.	720	0.514	0.320	0.602	0.302	0.467	0.302	0.558	0.281	
	Avg.	0.481	0.300	0.568	0.284	0.428	0.282	0.520	0.260	

Table 15: Full results for Table 4 on comparison of generalization capabilities of our data transfer protocol. 'A $\rightarrow$ B' indicates the raw data transformation is learned by Model A, and then fed to the Model B 

provide some insights, such as comparisons in sparsity. Here we also show the visualized weight matrices between w. and w./o. PTN in Figures 17-20 as an intuitive comparison.



Figure 17: Comparision between weights learned on ETTh1.

#### A.7.2 TRANING OF THE PROXIMAL TRANSFORMATION

As Section 5.2 showcases, the training processes of how the transformed data approaches the raw data are Channel-dependent. Here, we show all channels in the following Figures 21-27.





Figure 25: Channel 5 for PTN-RLI on the ETTh2 dataset.

