# PSFORMER: PARAMETER-EFFICIENT TRANSFORMER WITH SEGMENT ATTENTION FOR TIME SERIES FORE CASTING

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027 028

029

Paper under double-blind review

## ABSTRACT

Time series forecasting remains a critical challenge across various domains, often complicated by high-dimensional data and long-term dependencies. This paper presents a novel transformer architecture for time series forecasting, incorporating two key innovations: parameter sharing (PS) and Spatial-Temporal Segment Attention (SegAtt). We also define the time series segment as the concatenation of sequence patches from the same positions across different variables. The proposed model, PSformer, reduces the number of training parameters through the parameter sharing mechanism, thereby improving model efficiency and scalability. The introduction of SegAtt could enhance the capability of capturing local spatio-temporal dependencies by computing attention over the segments, and improve global representation by integrating information across segments. The combination of parameter sharing and SegAtt significantly improves the forecasting performance. Extensive experiments on benchmark datasets demonstrate that PSformer outperforms popular baselines and other transformer-based approaches in terms of accuracy and scalability, establishing itself as an accurate and scalable tool for time series forecasting.

## 1 INTRODUCTION

Time series forecasting is an important learning task with significant application values in a wide 031 range of domains, including the weather prediction (Ren et al., 2021; Chen et al., 2023a), traffic flow (Tedjopurnomo et al., 2020; Khan et al., 2023), energy consumption (Liu et al., 2020; Nti et al., 033 2020), anomaly detection (Zamanzadeh Darban et al., 2022) and the financial analysis (Nazareth & 034 Ramana Reddy, 2023), etc. With the advancement of artificial intelligence techniques, significant efforts have been devoted to developing innovative models that continue to improve the prediction performance (Liang et al., 2024; Wang et al., 2024b). In particular, the transformer-based model family has recently attracted more attention for its proved success in nature language processing 037 (OpenAI et al., 2024) and computer vision (Liu et al., 2021; Dosovitskiy et al., 2021). Moreover, pre-trained large models based on the transformer architecture have shown advantages in time series forecasting(Liu et al., 2024a; Jin et al., 2024; Chang et al., 2023; Woo et al., 2024), demonstrating 040 that increasing the amount of parameters in transformer models and the volume of training data can 041 effectively enhance the model capability. 042

On the other side, many simple linear models (Zeng et al., 2023; Li et al., 2023) also shown competi-043 tive performance compared to the more complex designs of transformer-based models. One possible 044 key reason for their success in time series forecasting is that they have low model complexities and are less likely to overfit on noisy or irrelevant signals. As a result, even with limited data, these 046 models can effectively capture robust information representations. To overcome the limitations of 047 modeling long-term dependencies and capturing complex temporal relationships, PatchTST (Nie 048 et al., 2023) process temporal information by combining patching techniques to extract local semantic information, leading to superior performance. However, it applies channel-independent designs and leaves the significant potential for improvement in effectively modeling across variables. More-051 over, the unique challenges of modeling multivariate time series data, where the temporal and spatial dimensions differ significantly from other data types, present many unexplored opportunities. While 052 past research (donghao & wang xue, 2024; Zhang & Yan, 2023; Ilbert et al., 2024) has largely treated these dimensions separately, the question of whether systematically mixing temporal and spatial inos4 formation can further enhance model performance remains an open area for future investigation.

056 To reduce the model complexity in deep learning, parameter sharing (PS) is a crucial technique that can significantly reduce the amount model parameters and enhance computational efficiency. 058 In CNNs, convolutional filters share weights across spatial locations, capturing local features with fewer parameters. Similarly, LSTM networks share weight matrices across time steps to manage 060 memory and control information flow. By studying the sharing of attention weights, (Xiao et al., 061 2019; Kitaev et al., 2020) improved performance in the fields of natural language processing and 062 computer vision. ALBERT (Lan et al., 2020) further extends parameter sharing in natural language processing by sharing weights across transformer layers, reducing parameter redundancy while 063 maintaining performance. In multi-task learning, the Task Adaptive Parameter Sharing (TAPS) 064 approach (Wallingford et al., 2022) selectively fine-tunes task-specific layers while maximizing pa-065 rameter sharing across tasks, achieving efficient learning with minimal task-specific modifications. 066 Those studies demonstrate that parameter sharing has the potential for model size reduction, gen-067 eralization ability enhancement and mitigating the over-fitting risks across various tasks. Besides, 068 employing advanced optimization technique is also essential in reducing the overfitting issue. (Il-069 bert et al., 2024) proposed a simple yet effective transformer-based network structure and equipped it with the SAM (Sharpness-Aware Minimization) optimization mechanism to achieve competitive 071 performance compared with large models (Foret et al., 2021).

In this work, we explore innovative designs of transformer-based model for time series forecasting 073 by incorporating the characteristics of time series tasks and the concept of parameter sharing. Un-074 like MOIRAI (Woo et al., 2024) that "flattens" multivariate time series by treating all variables as 075 a single sequence, or SAMformer (Ilbert et al., 2024) that applies attention to the channel dimen-076 sion to capture spatial dependencies, we propose a different and novel approach. We construct a 077 transformer encoder with two-stage segment attention structure, where each network layer consists of a parameter-sharing block. This parameter-sharing block is composed of three fully connected layers with residual connections, which keeps the overall number of parameters in the model very 079 small while enabling effective information sharing across different parts of the model. For segment attention, we use patching to divide the variables into different patches, then identify the patches 081 at the same position across different variables and merge them into a segment. As a result, each segment represents the spatial extension of a single-variable patch. In this way, we decompose the 083 multivariate time series into multiple segments. Attention is applied within each segment to enhance 084 the extraction of local spatial-temporal relationships, while information fusion across segments is 085 performed to improve the overall predictive performance. Additionally, by incorporating the SAM optimization method, we further reduce over-fitting while maintaining efficient training. We con-087 duct extensive experiments on long-term time series forecasting datasets to verify the effectiveness of this segment attention structure with parameter sharing. Our model remains competitive when compared to previous state-of-the-art models, achieving the best performance on 7 out of 8 mainstream long-term time series forecasting tasks. The contributions are summarized as follows:

- We developed a novel transformer-based model structure for time series forecasting, where the parameter sharing technique is applied in the transformer block to reduce the model complexity and improve the generalization ability.
- We proposed a segment attention mechanism tailored for multi-variate data, which merges the temporal sequence of different channels to form a local segment and applies attention within each segment to capture both temporal dependencies and cross-channel interactions.
- Through extensive experiments in long-term forecasting tasks and ablation studies, we verified the effectiveness and superior performance of our proposed framework.
- 099 100

102

104

091

092

093

095

096

097

098

2 RELATED WORK

## 103 2.1 TEMPORAL MODELING IN TIME SERIES FORECASTING

In recent years, time series analysis has received widespread attention, with more deep learning
 methods being applied to time series forecasting. These deep learning methods focus on estab lishing temporal dependencies within time series data to predict future trends. The models can be
 broadly categorized into RNN-based, CNN-based, MLP-based, and Transformer-based approaches.

108 RNNs and their LSTM variants were widely used for time series tasks in the past, with related 109 works such as DeepAR(Salinas et al., 2020). CNN-based methods like TCN (Bai et al., 2018) and 110 TimesNet (Wu et al., 2023) have been designed to adapt convolutional structures specifically for 111 temporal modeling. MLP-based approaches, such as N-BEATS (Oreshkin et al., 2020), RLinear 112 (Li et al., 2023), and TSMixer (Chen et al., 2023b), have demonstrated that even simple network structures can achieve solid predictive performance. Moreover, Transformer-based models have 113 become increasingly popular in time series forecasting due to the unique attention mechanism of 114 Transformers, which provides strong global modeling capabilities. Many recent works leverage this 115 to enhance time series modeling performance, such as Informer (Zhou et al., 2021), Autoformer 116 (Wu et al., 2021), Pyraformer (Liu et al., 2022), and Fedformer (Zhou et al., 2022). Additionally, 117 PatchTST (Nie et al., 2023) further divides time series data into different patches to enhance the abil-118 ity to capture local information. However, the aforementioned models primarily focus on temporal 119 modeling, with less emphasis on modeling the relationships between variables. Although PatchTST 120 attempted to incorporate cross-channel designs, it observed degraded performance in their model. 121

On the other hand, some pre-trained large models have been applied to time series forecasting tasks (Das et al., 2023; Liu et al., 2024a; Gao et al., 2024; Liu et al., 2024c; Zhou et al., 2023; Jin et al., 2024). For example, MOMENT (Goswami et al., 2024) uses the patching method and mask pretraining to build a pre-trained model for time series, while GPT4TS (Zhou et al., 2023) also adopts the patching method and uses GPT2 as the backbone. The increase in model parameters has provided them with greater expressive power but also increased the difficulty of training.

127 128 129

130

145

## 2.2 VARIATE MODELING IN TIME SERIES FORECASTING

131 In addition to modeling temporal dependencies, recent works have focused on modeling inter-132 variable dependencies (donghao & wang xue, 2024; Zhang & Yan, 2023; Ilbert et al., 2024; Woo 133 et al., 2024; Liu et al., 2024b). ModernTCN (donghao & wang xue, 2024) employs different 1-D 134 convolutions to capture the temporal and variable dimensions separately; Crossformer (Zhang & 135 Yan, 2023) constructs attention mechanisms across both time and spatial dimensions; SAMformer (Ilbert et al., 2024) applies attention to the variable dimension to model cross-channel dependen-136 cies; MORAI (Woo et al., 2024) "flattens" multivariate time series into univariate sequences to 137 merge information across variables; iTransformer (Liu et al., 2024b) represents multivariate time se-138 ries and captures global dependencies. All of these works emphasize the simultaneous modeling of 139 both variable and temporal dependencies as critical directions for improving multivariate time series 140 modeling, which helps establish global spatial-temporal dependencies. However, this may weaken 141 the ability to capture local spatial-temporal dependencies. Additionally, expanding the global recep-142 tive field of spatio-temporal dependencies could increase model complexity, which in turn may lead 143 to overfitting due to the larger number of parameters. Our work addresses these issues and proposes 144 solutions to mitigate them.



Figure 1: Data transformation in PSformer. Multivariate time series data first undergo patching and
 cross-channel merging before being fed into the PSformer Encoder. Then after inverse transforma tion and linear mapping for the prediction length, the final time series results are obtained.



Figure 2: PSformer network structure. The PSformer includes the PS Block for parameter sharing within the PS Encoder structure, as well as the two-stage segment attention.

## **3** The PSformer Framework

## 3.1 PROBLEM FORMULATION

195 We denote the input multivariate time series  $X \in \mathbb{R}^{M \times L}$  as with M variables and look-back window 196  $L: (x_1, x_2, ..., x_L)$ , where  $x_t$  represents the M-dimensional vector at time step t. L will be equally 197 divided into N non-overlapping patches of size P.  $P^{(i)}$  denote the i-th patch with size P, where 198 i = 1, 2, 3, ..., N. The  $P^{(i)}$  of the M variables forms the i-th segment, which denote cross-channel 199 patch of length C, where  $C = M \times P$ . As shown in Figure 1, we transform X from  $X \in \mathbb{R}^{M \times L}$  to 200  $X \in \mathbb{R}^{C \times N}$  before feeding it into the model. We aim to predict the future values of F time steps, 201 e.g.,  $(x_{L+1}, ..., x_{L+F})$ . Beside, we denote  $X^{in}$  and  $X^{out}$  as the input and output signals for the 202 specified layers.

203 204

205

187

188

189 190 191

192 193

194

3.2 MODEL STRUCTURE

The PSformer model we constructed is depicted in Figure 2. The key components of the model include the PSformer Encoder, segment attention, and PS Block. The PSformer Encoder serves as the backbone of the model and contains both the segment attention and the PS Block. The PS Block provides the parameters for all layers within the Encoder, utilizing the parameter sharing technique.

Forward Process The univariate time series of length L for the *i*-th variable, starting at time index 1, is denoted as  $x_1^{(i)} = (x_1^{(i)}, ..., x_L^{(i)})$ , where i = 1, ..., M. Then the input  $(x_1, ..., x_L)$  with Mdimensions is presented as  $x_1 \in \mathbb{R}^{M \times L}$ , and  $x_1$  is used as the input to the transformer network structure. Similar to other time series forecasting methods, we use a RevIN layer Kim et al. (2022), which is added at both the input and output of the model.

## Spatial-Temporal Segment Attention

We introduce spatial-temporal segment attention (SegAtt), which merges patches from different channels at the same position into a segment and establishes spatial-temporal relationships across different segments. Specifically, the input time series  $X \in \mathbb{R}^{M \times L}$  is first divided into patches, where  $L = P \times N$ , and then transformed into  $X \in \mathbb{R}^{(M \times P) \times N}$ . By merging the M and P dimensions, the input becomes  $X \in \mathbb{R}^{C \times N}$ , where  $C = M \times P$ , facilitating the subsequent cross-channel information fusion.

222 In this transformed space, identical  $Q \in \mathbb{R}^{C \times N}$ ,  $K \in \mathbb{R}^{C \times N}$ , and  $V \in \mathbb{R}^{C \times N}$  matrices are 223 generated by applying a shared block's non-linear projection (applied by Parameter Shared Block), 224 with weights  $W^{\check{Q}} \in \mathbb{R}^{N \times N}$  and  $W^{K} \in \mathbb{R}^{N \times N}$  used to project the input  $X^{in} \in \mathbb{R}^{C \times N}$ . The 225 Q and K matrices are then multiplied using a dot-product operation to form the attention matrix 226  $QK^T \in \mathbb{R}^{C \times C}$ , which captures relationships between different spatial-temporal segments (in the C dimension) and is used to act on the V matrix. While the computation of Q, K and V involves 227 non-linear transformations of the input  $X^{in}$  across segments in the N dimension (corresponding to 228 the  $d_{model}$  dimension of attention in the NLP domain), the scaled dot-product attention primarily 229 applies attention across the C dimension, allowing the model to focus on dependencies between 230 spatial-temporal segments across channels and time. 231

This mechanism ensures that information from different segments is integrated through the computation of Q, K, and V. It also captures local spatial-temporal dependencies within individual segments by applying attention to the internal structure of each segment. Additionally, it captures long-term dependencies across segments over larger time steps. The final output is  $X^{out} \in \mathbb{R}^{C \times N}$ , completing the attention process.

**Parameter Shared Block** In our work, we propose a novel Parameter Shared Block (PS Block), which consists of three fully connected layers with residual connections, as illustrated in Figure 2. Specifically, we construct three trainable linear mappings  $W^{(j)} \in \mathbb{R}^{N \times N}$  with  $j \in \{1, 2, 3\}$ . The output of the first two layers is computed as:

$$\boldsymbol{X}^{out} = \text{GeLU}(\boldsymbol{X}^{in}\boldsymbol{W}^{(1)})\boldsymbol{W}^{(2)} + \boldsymbol{X}^{in}, \tag{1}$$

which follows a similar structure as the feed-forward network (FFN) with residual connections. This intermediate output  $X^{out}$  is then used as the input for the third transformation, yielding:

$$\boldsymbol{X}^{out} = \boldsymbol{X}^{in} \boldsymbol{W}^{(3)}$$

248 Therefore, the PS Block as a whole can be expressed as:

241 242 243

244

245 246 247

249 250

251

$$X^{out} = (\text{GeLU}(X^{in}W^{(1)})W^{(2)} + X^{in})W^{(3)},$$
(2)

and we denote PS Block output as  $X^{out} = X^{in} W^S$ , where  $W^S \in \mathbb{R}^{N \times N}$  and  $X^{out} \in \mathbb{R}^{C \times N}$ . The 252 structure of the PS Block allows it to perform nonlinear transformations while preserving a linear 253 transformation path. Although the three layers within the PS Block have different parameters, the 254 entire PS Block is reused across different positions in the PS former Encoder, ensuring that the same 255 block parameters  $W^S$  are shared across those positions, as illustrated in Figure 2. Specifically, PS 256 Block share parameters across three parts of each PSformer Encoder, which including two segment 257 attention layer and the final PS Block. In the segment attention layer, the PS block outputs are used 258 as the Q, K, and V matrices to construct the attention mechanism. This parameter-sharing strategy 259 reduces the overall number of parameters while maintaining the network's expressive capacity. 260

261 **PSformer Encoder** In the PSformer Encoder, as Figure 2, each layer shares the same parameters 262  $W^S$  of PS Block. For the input  $X^{in}$ , the transformation in the PSformer Encoder can be expressed 263 as follows:

264 SegAtt stage one is represented as:  $Q^{(1)} = X^{in}W^S$ ,  $K^{(1)} = X^{in}W^S$ ,  $V^{(1)} = X^{in}W^S$ , 265 Therefore,  $Q^{(1)}$ ,  $K^{(1)}$ ,  $V^{(1)} \in \mathbb{R}^{C \times N}$ . Dot-product attention operation

$$O^{(1)} = \text{Attention}^{(1)}(Q^{(1)}, K^{(1)}, V^{(1)}) = \text{Softmax}(\frac{Q^{(1)}K^{(1)}}{\sqrt{d_k}})V^{(1)}.$$

then with ReLU activation:  $O_{act}^{(1)} = ReLU(O^{(1)})$ .

270 271 SegAtt stage two is represented as:  $Q^{(2)} = O_{act}^{(1)}W^S$ ,  $K^{(2)} = O_{act}^{(1)}W^S$ ,  $V^{(2)} = O_{act}^{(1)}W^S$ , Therefore,  $Q^{(2)}$ ,  $K^{(2)}$ ,  $V^{(2)} \in \mathbb{R}^{C \times N}$ . Dot-product attention operation

$$O^{(2)} = \text{Attention}^{(2)}(Q^{(2)}, K^{(2)}, V^{(2)}) = \text{Softmax}(\frac{Q^{(2)}K^{(2)}}{\sqrt{d_k}})V^{(2)}$$

274 275

282

291

292 293

315 316

317

273

The two-stage SegAtt mechanism can be viewed as analogous to an FFN layer, where the MLP is replaced with attention operations. Additionally, residual connections are introduced between the input and output, and the result is then fed into the final PS Block.

The transformation in the final PS Block is represented as  $O^{out} = O^{in}W^S$ . Therefore, the entire encoder can be expressed as

$$\boldsymbol{X}^{out} = (\text{Attention}^{(2)}(\text{ReLU}(\text{Attention}^{(1)}(\boldsymbol{X}^{in}))) + \boldsymbol{X}^{in})\boldsymbol{W}^{S},$$

with  $X^{out} \in \mathbb{R}^{C \times N}$ . Then we apply a dimensionality transformation to obtain  $X^{out} \in \mathbb{R}^{M \times L}$ , since  $C = M \times P$  and  $L = P \times N$ .

After passing through *n* layers of the PSformer Encoder, the final output is  $X^{pred} = X^{out}W^F$ , where  $X^{pred} \in \mathbb{R}^{M \times F}$ , and  $W^F \in \mathbb{R}^{L \times F}$  is a linear mapping, where *F* is the prediction length. The  $X^{pred}$  is the final output of the PSformer model. The PSformer structure does not use positional encoding, as the segment attention mixes local spatiotemporal information. We discuss this in more detail in Appendix A.6 and Appendix B.8.

## 4 EXPERIMENT

Datasets In this paper, we focus on the long-term time series forecasting. We follow the time series forecasting work in Ilbert et al. (2024) and use 8 mainstream datasets to evaluate the performance of our proposed PSformer model. As shown in Table 1, these datasets include 4 ETT datasets (ETTh1, ETTh2, ETTm1, ETTm2), as well as Weather, Traffic, Electricity, and Exchange. These datasets have been used as benchmark evaluations in many previous time series forecasting studies.

299 Baselines We select state-of-the-art (SOTA) models in the field of long-term time series forecasting, including not only Transformer-based models but also large models and other SOTA mod-300 els. Specifically, baselines include (1) Transformer-based model: SAMformer (Ilbert et al., 2024), 301 iTransformer (Liu et al., 2024b), PatchTST (Nie et al., 2023), FEDformer (Zhou et al., 2022), Auto-302 former (Wu et al., 2021), (2) Pretrained Large model: MOMENT (Goswami et al., 2024), GPT4TS 303 (Zhou et al., 2023), (3) TCN-based model: ModernTCN (donghao & wang xue, 2024) and (4) MLP-304 based methods: TSMixer (Chen et al., 2023b), RLinear (Li et al., 2023). Additionally, we provide 305 more baselines for a comprehensive comparison, including TimeMixer (Wang et al., 2024a), Cross-306 GNN (Huang et al., 2023), MICN (Wang et al., 2023), TimesNet (Wu et al., 2023), FITS (Xu et al., 307 2024), Crossformer (Zhang & Yan, 2023), PDF (Dai et al., 2024), and TimeLLM (Jin et al., 2024) 308 Further details about these baselines can be found in Appendix B.9.

**Experimental Settings** The input time series length T is set to 512, and four different prediction lengths  $H \in \{96, 192, 336, 720\}$  are used. Evaluation metrics include Mean Squared Error (MSE) and Mean Absolute Error (MAE). We train our constructed models using the SAM optimization technique as in Ilbert et al. (2024). We set the look-back window of RevIN to 16 for the Exchange dataset, more details about this setting can be found in Appendix B.7. We collect baseline results of SAMformer, TSMixer, FEDformer and Autoformer from Ilbert et al. (2024), in which SAMformer

Table 1: Datasets for long-term forecasting. Dataset size is structured as (Train, Validation, Test).

Dataset	Variate	Predict Length	Frequency	Dataset Size	Information
ETTh1,ETTh2	7	{96,192,336,720}	Hourly	(8545, 2881, 2881)	Electricity
ETTm1,ETTm2	7	{96,192,336,720}	10 mins	(34465, 11521, 11521)	Electricity
Weather	21	{96,192,336,720}	15 mins	(36792, 5271, 10540)	Weather
Electricity	321	{96,192,336,720}	Hourly	(18317, 2633, 5261)	Electricity
Exchange	8	{96,192,336,720}	Daily	(5120, 665, 1422)	Exchange rate
Traffic	862	{96,192,336,720}	Hourly	(12185, 1757, 3509)	Transportation

Metric	PSfo	rmer	SAM	ormer	TSM	fixer	Pate	hTST	MON	1ENT	Mode	mTCN	FED	ormer	GPI	4TS	Autof	ormer	RLi	near	iTrans	forme
	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAI
ETTh1	0.397	0.418	<u>0.410</u>	<u>0.424</u>	0.420	0.431	0.468	0.455	0.418	0.436	0.421	0.432	0.428	0.454	0.428	0.426	0.473	0.477	0.446	0.434	0.454	0.44
ETTh2	0.338	0.390	0.344	<u>0.391</u>	0.354	0.400	0.387	0.407	0.352	0.394	<u>0.343</u>	0.393	0.388	0.434	0.355	0.395	0.422	0.443	0.374	0.398	0.383	0.40
ETTm1	0.342	0.372	0.373	0.388	0.378	0.392	0.387	0.400	<u>0.344</u>	<u>0.379</u>	0.361	0.384	0.382	0.422	0.351	0.383	0.515	0.493	0.414	0.407	0.407	0.41
ETTm2	0.251	0.313	0.269	0.327	0.283	0.339	0.281	0.326	<u>0.259</u>	<u>0.318</u>	0.262	0.322	0.292	0.343	0.267	0.326	0.310	0.357	0.286	0.327	0.288	0.33
Weather	0.225	0.264	0.261	0.293	0.255	0.289	0.259	0.281	<u>0.228</u>	<u>0.270</u>	0.237	0.274	0.310	0.357	0.237	0.271	0.335	0.379	0.272	0.291	0.258	0.27
Electricity	<u>0.162</u>	0.255	0.181	0.275	0.198	0.296	0.205	0.290	0.165	<u>0.260</u>	0.160	0.255	0.207	0.321	0.167	0.263	0.214	0.327	0.219	0.298	0.178	0.27
Exchange	0.358	0.399	0.445	0.470	0.532	0.523	0.367	0.404	0.437	0.446	0.555	0.536	0.478	0.477	0.371	0.409	0.613	0.539	0.378	0.417	<u>0.360</u>	<u>0.40</u>
Traffic	0.400	0.274	0.425	0.297	0.439	0.315	0.481	0.304	0.415	0.293	<u>0.414</u>	0.283	0.604	0.372	<u>0.414</u>	0.295	0.617	0.384	0.626	0.378	0.428	0.28
Count	7	8	0	0	0	0	1	0	0	0	1	1	0	0	0	0	0	0	0	0	1	1

324 Table 2: Long-term forecasting task. All the results are averaged from 4 different prediction lengths 325 {96, 192, 336, 720}. A lower MSE or MAE indicates a better performance. See Table B.2 in 326 Appendix for the full results with more baselines.

and TSMixer are also trained with SAM, and the results of iTransformer, RLinear and PatchTST from Liu et al. (2024b). We test ModernTCN with the fixed input length of 512 following donghao & wang xue (2024). For the pre-trained large models MOMENT and GPT4TS, the results are collected from (Goswami et al., 2024). More details about the baseline settings can be found in Appendix A.2. Besides, We also provide results under uni-variate series and unrelated variate in the Appendix B.6.

4.1 **RESULTS AND ANALYSIS** 

348 Summary of Experimental Results The main experimental results are reported in Table 2. PS-349 former achieved the best performance on 7 out of 8 major datasets in long-term time series fore-350 casting tasks, demonstrating its strong predictive capabilities across various time series prediction 351 problems. This success is attributed to its segment attention mechanism, which enhances the ex-352 traction of spatial-temporal information, and the parameter-sharing structure, which improves the 353 model's robustness. The complete experimental results are in Appendix B.2. We also provide additional visualization results to analyze the attention maps across spatial and temporal dimensions in 354 the Appendix B.5. 355

356 **Comparison with Other SOTA models** Compared to other Transformer-based models, PSformer 357 demonstrates higher predictive accuracy, reflecting the effectiveness of dividing multivariate time 358 series data into spatial-temporal segments for attention calculation. The neural network's ability to extract information from all spatial-temporal segments enhances the prediction performance. In 359 contrast to current large pre-trained models, PSformer not only achieves better accuracy but also 360 reduces the amount of parameters through parameter sharing. Although linear models are simpler 361 and have fewer parameters with satisfactory performance in some cases, the ability to extract rich 362 information from complex data is limited. In contrast, PSformer integrates the residual connections, 363 thus enabling a linear path while retaining the capability to process complex nonlinear information. 364 Moreover, the ConvFFN component in ModernTCN tailored for temporal data also confirms that the convolutional mechanism, which actually embodies the idea of parameter sharing, is also effective 366 in the time series domain. With the same spirit, we have successfully applied the parameter sharing 367 to the transformer-based models in the time series field and achieved superior performance. 368

Comparison with PatchTST and SAMformer PatchTST employs a channel-independent design 369 and divides time series data into multiple patches, which demonstrates the effectiveness of the patch-370 ing method in time series processing. However, its channel-independent approach does not fully 371 consider the relationships between different channels, focusing only on processing each channel 372 individually. On the other hand, SAMformer applies attention directly to the channel dimension 373 via dimension transformations and utilizes a simplified model structure, achieving good predic-374 tive performance. However, it may fail to capture valuable local information without the patching 375 method. PSformer combines the advantages of both models while addressing their limitations. By using spatial-temporal segment attention mechanism, PSformer effectively captures local temporal 376 information and handles relationships among different channels. This design enables PS former to 377 outperform them in various time series forecasting tasks as validated by extensive experiments.

327 328

340

341

342

343

344

345 346

347

Num of comonte		ET	Th1		ETTm1					
Num. of segments	96	192	336	720	96	192	336	720		
8	0.362	0.406	0.670	0.451	0.290	0.325	0.357	0.412		
16	0.352	0.417	0.424	<u>0.446</u>	<u>0.285</u>	<u>0.323</u>	<u>0.353</u>	<u>0.412</u>		
32	0.352	0.386	0.410	0.440	0.282	0.321	0.352	0.413		
64	<u>0.354</u>	<u>0.389</u>	<u>0.412</u>	<u>0.446</u>	0.288	0.325	0.355	0.417		

Table 3: The MSE results of different number of segments (N) on the ETTh1 and ETTm1 dataset.

## 4.2 Ablation studies

378

388

389 The Effect of segments Numbers Since PS former employs a non-overlapping patching to construct 390 segments, the model's performance is affected by the number of segments. Therefore, we tested the 391 model's performance with different segment numbers on two datasets, ETTh1 and ETTm1. Given 392 that the input sequence length is fixed at 512, the number of segments must be a divisor of the 393 sequence length. Consequently, we set the number of segments to 8, 16, 32, and 64, and test the 394 model on four different forecasting horizons. The test results are shown in Table 3, which indicate 395 that the number of segments impacts the model's prediction accuracy to some extent. Across both 396 datasets, a moderate number of segments (such as 32) tends to perform the best, particularly for 397 both short and long forecasting horizons. For shorter horizons (96, 192), using 16 or 32 segments generally yields the lowest MSE, while fewer segments (8) often result in poorer performance. As 398 the horizon increases (336, 720), 32 segments consistently lead to optimal results, indicating that 399 this number balances the extraction of both local and global temporal features effectively. Besides, 400 we provide results and analysis about the impact of the segment number across more datasets in the 401 Appendix B.13. 402

The Effect of PSformer Encoder Numbers Since PSformer adopts the segment attention mecha-403 nism, with non-shared PS Block parameters across different Encoders (shared within Encoder), we 404 tested the impact of varying the number of encoders on model performance. We conducted tests 405 on the ETTh1 and ETTm1 datasets, varying the number of encoders from 1 to 4 with four differ-406 ent forecasting horizons. The experimental results are shown in Table 4. The results indicate that 407 ETTm1 performs best with 3 encoders, while ETTh1 achieves better performance with just 1 en-408 coder. This may suggest that for smaller datasets, fewer encoders result in better performance, as 409 reducing the number of encoders decreases the amount of model parameters, thereby mitigating the 410 risk of over-fitting. 411

Table 4: The MSE results of different number of encoders on the ETTh1 and ETTm1 dataset.

Num of Encodore		ET	Th1		ETTm1				
Num. of Encouers	96	192	336	720	96	192	336	720	
1	0.352	0.385	0.411	0.440	0.288	0.324	0.356	0.414	
2	0.355	0.392	0.418	<u>0.443</u>	0.284	0.323	0.356	0.415	
3	0.355	0.391	<u>0.416</u>	<u>0.443</u>	0.282	0.321	0.352	0.413	
4	<u>0.355</u>	<u>0.389</u>	<u>0.416</u>	0.440	0.282	0.321	<u>0.353</u>	<u>0.416</u>	

419 420 421

412

Ablation of Parameter Sharing methods We investigate the impact of parameter-sharing mecha-422 nism on the model performance. In addition to the default parameter-sharing approach, which shares 423 parameters only within encoder (In-Encoder), we also test the following approaches: a. no pa-424 rameter sharing, i.e., None; b. sharing parameters only across encoders (with different parameters 425 used for the PS blocks within each encoder), i.e., Cross-Encoders; and c. sharing parameters 426 both within and across encoders, i.e., ALL. We conducted experiments on the ETTm1 and Weather 427 datasets, and the results are shown in Table 5. As can be seen, the In-Encoder method per-428 forms the best, followed by ALL, while None shows the worst performance. This indicates that the 429 parameter-sharing mechanism contributes to improving the model performance. Furthermore, we provide a comparison of convergence rates between parameter sharing and non-parameter sharing 430 in Appendix B.10, as well as a comparison of parameter savings achieved by parameter sharing in 431 Appendix B.11 and Appendix A.4.

Sharing N	Methods	Cross-Encoders	In-Encoder	ALL	None
	96	0.297	0.282	0.288	0.295
ETTm 1	196	0.329	0.321	0.326	0.338
EIImI	336	0.360	0.352	0.358	0.372
	720	0.420	0.413	0.414	0.425
	96	0.158	0.149	0.151	0.154
Waathan	196	0.201	0.193	0.196	0.198
weather	336	0.252	0.245	0.245	0.245
	720	0.319	0.314	<u>0.316</u>	0.319

 Table 5: Parameter Sharing with different methods

Table 6: Ablation analysis of SegAtt on ETTh1 and ETTh2 (MSE reported). Red: the best.

Variant	1	ET	Гh1		ETTh2					
variant	96	192	336	720	96	192	336	720		
w SegAtt (Default)	0.352	0.385	0.411	0.440	0.272	0.335	0.356	0.389		
w/o SegAtt	<u>0.369</u>	<u>0.397</u>	<u>0.414</u>	<u>0.448</u>	0.288	0.365	0.373	0.398		
w channel independent	0.376	0.407	0.427	0.455	0.285	0.382	0.369	<u>0.395</u>		
only mlp	0.379	0.399	0.426	0.450	<u>0.282</u>	<u>0.352</u>	<u>0.358</u>	0.398		

Ablation of Segment Attention methods As an important component of PSformer, we designed segment attention to effectively capture temporal information while fully utilizing the correlations between channels. In the ablation study, we compared different attention mechanisms without using SegAtt under parameter sharing, as well as the performance of using only MLP. We conducted experiments on the ETTh1 and ETTh2 datasets. In this study, w/o SegAtt means applying the attention mechanism to the input  $x \in \mathbb{R}^{B \times M \times L}$  on the  $M \times L$  dimensions. For the setting with channel independence' applies the attention mechanism on the  $N \times P$  dimensions. The results in Table 6 show the performance of PSformer variants with different attention mechanisms on the ETTh1 and ETTh2 datasets across various forecasting horizons:

- The default PSformer configuration (with SegAtt) consistently achieves the lowest MSE across all horizons, demonstrating the effectiveness of segment attention.
- When SegAtt is removed, the cross-channel attention is less effective at capturing both local and global temporal dependencies, resulting in slightly increased MSE.
- The variant with channel-independent attention shows further degradation, suggesting that neglecting inter-channel correlations impacts the model's ability to capture temporal features.
- Using only MLP layers also results in higher MSE. Although it performs second best for the three forecasting horizons on ETTh2, it still falls short compared to the performance with SegAtt, highlighting the necessity of applying SegAtt.

## 5 CONCLUSION AND FUTURE WORK

In this work, we proposed the PS former model for multivariate time series forecasting, which lever-ages segment attention technique to facilitate information transfer among time series variables and capture spatial-temporal dependencies. By employing parameter sharing, the model effectively improves parameter efficiency and reduces the risk of overfitting when the data size is relatively limited. This network structure, which combines parameter sharing with segment attention mecha-nism, achieves state-of-the-art performance on long-term multivariate forecasting tasks by enhanc-ing model parameter efficiency and improving the utilization of channel-wise information. Future work can focus on applying this parameter sharing and segment attention technique to the develop-ment of pre-trained large models for time series forecasting, to overcome the issue of excessively large parameter counts in existing pre-trained models, and improve the capability of extracting in-formation from multivariate time series. Additionally, exploring model architecture designs with parameter sharing to improve prediction performance still holds potential value.

# 486 REFERENCES

- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- Ching Chang, Wen-Chih Peng, and Tien-Fu Chen. LLM4TS: Aligning pre-trained llms as dataefficient time-series forecasters. *arXiv preprint arXiv:2308.08469*, 2023.
- Liuyi Chen, Bocheng Han, Xuesong Wang, Jiazhen Zhao, Wenke Yang, and Zhengyi Yang. Machine
  learning methods in weather and climate applications: A survey. *Applied Sciences*, 13(21), 2023a.
  ISSN 2076-3417. doi: 10.3390/app132112019.
- 496
   497
   497
   498
   498
   498
   498
   498
   499
   499
   499
   499
   490
   490
   490
   491
   492
   493
   494
   495
   495
   495
   496
   496
   497
   498
   498
   499
   499
   490
   490
   490
   490
   490
   491
   491
   492
   493
   494
   494
   495
   495
   495
   496
   496
   497
   498
   498
   498
   499
   499
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
   400
- Tao Dai, Beiliang Wu, Peiyuan Liu, Naiqi Li, Jigang Bao, Yong Jiang, and Shu-Tao Xia. Periodicity decoupling framework for long-term series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=dp27P5HBBt.
- Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
   time-series forecasting. *arXiv preprint arXiv:2310.10688*, 2023.
- Luo donghao and wang xue. ModernTCN: A modern pure convolution structure for general time series analysis. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=vpJMJerXHU.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
   Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recogni tion at scale. In *International Conference on Learning Representations*, 2021. URL https:
   //openreview.net/forum?id=YicbFdNTTy.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimiza tion for efficiently improving generalization. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=6TmlmposlrM.
- Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and Marinka Zitnik. UniTS: Building a unified time series model. *arXiv preprint arXiv:2403.00131*, 2024.
- Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski.
   MOMENT: A family of open time-series foundation models. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=
   FVvf69a5rx.
- Qihe Huang, Lei Shen, Ruixin Zhang, Shouhong Ding, Binwu Wang, Zhengyang Zhou, and Yang Wang. CrossGNN: Confronting noisy multivariate time series via cross interaction refinement. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL https: //openreview.net/forum?id=x0z1W2vUYc.
- 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. URL https://openreview.net/forum?id=8kLzL5QBh2.
- Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y. Zhang, Xiaoming Shi, Pin-Yu Chen,
   Yuxuan Liang, Yuan-Fang Li, Shirui Pan, and Qingsong Wen. Time-LLM: Time series forecasting
   by reprogramming large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=Unb5cVPtae.

540 Anwar Khan, Mostafa M Fouda, Dinh-Thuan Do, Abdulaziz Almaleh, and Atiq Ur Rahman. Short-541 term traffic prediction using deep learning long short-term memory: Taxonomy, applications, 542 challenges, and future trends. IEEE Access, 11:94371-94391, 2023. 543 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-544 versible instance normalization for accurate time-series forecasting against distribution shift. In 545 International Conference on Learning Representations, 2022. URL https://openreview. 546 net/forum?id=cGDAkQo1C0p. 547 548 Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In 549 International Conference on Learning Representations, 2020. URL https://openreview. 550 net/forum?id=rkqNKkHtvB. 551 Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Sori-552 cut. Albert: A lite bert for self-supervised learning of language representations. In International 553 Conference on Learning Representations, 2020. URL https://openreview.net/forum? 554 id=H1eA7AEtvS. 555 556 Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting long-term time series forecasting: An investigation on linear mapping, 2023. URL https://arxiv.org/abs/2305.10721. 558 Yuxuan Liang, Haomin Wen, Yuqi Nie, Yushan Jiang, Ming Jin, Dongjin Song, Shirui Pan, and 559 Qingsong Wen. Foundation models for time series analysis: A tutorial and survey. In Proceedings 560 of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 6555-561 6565, 2024. 562 563 Che Liu, Bo Sun, Chenghui Zhang, and Fan Li. A hybrid prediction model for residential electricity 564 consumption using holt-winters and extreme learning machine. Applied Energy, 275:115383, 2020. ISSN 0306-2619. 565 566 Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dust-567 dar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and 568 forecasting. In International Conference on Learning Representations, 2022. URL https: 569 //openreview.net/forum?id=0EXmFzUn5I. 570 571 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 572 Proceedings of the ACM on Web Conference 2024, pp. 4095–4106, 2024a. 573 574 Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. 575 iTransformer: Inverted transformers are effective for time series forecasting. In The Twelfth In-576 ternational Conference on Learning Representations, 2024b. URL https://openreview. 577 net/forum?id=JePfAI8fah. 578 Yong Liu, Haoran Zhang, Chenyu Li, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. 579 Timer: Transformers for time series analysis at scale. arXiv preprint arXiv:2402.02368, 2024c. 580 581 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 582 Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of 583 the IEEE/CVF International Conference on Computer Vision (ICCV), pp. 10012–10022, October 584 2021. 585 Noella Nazareth and Yeruva Venkata Ramana Reddy. Financial applications of machine learning: 586 A literature review. Expert Systems with Applications, 219:119640, 2023. ISSN 0957-4174. 587 588 Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 589 64 words: Long-term forecasting with transformers. In International Conference on Learning 590 Representations, 2023. 591 Isaac Kofi Nti, Moses Teimeh, Owusu Nyarko-Boateng, and Adebayo Felix Adekoya. Electricity 592 load forecasting: a systematic review. Journal of Electrical Systems and Information Technology, 593 7:1-19, 2020.

- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, and Diogo Almeida. Gpt-4 technical report, 2024.
- Boris N. Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis expansion analysis for interpretable time series forecasting. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=rlecqn4YwB.
- Xiaoli Ren, Xiaoyong Li, Kaijun Ren, Junqiang Song, Zichen Xu, Kefeng Deng, and Xiang Wang.
   Deep learning-based weather prediction: A survey. *Big Data Research*, 23:100178, 2021. ISSN 2214-5796. doi: https://doi.org/10.1016/j.bdr.2020.100178.
- David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International journal of forecasting*, 36(3):1181– 1191, 2020.
- David Alexander Tedjopurnomo, Zhifeng Bao, Baihua Zheng, Farhana Murtaza Choudhury, and
  Alex Kai Qin. A survey on modern deep neural network for traffic prediction: Trends, methods and challenges. *IEEE Transactions on Knowledge and Data Engineering*, 34(4):1544–1561,
  2020.
- Matthew Wallingford, Hao Li, Alessandro Achille, Avinash Ravichandran, Charless Fowlkes, Rahul
   Bhotika, and Stefano Soatto. Task adaptive parameter sharing for multi-task learning. In 2022
   *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7551–7560,
   2022. doi: 10.1109/CVPR52688.2022.00741.
- Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. MICN: Multi-scale local and global context modeling for long-term series forecasting. In *The Eleventh Interna-tional Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=zt53IDUR1U.
- Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y. Zhang, and JUN ZHOU. Timemixer: Decomposable multiscale mixing for time series forecasting. In *The Twelfth International Conference on Learning Representations*, 2024a. URL https://openreview.net/forum?id=7oLshfEIC2.
- Yuxuan Wang, Haixu Wu, Jiaxiang Dong, Yong Liu, Mingsheng Long, and Jianmin Wang. Deep time series models: A comprehensive survey and benchmark. *arXiv preprint arXiv:2407.13278*, 2024b.
- Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sa hoo. Unified training of universal time series forecasting transformers. arXiv preprint arXiv:2402.02592, 2024.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with Auto-Correlation for long-term series forecasting. In *Advances in Neural Information Processing Systems*, 2021.
- Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet:
   Temporal 2d-variation modeling for general time series analysis. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?
   id=ju\_Uqw3840q.
- Tong Xiao, Yinqiao Li, Jingbo Zhu, Zhengtao Yu, and Tongran Liu. Sharing attention weights for fast transformer. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 5292–5298. International Joint Conferences on Artificial Intelligence
   Organization, 7 2019. doi: 10.24963/ijcai.2019/735. URL https://doi.org/10.24963/ ijcai.2019/735.
- <sup>643</sup>
   <sup>644</sup>
   <sup>644</sup>
   <sup>645</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>647</sup>
   <sup>647</sup>
   <sup>648</sup>
   <sup>648</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>641</sup>
   <sup>641</sup>
   <sup>642</sup>
   <sup>642</sup>
   <sup>643</sup>
   <sup>643</sup>
   <sup>644</sup>
   <sup>645</sup>
   <sup>645</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>647</sup>
   <sup>647</sup>
   <sup>648</sup>
   <sup>648</sup>
   <sup>648</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>640</sup>
   <sup>641</sup>
   <sup>642</sup>
   <sup>645</sup>
   <sup>645</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>647</sup>
   <sup>647</sup>
   <sup>648</sup>
   <sup>648</sup>
   <sup>648</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>641</sup>
   <sup>642</sup>
   <sup>642</sup>
   <sup>645</sup>
   <sup>645</sup>
   <sup>645</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>647</sup>
   <sup>647</sup>
   <sup>648</sup>
   <sup>648</sup>
   <sup>649</sup>
   <sup>641</sup>
   <sup>642</sup>
   <sup>642</sup>
   <sup>645</sup>
   <sup>645</sup>
   <sup>645</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>646</sup>
   <sup>647</sup>
   <sup>647</sup>
   <sup>648</sup>
   <sup>648</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>649</sup>
   <sup>641</sup>
   <sup>642</sup>
   <sup>642</sup>
   <sup>643</sup>
   <sup>644</sup>
- 647 Zahra Zamanzadeh Darban, Geoffrey I Webb, Shirui Pan, Charu Aggarwal, and Mahsa Salehi. Deep learning for time series anomaly detection: A survey. *ACM Computing Surveys*, 2022.

648 649 650 651 652	Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series fore- casting? In Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Sym- posium on Educational Advances in Artificial Intelligence, AAAI'23/IAAI'23/EAAI'23. AAAI Press, 2023. ISBN 978-1-57735-880-0. doi: 10.1609/aaai.v37i9.26317.
653 654 655 656	Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=vSVLM2j9eie.
657 658 659	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, number 12, pp. 11106–11115, 2021.
660 661 662	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> , pp. 27268–27286. PMLR, 2022.
664 665 666	Tian Zhou, Peisong Niu, Xue Wang, Liang Sun, and Rong Jin. One fits all: Power general time series analysis by pretrained LM. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. URL https://openreview.net/forum?id=gMS6FVZvmF.
667 668 669	
670 671 672	
673 674 675	
676 677 678	
679 680 681 682	
683 684 685	
686 687 688	
689 690 691	
692 693 694	
695 696 697	
698 699 700 701	

#### 702 EXPERIMENTAL CONFIGURATION А 703

### A.1 HARDWARE 705

All experiments were conducted on two servers, each equipped with an 80GB NVIDIA A100 GPU and 4 Intel Xeon Gold 5218 CPUs.

## A.2 DETAILS OF BASELINE SETTINGS

711 We conducted all of our experiments with look-back window L = 512 and prediction horizons  $H \in \{96, 192, 336, 720\}$ . Results of PSformer and ModernTCN reported in Table 11 come from 712 our own experiments. The difference between our experiment with ModernTCN and its official 713 code is that we standardized the look-back window to 512 and set drop\_last=False for the test set 714 in the dataloader to ensure consistency with our experimental settings for a fair comparison. For 715 MOMENT and GPT4TS, the results are collected from (Goswami et al., 2024), except Exchange 716 rate (which is tested on official code repository). The results of SAMformer, TSMixer, FEDformer 717 and Autoformer are obtained from (Ilbert et al., 2024), while the results of iTransformer, PatchTST, 718 and RLinear are taken from (Liu et al., 2024b).

719 720

721

704

706

708 709

710

A.3 SETTINGS FOR PSFORMER

722 We provides an overview of the experimental configurations for the PSFormer model across various 723 tasks and datasets in Table 7. The experiments cover multiple time-series datasets, including ETTh1, 724 ETTh2, ETTm1, ETTm2, Weather, Electricity, Traffic, and Exchange. In all experiments, the input 725 sequence length (Input Length T) is set to 512, and each input is equally divided into 32 nonoverlapping segments (Segments Num. N). The model architecture uses 3 Encoders for tasks, 726 while for the ETTh1, ETTh2, ETTm2 and Exchange, 1 Encoder are used. 727

728 The learning rate adjustment strategy (schedule) is set to "constant" for all experiments, with a fixed 729 learning rate (LR) of  $10^{-4}$ . The loss function used in all experiments is Mean Squared Error (MSE). 730 The batch size is set to 16 for most tasks, except for the Traffic dataset, where the batch size is 8. 731 Each experiment runs for 300 epochs, with a patience value of 30 for early stopping. A fixed random 732 seed of 1 is applied across all experiments to ensure reproducibility.

733 734

Table 7: An o	overview of the e	experimental c	configurations	for <b>PSFORMER</b>
---------------	-------------------	----------------	----------------	---------------------

Task-Dataset	Encoder Num.	Input Length T	Segment Num. $N$	schedule	LR*	Loss	Batch Size	Epochs	patient	random seed
ETTh1	1	512	32	constant	$10^{-4}$	MSE	16	300	30	1
ETTh2	1	512	32	constant	$10^{-4}$	MSE	16	300	30	1
ETTm1	3	512	32	constant	$10^{-4}$	MSE	16	300	30	1
ETTm2	1	512	32	constant	$10^{-4}$	MSE	16	300	30	1
Weather	3	512	32	constant	$10^{-4}$	MSE	16	300	30	1
Electricity	3	512	32	constant	$10^{-4}$	MSE	16	300	30	1
Traffic	3	512	32	constant	$10^{-4}$	MSE	8	300	30	1
Exchange	1	512	32	constant	$10^{-4}$	MSE	16	300	30	1

741 742

## 743

744 745

## A.4 MODEL SIZE COMPARISON

Table 8 presents a comparison of the parameter size across different models, including the PS-746 Former and other baseline models such as SAMformer, TSMixer, ModernTCN, and RLinear. 747 The comparison is conducted on ETTh1, Weather, and Traffic datasets, with prediction horizons 748  $H \in \{96, 192, 336, 720\}$ . PSFormer is evaluated in two configurations: the full model and the en-749 coder part. The parameters of the encoder part refer to the number of parameters after excluding the 750 linear mapping in the output layer. The table denote that both PSFormer (full) and SAMformer have 751 parameter sizes that are close to RLinear, where RLinear's parameters mainly stem from the linear 752 mapping between input and output. Notably, the parameter sizes of these three models are relatively 753 unaffected by the number of input channels. In contrast, TSMixer and ModernTCN exhibit significantly larger parameter sizes, with the number of input channels playing a major role in the overall 754 parameter burden. The relative size comparison shows that TSMixer and ModernTCN have several 755 times, or even thousands of times, more parameters than PSFormer(full). Finally, the parameter size

Dataset	Horizon	PSfo	rmer	SAMformer	TSMixer	ModernTCN	RLinear
2		Full	Encoder				
	H=96	52,416	3,168	50,272	124,142	876,064	49248
<b>ETTb</b> 1	H=192	101,664	3,168	99,520	173,390	1,662,592	98496
CIIII	H=336	175,536	3,168	173,392	247,262	2,842,384	172368
	H=720	372,528	3,168	369,904	444,254	5,988,496	369360
Relative S	Size (Avg)	1.0	0.014	0.987	1.408	16.192	0.948
	H=96	58,752	9,504	50,272	121,908	2,709,280	49248
W	H=192	108,000	9,504	99,520	171,156	3,495,808	98496
weather	H=336	181,872	9,504	173,392	245,028	4,675,600	172368
	H=720	378,864	9,504	369,904	442,020	7,821,712	369360
Relative S	Size (Avg)	1.0	0.039	0.953	1.347	25.708	0.948
	H=96	58,752	9,504	50,272	793,424	822,018,208	49248
Troffee	H=192	108,000	9,504	99,520	842,672	822,804,736	98496
Traffic	H=336	181,872	9,504	173,392	916,544	823,984,528	172368
	H=720	378,864	9,504	369,904	1,113,536	827,130,640	369360
Relative S	Size (Avg)	1.0	0.039	0.953	5.040	4530.574	0.948

Table 8: Cor	nparison o	of the	trainable	model	parameters
--------------	------------	--------	-----------	-------	------------

of PSFormer(Encoder) is much smaller, indicating that optimizing the linear mapping layer in the output could further reduce the overall parameter count.

## A.5 RUNNING TIME COMPARISON

The average running time per iteration (s/iter) of different models on the ETTh1 and Weather datasets with varying prediction horizons is shown in Table 9. PSformer demonstrates relatively stable running times across different horizons. For the ETTh1 dataset, the running time remains between 0.011 and 0.012 seconds, while for the Weather dataset, it varies slightly between 0.026 and 0.027 seconds. PSformer also shows comparatively efficient running times across the datasets, with performance that remains competitive even as the prediction horizon increases. This indicates that PS former manages computational costs effectively, especially when compared to more complex models. 

Table 9: Comparison of the running time (s/iter). We test the average running time per iteration ofdifferent models across the first five epochs on the ETTh1 and Weather datasets.

Dataset	Horizon	PSformer	PatchTST	ModernTCN	TSMixer	RLinear	iTransformer
	H=96	0.012	0.049	0.241	0.013	0.032	0.020
	H=192	0.011	0.049	0.244	0.016	0.033	0.023
EIINI	H=336	0.012	0.054	0.245	0.015	0.034	0.022
	H=720	0.012	0.063	0.279	0.019	0.037	0.025
	H=96	0.026	0.165	0.388	0.013	0.030	0.022
Waathar	H=192	0.026	0.166	0.361	0.016	0.032	0.022
weather	H=336	0.027	0.176	0.730	0.016	0.033	0.027
	H=720	0.027	0.185	0.788	0.024	0.034	0.035

# 

## A.6 DISCUSS ABOUT POSITIONAL ENCODING

805 The reasons of eliminating positional encoding Similarly, SAMformer does not use positional 806 encoding because it applies attention to the channel dimension, where there is no strict sequential 807 relationship between channels. Although SegAtt is also a cross-channel structure, it involves local 808 time series constructed by patches. We consider such local sequences as local representations (or 809 tokens) rather than short sequences. The experimental results also demonstrate that this structure 808 can similarly reduce the dependency on positional encoding while achieving good performance.

810A.7SHARPNESS-AWARE MINIMIZATION (SAM)811

Optimization steps SAM optimizer (Foret et al., 2021) modifies the typical gradient descent update
 to seek a flatten optimum. Below is the mathematical formulation:

Let  $\theta$  be the model parameters,  $\mathcal{L}(\theta)$  be the loss function, and  $\epsilon$  be a small perturbation applied to the parameters. The SAM optimization process can be described in two steps:

• Find the adversarial perturbation that maximizes the loss in a neighborhood of the current weights  $\theta$ :

 $\hat{\epsilon} = \arg \max_{\|\epsilon\| \le \rho} \mathcal{L}(\theta + \epsilon)$ 

where  $\rho$  is a small constant that controls the size of the neighborhood.

- Update the parameters in the direction that minimizes the loss with respect to the perturbed parameters:
  - $\theta \leftarrow \theta \eta \nabla_{\theta} \mathcal{L}(\theta + \hat{\epsilon})$
- where  $\eta$  is the learning rate.

As used in SAMformer, we also employ this optimization technique to train our models, which can generate promising results.

**B** MORE RESULTS AND ANALYSIS

**B.1** INVESTIGATION OF HYPER-PARAMETER  $\rho$ 

833 The Effect of  $\rho$  Since we employed SAM to ensure training stability, we also conducted sensitivity 834 tests on the hyper-parameter  $\rho$  in SAM. We divided the range of  $\rho$  from 0 to 1 into 10 equal parts 835 and tested its effect on model prediction performance across the ETTh1, ETTm1, ETTm2, and 836 Weather datasets. The results are shown in Figure 3. It can be observed that as the parameter 837  $\rho$  gradually increases in SAM can improve the model's prediction performance to some extent. 838 However, if  $\rho$  is set too large, it may degrade the model's performance. When selecting  $\rho$ , it's 839 important to consider the dataset's complexity and noise levels, as well as the model's architecture. 840 For more complex datasets or larger models, a slightly larger  $\rho$  can help smooth the loss landscape 841 and improve generalization. Further,  $\rho$  should also be balanced with the learning rate to avoid 842 instability or performance degradation. As a comparison, we also report the  $\rho^*$  used by *PSformer*, 843 SAMformer and TSMixer in the Table 10.

Table 10: Neighborhood size  $\rho^*$  used by *PSformer*, *SAMformer* and *TSMixer* for SAM optimization to achieve their best performance on the benchmarks.

Н	Model	ETTh1	ETTh2	ETTm1	ETTm2	Electricity	Exchange	Traffic	Weather
96	PSformer	0.6	0.1	0.4	0.0	0.0	0.2	0.1	0.1
	SAMformer	0.5	0.5	0.6	0.2	0.5	0.7	0.8	0.4
	TSMixer	1.0	0.9	1.0	1.0	1.0	0.9	0.0	0.5
192	PSformer	0.8	0.0	0.4	0.2	0.1	0.1	0.1	0.1
	SAMformer	0.6	0.8	0.9	0.9	0.6	0.8	0.1	0.4
	TSMixer	0.7	0.1	0.6	0.9	0.6	0.9	0.9	0.4
336	PSformer	0.9	0.6	0.4	0.3	0.1	0.2	0.2	0.2
	SAMformer	0.9	0.6	0.9	0.8	0.5	0.5	0.5	0.6
	TSMixer	0.7	0.7	0.9	1.0	0.4	1.0	0.6	0.6
720	PSformer	0.6	0.5	0.4	0.3	0.1	0.2	0.3	0.3
	SAMformer	0.9	0.8	0.9	0.9	1.0	0.9	0.7	0.5
	TSMixer	0.3	0.4	0.5	1.0	0.9	0.1	0.9	0.3

859 860

817

818

819 820

821

822

823

824 825

826

827

828 829 830

831

844

845

861

862



Figure 3: Ablation analysis on hyper-parameter  $\rho$ . When taking  $\rho$  values from 0 to 1 in steps of 0.1, the prediction loss will slightly decrease first and then increase significantly if the  $\rho$  exceeds a threshold, which means the selection of  $\rho$  should be careful.

## B.2 FULL RESULTS

Table 11 presents the detailed experimental results of different models and forecasting horizons, providing a comprehensive evaluation of their performance in long-term time series forecasting tasks. The performance is measured using Mean Squared Error (MSE) and Mean Absolute Error (MAE). In the table, red values represent the best performance in the respective task and metric, while blue-lined values indicate the second-best performance.

The last row of the table summarizes the number of first-place results for each model across all tasks. As can be seen, PSformer achieved the best MSE performance in 20 out of 32 prediction tasks, and ranked second in 8 tasks. In terms of MAE, PSformer achieved the best results in 22 tasks and came second in 5 tasks. This clearly demonstrates the superior performance of PSformer compared to other baseline models in long-term time series forecasting tasks.

The next best-performing model is ModernTCN, which achieved the best MSE results in 6 tasks and the best MAE results in 3 tasks. While other models such as SAMformer and PatchTST also showed competitive performance in some tasks, their overall results are not as strong as those of PSformer and ModernTCN. In summary, PSformer's strong performance across multiple benchmark tasks suggests its potential effectiveness in long-term forecasting.

910

888

889

890 891 892

893

899

## B.3 TRAINING LOSS

911 912

Figure 4 illustrates the training and validation loss curves of the ETTh1 and ETTm1 datasets with prediction horizons  $H = \{96, 192\}$ . In this experiment, we set the number of epochs to 200 and disabled early stopping to observe the complete training process. As shown in the plots, despite the model reaching convergence early in the training (as indicated by the flattening of the training loss curve), the validation loss remains consistently low throughout the training duration. This indicates the model's stability and its ability to generalize well to unseen data over extended epochs.

Table 11: Full long-term forecasting results. We set the forecasting horizons  $H \in \{96, 192, 336, 720\}$ . A lower value indicates better performance. Avg means the average results from all prediction lengths. Red: the best, Blue lined: the second best.

932	Model	5	PSfo	rmer	SAM	ormer	TSM	lixer	Patc	hTST	MON	AENT	Moder	nTCN	FEDf	ormer	GPT	'4TS	Autof	ormer	RLi	near	iTrans	former
033	Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
004		96	0.352	0.385	0.381	0.402	0.388	0.408	0.414	0.419	0.387	0.410	<u>0.373</u>	0.399	0.376	0.415	0.376	0.397	0.435	0.446	0.386	<u>0.395</u>	0.386	0.405
934	ETTh1	192	0.385	0.406	0.409	0.418	0.421	0.426	0.460	0.445	0.410	0.426	0.436	0.419	0.423	0.446	0.416	0.418 0.433	0.456	0.457	0.437	0.424	0.441	0.436
935		720	0.411	0.456	0.425	0.449	0.440	0.459	0.500	0.488	0.454	0.472	0.467	0.474	0.469	0.492	0.477	0.515	0.515	0.517	0.481	0.470	0.503	0.491
936		Avg	0.397	0.418	<u>0.410</u>	<u>0.424</u>	0.420	0.431	0.468	0.455	0.418	0.436	0.421	0.432	0.428	0.454	0.428	0.426	0.473	0.477	0.446	0.434	0.454	0.448
937		96	<u>0.272</u>	0.337	0.295	0.358	0.305	0.367	0.302	0.348	0.288	0.345	0.271	<u>0.339</u>	0.332	0.374	0.285	0.342	0.332	0.368	0.288	0.338	0.297	0.349
029	ETTb2	192	0.335	0.379	0.340	0.386	0.350	0.393	0.388	0.400	0.349	0.386	0.332	0.382	0.407	0.446	0.354	0.389	0.426	0.434	0.374	0.390	0.380	0.400
930	21112	720	0.389	0.431	0.391	0.393	0.402	0.435	0.420	0.446	0.309	0.439	0.402	0.441	0.412	0.469	0.406	0.441	0.453	0.490	0.420	0.440	0.428	0.445
939		Avg	0.338	0.390	0.344	<u>0.391</u>	0.354	0.400	0.387	0.407	0.352	0.394	<u>0.343</u>	0.393	0.388	0.434	0.355	0.395	0.422	0.443	0.374	0.398	0.383	0.407
940		96	0.282	0.336	0.329	0.363	0.327	0.363	0.329	0.367	0.293	0.349	0.310	0.356	0.326	0.390	<u>0.292</u>	<u>0.346</u>	0.510	0.492	0.355	0.376	0.334	0.368
941		192	0.321	0.360	0.353	0.378	0.356	0.381	0.367	0.385	<u>0.326</u>	<u>0.368</u>	0.340	0.373	0.365	0.415	0.332	0.372	0.514	0.495	0.391	0.392	0.377	0.391
942	ETIMI	336 720	0.352	0.380	0.382	0.394	0.387	0.397	0.399	0.410	0.352	0.384	0.373	0.392	0.392	0.425	0.417	0.394	0.510	0.492	0.424	0.415	0.426	0.420
040		Avg	0.342	0.372	0.373	0.388	0.378	0.392	0.387	0.400	0.344	0.379	0.361	0.384	0.382	0.422	0.351	0.383	0.515	0.493	0.414	0.407	0.407	0.410
943		96	0.167	0.258	0.181	0.274	0.190	0.284	0.175	0.259	0.170	0.260	<u>0.168</u>	0.261	0.180	0.271	0.173	0.262	0.205	0.293	0.182	0.265	0.180	0.264
944		192	0.219	0.292	0.233	0.306	0.250	0.320	0.241	0.302	<u>0.227</u>	<u>0.297</u>	0.231	0.305	0.252	0.318	0.229	0.301	0.278	0.336	0.246	0.304	0.250	0.309
945	ETTm2	336	0.269	0.325	0.285	0.338	0.301	0.350	0.305	0.343	0.275	<u>0.328</u>	<u>0.272</u>	<u>0.328</u>	0.324	0.364	0.286	0.341	0.343	0.379	0.307	0.342	0.311	0.348
946		Avg	0.347	0.376	0.375	0.390	0.389	0.402	0.402	0.400	0.363	0.318	0.375	0.394	0.410	0.420	0.378	0.401	0.414	0.419	0.407	0.398	0.412	0.407
047		96	0.149	0.200	0.197	0.249	0.189	0.242	0.177	0.218	0.154	0.209	0.154	0.209	0.238	0.314	0.162	0.212	0.249	0.329	0.192	0.232	0.174	0.214
947		192	0.193	0.243	0.235	0.277	0.228	0.272	0.225	0.259	0.197	0.248	0.207	0.257	0.275	0.329	0.204	0.248	0.325	0.370	0.240	0.271	0.221	0.254
948	Weather	336	0.245	0.282	0.276	0.304	0.271	0.299	0.278	0.297	<u>0.246</u>	<u>0.285</u>	0.248	0.289	0.339	0.377	0.254	0.286	0.351	0.391	0.292	0.307	0.278	0.296
949		720	0.314	0.332	0.334	0.342	0.331	0.341	0.354	0.348	0.315	0.336	0.337	0.342	0.389	0.409	0.326	0.337	0.415	0.426	0.364	0.353	0.358	0.347
950		Avg	0.122	0.204	0.201	0.253	0.255	0.289	0.239	0.281	0.126	0.222	0.237	0.274	0.310	0.357	0.237	0.271	0.335	0.379	0.272	0.291	0.149	0.278
051		192	0.133 0.149	0.225	0.155	0.252	0.191	0.273	0.181	0.270	0.150	0.233	0.133	0.228	0.180	0.302	0.153	0.258	0.190	0.313	0.201	0.281	0.148	0.240
551	Electricity	336	<u>0.164</u>	0.258	0.183	0.277	0.198	0.297	0.204	0.293	0.167	0.264	0.163	<u>0.260</u>	0.213	0.328	0.169	0.266	0.214	0.327	0.215	0.298	0.178	0.269
952		720	<u>0.203</u>	<u>0.291</u>	0.219	0.306	0.230	0.321	0.246	0.324	0.205	0.295	0.194	0.289	0.233	0.344	0.206	0.297	0.236	0.342	0.257	0.331	0.225	0.317
953		Avg	0.162	0.255	0.181	0.275	0.198	0.296	0.205	0.290	0.165	0.260	0.160	0.255	0.207	0.321	0.167	0.263	0.214	0.327	0.219	0.298	0.178	0.270
954		96 192	0.081	0.197	0.161	0.306	0.233	0.363	0.088	0.205	0.098	0.224	0.207	0.342	0.139	0.276	0.091	0.212	0.197	0.323	0.093	0.217	0.177	0.206
955	Exchange	336	0.328	0.412	0.368	0.453	0.474	0.515	0.301	0.397	0.387	0.454	0.520	0.553	0.426	0.464	0.328	0.417	0.509	0.524	0.351	0.432	0.331	0.417
056		720	0.842	0.689	1.003	0.750	1.078	0.777	0.901	0.714	1.062	0.783	1.154	0.810	1.090	0.800	0.880	0.704	1.447	0.941	0.886	0.714	<u>0.847</u>	<u>0.691</u>
900		Avg	0.358	0.399	0.445	0.470	0.532	0.523	0.367	0.404	0.437	0.446	0.555	0.536	0.478	0.477	0.371	0.409	0.613	0.539	0.378	0.417	<u>0.360</u>	<u>0.403</u>
957		96	0.367	0.257	0.407	0.292	0.409	0.300	0.462	0.295	0.391	0.282	0.391	0.271	0.576	0.359	0.388	0.282	0.597	0.371	0.649	0.389	0.395	0.268
958	Traffic	336	0.390	0.272	0.415	0.294	0.433	0.299	0.400	0.304	0.404	0.287	0.410	0.275	0.608	0.330	0.407	0.290	0.623	0.382	0.609	0.369	0.433	0.270
959		720	0.439	0.294	0.456	0.311	0.488	0.344	0.514	0.322	<u>0.450</u>	0.310	0.451	0.305	0.621	0.375	0.450	0.312	0.639	0.395	0.647	0.387	0.467	<u>0.302</u>
060		Avg	0.400	0.274	0.425	0.297	0.439	0.315	0.481	0.304	0.415	0.293	<u>0.414</u>	0.283	0.604	0.372	<u>0.414</u>	0.295	0.617	0.384	0.626	0.378	0.428	<u>0.282</u>
500	1 <sup>st</sup> Cou	nt	20	22	2	3	0	0	2	3	2	0	6	3	0	0	0	0	0	0	0	0	2	2



Figure 4: Training and validation loss curves of the ETTh1 and ETTm1 datasets.

## 1002 B.4 VISUALIZATION

To better understand our method, we present the forecast curves and the corresponding attention maps for several selected samples. The attention maps visualize the attention weights in different stages of the model, providing insight into how the model processes the input data at each stage.

1007 For the ETTh1 in Figure 5, the input sequence length is 512, and we display the last 100 time steps 1008 of the input. The prediction length is set to 96. We select five samples from the ETTh1 dataset, and 1009 for each sample, we visualize the attention maps for stage 1 and stage 2 of the SegAtt. From the 1010 attention maps, it is evident that there are significant variations in the attention distributions across 1011 different samples. Additionally, the attention maps from stage 1 and stage 2 also show noticeable 1012 differences, despite both stages sharing the same PS block parameters. This indicates that while the 1013 two-stage share parameters, they are able to handle and process the information differently, capturing different aspects of the input data at each stage of the model. 1014

For the Weather in Figure 6, the input length is also 512, and the last 100 time steps are displayed, while the prediction length is 192. Since the model for this dataset employs a three-layer Encoder structure, we display the attention maps for both stages across each layer. Specifically, the notation"1-2" represents the attention map for layer 1, stage 2, and similarly for the other layers. The first two rows of attention maps correspond to the attention distributions from the three Encoder layers. Following that, the prediction curves for nine selected variates are plotted, providing a detailed view of the model's forecast performance across different variates.

1022

1024

## 1023 B.5 Additional Visualization of Attention Map

Here, We further analyze the SegAtt attention matrix of PS former in both the temporal and spatial dimensions. Using samples from ETTh1, the key parameters include: input length L = 512, segment







# 1134 B.5.1 Additional Visualization of the SegAtt Attention

Since stage one of SegAtt operates on the input data x, we print out the attention matrix to observe how SegAtt captures local spatio-temporal information. The attention matrix is shown in the Figure 7 below. To better display the colors, we choose the plasma color style and apply a clipping operation to the top 1.



Figure 7: Comparison of SegAtt map and cross-channel map.

1161 1162

1140

1163 Since both the x-axis and y-axis correspond to a specific point within a segment of length C, the 1164 brighter yellow points in this attention matrix indicate higher attention weights between the cor-1165 responding x-axis and y-axis positions. Taking this attention map as an example, we can observe 1166 distinct high-attention regions (e.g., the top-left corner) and high-attention areas between different 1167 time steps and variables (e.g., the non-diagonal symmetric grid section in the top-left corner). From this, we can observe that the coordinate positions within the spatial attention submatrix (inside the 1168 grids) exhibit the same or gradually changing attention weights across different time steps (between 1169 the grids). This suggests that the model may capture temporal consistency between variables, local 1170 stationary features, and long-term dependencies across time steps. 1171

1172

# 1173 B.5.2 SINGLE-CHANNEL ATTENTION MAP

For better visualization of SegAtt's capture of the local temporal dimension, we extract the singlechannel attention submatrix from the attention matrix, as shown in Figure 8a, and visualize the corresponding single-channel local time sequence in Figure 8b. This allows us to observe the highattention weights at specific temporal local positions. We present three samples of attention along the temporal dimension for a single sequence.

From Figure 8a, the high attention weight at coordinates (x=14, y=7) corresponds to time steps 7 and 14 in the Figure 8b. These two time points may represent a pattern of significant change in the sequence, which the model identifies as a noteworthy relationship.

- Using the same analysis method, we can observe in Figure 9 that high attention occurs between time steps 2 and 3, as well as between time steps 9 and 10. This includes both consecutive time points with related attention (e.g. 2 and 3) and attention with intervals (e.g. 2-3 and 9-10).
- 1187 At the same time, as shown in Figure 10, there are some more complex temporal attention patterns, which are difficult to intuitively understand.



#### 1242 **B**.5.3 **CROSS-CHANNELS ATTENTION MAP** 1243

1244 To visualize the cross-channel attention relationships in SegAtt, we extract the attention weights between two channels, forming the corresponding cross-channel attention submatrix map and the 1245 time series plots for the two variables, as shown in Figure 11. 1246



Figure 11: Comparison between Variate 0 and Variate 5

1264 Besides, The x-axis positions in Figure 11a correspond to the time series of variate 5 (yellow line) in 1265 Figure 11b, and the y-axis positions correspond to the time series of variate 0 (blue line). Therefore, 1266 after establishing the correspondence, we can observe that time steps 1-2 of variate 5 and time steps 4-5, 11-12 of variate 0 have higher attention weights.



1282

1262

1263

1267

1283

1294

1284

Figure 12: Comparison between Variate 0 and Variate 2

1285 From Figure 12, we can see that high attention is mainly concentrated at coordinates (3, 14) and (5, 1286 14), which correspond to time step 14 of variate 2 and time steps 3 and 5 of variate 0. These three 1287 positions may reflect a reversal relationship for the two time series.

Additionally, in Figure 13, we can also observe high attention weights between the two channels 1289 at the same time step. This occurs when both sequences undergo significant changes in opposite 1290 directions simultaneously. 1291

#### **B.5.4** ADDITIONAL VISUALIZATION OF THE SEGATT ATTENTION MAP WITHOUT 1293 PARAMETER SHARING

We further analyzed how parameter sharing affects the model's ability to capture temporal pat-1295 terns by comparing attention maps with and without parameter sharing. To achieve this, we also



the interpretability in temporal models is challenging and remains an open research question worthy of further exploration.



Variate	Model	96	192	336	720
1	PSformer	0.057	0.147	0.360	0.803
-	PatchTST	0.063	0.152	0.461	-
2	PSformer	0.042	0.094	0.152	0.216
2	PatchTST	0.052	0.140	0.181	-
3	PSformer	0.034	0.093	0.169	0.476
3	PatchTST	0.055	0.116	0.176	-
4	PSformer	0.041	0.066	0.090	0.146
4	PatchTST	0.058	0.085	0.111	-
5	PSformer	0.007	0.009	0.012	0.073
3	PatchTST	0.016	0.025	0.022	-
6	PSformer	0.077	0.169	0.532	1.125
U	PatchTST	0.099	0.215	0.505	-
7	PSformer	0.033	0.069	0.120	0.342
1	PatchTST	0.039	0.079	0.140	-
ОТ	PSformer	0.047	0.101	0.190	0.510
01	PatchTST	0.097	0.154	0.230	-

Table 12: Performance comparison between PSformer and PatchTST across various variates.

Table 13: Performance comparison with different norm window sizes on the Exchange Rate dataset.

Dataset	Horizon	16	64	128	256	512
	96	0.081	0.085	0.090	0.092	0.091
	192	0.179	0.187	0.189	0.191	0.197
Exchange Rate	336	0.328	0.338	0.356	0.362	0.345
	720	0.842	0.900	0.976	1.003	1.036
	Avg	0.358	0.378	0.403	0.412	0.417

Table 14: Performance comparison with different positional embedding methods.

Dataset	Horizon	Pos emb (time series)	Pos emb (segment)	No pos emb
<b>Б</b> ТТЬ1	96	0.378	0.353	0.352
EIIII	192	0.412	0.388	0.385
ETTL3	96	0.298	0.276	0.272
E11112	192	0.352	0.338	0.335
	96	0.189	0.095	0.091
Evolongo Doto	192	0.314	0.201	0.197
Exchange Kate	336	0.525	0.370	0.345
	720	1.574	1.041	1.036

1459				-				
1/60		Models	PSformer (F, noDL)	TimeMixer (F, noDL)	CrossGNN (F, noDL)	MICN (F, DL)	TimesNet (F, DL)	FITS (M, noDL)
1400	ETTh1	96	0.352	0.375	0.382	\	0.384	0.372
1461		192	0.385	0.429	0.427	N	0.436	0.404
1462		336	0.411	0.484	0.465	N N	0.491	0.427
1402		720 Avg	0.44	0.498	0.472		0.521	0.424
1463	ETTLO		0.377	0.390	0.300		0.450	0.271
1464	ETINZ	192	0.335	0.372	0.39		0.34	0.331
1/65		336	0.356	0.386	0.426	Ň	0.452	0.354
1405		720	0.389	0.412	0.445	Ň	0.462	0.377
1466		Avg	0.338	0.365	0.393	\	0.414	0.333
1467	ETTm1	96	0.282	0.32	0.335	\	0.338	0.303
1.100		192	0.321	0.361	0.372	)	0.372	0.337
1468		336	0.352	0.39	0.403		0.41	0.366
1469		Avg	0.342	0.381	0.393	\`	0.4	0.355
1470	ETTm2	96	0.167	0.175	0.176	0.179	0.187	0.162
		192	0.219	0.237	0.24	0.307	0.249	0.216
1471		336	0.269	0.298	0.304	0.325	0.321	0.268
1472		720	0.347	0.391	0.406	0.502	0.408	0.348
		Avg	0.251	0.275	0.282	0.328	0.291	0.249
1473	Weather	96	0.149	0.163	0.159	\	0.172	0.143
1474		192	0.193	0.208	0.211	N N	0.219	0.186
1/75		336	0.245	0.251	0.267		0.28	0.236
1470		120	0.314	0.339	0.352		0.365	0.307
1476	Flootnioity	04	0.122	0.152	0.172	0.164	0.169	0.124
1477	Electricity	192	0.133	0.155	0.175	0.177	0.184	0.134
		336	0.164	0.185	0.206	0.193	0.198	0.165
1478		720	0.203	0.225	0.231	0.212	0.22	0.203
1479		Avg	0.162	0.182	0.201	0.187	0.192	0.163
1/180	Exchange rate	96	0.091	0.09	0.084	0.102	0.107	\ \
1400		192	0.197	0.187	0.171	0.172	0.226	\
1481		336	0.345	0.353	0.319	0.272	0.367	\
1482		Avg	0.417	0.391	0.345	0.315	0.416	Ì
1/02	Traffic	96	0.367	0.462	0.57	0.519	0.593	0.385
1403		192	0.39	0.473	0.577	0.537	0.617	0.397
1484		336	0.404	0.498	0.588	0.534	0.629	0.41
1485		720 Avg	0.439	0.506	0.597	0.577	0.64	0.448

Table 15: Comparasion with additional baselines (Part-I).

1458

To highlight the differences between seasonal vs. non-seasonal and stationary vs. non-stationary characteristics, we selected the ETTh1 and ETTh2 datasets (relatively seasonal and stable) as well as the Exchange dataset (relatively non-seasonal and non-stationary). The experimental results are shown in the Table 14.

The degraded performance of pos emb (time series) might be due to the incompatibility of the positional encoding with dimension transformation, as the original temporal order is lost in the segment dimension, making it unsuitable for dot-product attention calculations. On the other hand, pos emb (segment) shows smaller changes compared to the No pos emb case, but the performance still deteriorates slightly. This suggests that the significance of positional encoding in the context of multivariate time series forecasting might need to be re-evaluated, as there are fundamental differences between NLP and time series data when applying attention mechanisms.

1499 1500

1502

1504

1501 B.9 COMPARASION WITH ADDITIONAL BASELINES (PART-I).

1503 B.9.1 COMPARISON WITH EXTENDED BASELINES

We have collected the Extended experimental MSE loss results of the relevant models in Table 15. Although there are differences in experimental setups across each work, which may affect the results and prevent a completely fair comparison, we have provided some key settings to help better understand the model performance.

For each model, the name is followed by two values in parentheses: the first value represents the look-back window mode (with "F" indicating fixed input length and "M" indicating a grid search was performed across different input lengths), and the second value indicates whether the test set dataloader drops the last batch of data (noDL: does not drop last; DL: drops last).

Table 16: Comparison with additional baselines (Part-II).

	Models	PSformer (512, noDL)	Crossformer (96,noDL)	PDF (512, noDL)	PatchTST (512, noDL)	TimeLLM (512, DL)
ETTh1	96	0.352	0.384	0.361	0.374	0.362
	192	0.385	0.438	0.391	0.413	0.398
	336	0.411	0.495	0.415	0.434	0.43
	720	0.44	0.522	0.468	0.455	0.442
	Avg	0.397	0.46	0.409	0.419	0.408
ETTh2	96	0.272	0.347	0.272	0.274	0.268
	192	0.335	0.419	0.334	0.341	0.329
	336	0.356	0.449	0.357	0.364	0.368
	720	0.389	0.479	0.397	0.39	0.372
	Avg	0.338	0.424	0.34	0.342	0.334
ETTm1	96	0.282	0.349	0.284	0.29	0.272
	192	0.321	0.405	0.327	0.333	0.31
	336	0.352	0.432	0.351	0.37	0.352
	720	0.413	0.487	0.409	0.416	0.383
	Avg	0.342	0.418	0.343	0.352	0.329
ETTm2	96	0.167	0.208	0.162	0.166	0.161
	192	0.219	0.263	0.224	0.223	0.219
	336	0.269	0.337	0.277	0.273	0.271
	720	0.347	0.429	0.354	0.363	0.352
	Avg	0.251	0.309	0.254	0.256	0.251
Weather	96	0.149	0.191	0.147	0.152	0.147
	192	0.193	0.219	0.191	0.196	0.189
	336	0.245	0.287	0.243	0.247	0.262
	720	0.314	0.368	0.317	0.315	0.304
	Avg	0.225	0.266	0.225	0.228	0.226

From the results, compared to models with fixed windows, PS former performed best on 7/8 of the prediction tasks. This further highlights PS former's competitive performance in forecasting. Even when compared to non-fixed window models like FITS, PSformer performed best on 4/7 of the prediction tasks. 

## **B.9.2** COMPARISON WITH ADDITIONAL BASELINES (PART-II)

We have collected the Extended experimental results of the relevant models in Table 16. Compari-son with PDF and PatchTST. While PDF demonstrates strong predictive performance, as shown in its paper, it also relies on a more complex hyperparameter search process. Additionally, PDF ad PatchTST set  $drop_last = True$  in the test set dataloader, which results in the model inadvertently using incomplete batches during evaluation. To ensure a fair comparison, we modified PDF and PatchTST by correcting this DL issue (setting  $drop\_last = True$ ) and set the input length to 512 to align with our experimental setup. Under these adjustments, we have included the updated results in the table below, which provide further insights into PS former performance under fair conditions. We conducted tests on the following five datasets. From the experimental results, both PS former and PDF significantly outperform PatchTST in predictive performance, and PSformer achieving better predictions than PDF. However, the average prediction loss reduction relative to PDF is not substan-tial. Therefore, we consider PSformer and PDF to exhibit equally excellent predictive performance under the same settings. 

Comparison with Time-LLM. For Time-LLM, We listed the predictive performance from the origi-nal paper of the model in the table below. However, we check its offical repository and find Time-LLM also faces DL issue in the test set dataloader, which may affects its reported results. Besides, due to the computational demands of large-scale models, we decided not to execute its code directly in our experiments. When selecting the baseline large models, we considered both MOMENT and TimeLLM. We ultimately chose MOMENT for two main reasons: the MOMENT paper includes a direct comparison with TimeLLM, and it is relatively lightweight (428M parameters). In sum-mary, despite the significant difference in parameter scale and the DL issue present in TimeLLM, PSformer still achieves equal or better average predictive performance than TimeLLM on 3 out of 5 datasets. 

Comparison with Crossformer. For Crossformer, we used experimental results based on the Cross-GNN[7] paper, as there is an inconsistency in prediction lengths between the original Crossformer paper and current mainstream work. Therefore, we aligned the experimental setup with the Cross-GNN results to ensure consistency in the comparison. From the results, PSformer performs better than Crossformer.

	methods	96	192	336	720
ETTh1	w parameter sharing	83/0.352	70/0.385	53/0.411	35/0.440
	w/o parameter sharing	46/0.359	108/0.392	39/0.423	36/0.441
Exchange	w parameter sharing	71/0.081	51/0.179	82/0.328	32/0.842
	w/o parameter sharing	47/0.084	32/0.183	43/0.333	31/0.855

Table 17: Comparisons of convergence rates.

Table 18: Comparisons of Parameter-Saving Capacity.

	1	3	12	24	36	48
Parameter Sharing	52,416	58,752	87,264	125,280	163,296	201,312
No Parameter Sharing	71,424	115,776	315,360	581,472	847,584	1,113,696

## 1584 B.10 COMPARISONS OF CONVERGENCE RATES WITH AND WITHOUT PARAMETER SHARING.

We compare the impact of parameter sharing on the convergence rate using the ETTh1 and Exchange datasets, recording the total number of epochs and the MSE loss under the same settings.

The experimental results are shown in the Table 17, where the values represent epochs/MSE loss. We observe that the number of epochs required with or without parameter sharing on the ETTh1 dataset varies depending on the prediction length. However, for the Exchange dataset, the convergence rate is faster without parameter sharing.

Additionally, in terms of MSE loss, using parameter sharing leads to greater reductions in loss and also results in fewer parameters. Therefore, there exists a trade-off between convergence rates, loss reduction, and parameter efficiency.

## 1596 B.11 FURTHER TESTING OF PARAMETER-SAVING CAPACITY FOR PRE-TRAINED MODELS

For further validating the framework's parameter-saving capacity, we compared the parameter count of the PSformer under Parameter Sharing and No Parameter Sharing scenarios, including the comparison for 1-layer and 3-layer Encoders, as well as for different layers same as GPT2 models (GPT2-small, GPT2-medium, GPT2-large, and GPT2-xl) at 12-layer, 24-layer, 36-layer, and 48-layer configurations. The results are reported in the Table 18. In addition, if the hidden layer dimension is expanded from 32 to 1024 (as in GPT2), or if multi-head attention is adopted, the total number of parameters will also increase significantly.

1604

1595

1566

1579 1580 1581

## 5 **B.12** More details about parameter search space.

The main hyperparameters of PS former include: 1. the number of encoders, 2. the number of segments, and 3. the SAM hyperparameter rho. For the number of encoders, we primarily searched 1608 within 1 to 3 layers. For the number of segments, we maintain the same as the number of patches 1609 in PatchTST, we also analyzed values that divide the input length evenly (specifically: 2, 4, 8, 16, 1610 32, 64, 128, 256), and ultimately set all prediction tasks to 32 to avoid performance improvements 1611 caused by complex hyperparameter tuning. For the SAM hyperparameter rho, we referred to SAM-1612 former and performed a search across 11 parameter points evenly spaced in the range 0 to 1. The 1613 number of encoders and segments can be found in the ablation study in section 4.2. Additionally, for 1614 the learning rate, we mainly tested values of 1e - 3 and 1e - 4, and for the learning rate scheduler, 1615 we tested OneCycle and MultiStepLR.

1616

1618

## 1617 B.13 MORE RESULTS FOR SEGMENT NUMBER SELECTION.

1619 We further conducted tests on the selection of segmentation numbers using the ETTh2 and Exchange Rate datasets. The results are shown in the Appendix B.13. For prediction lengths of 192 and 336, a

	segment number	8	16	32	64
ETTh1	96	0.273	0.273	0.272	0.275
	192	0.333	0.333	0.335	0.337
	336	0.353	0.357	0.356	0.357
	720	0.387	0.387	0.389	0.394
Exchange	96	0.091	0.092	0.091	0.098
	192	0.183	0.2	0.197	0.235
	336	0.34	0.345	0.345	0.417
	720	1.071	1.071	1.036	1.037
	ETTh1 Exchange	segment number           ETTh1         96           192         336           720         720           Exchange         96           192         336           720         720	segment number         8           ETTh1         96         0.273           192         0.333         336           336         0.353         720           Exchange         96         0.091           192         0.183         336           336         0.34         720	segment number         8         16           ETTh1         96         0.273         0.273           192         0.333         0.333           336         0.353         0.357           720         0.387         0.387           Exchange         96         0.091         0.092           192         0.183         0.2           336         0.34         0.345           720         1.071         1.071	segment number         8         16         32           ETTh1         96         0.273         0.273         0.272           192         0.333         0.333         0.335           336         0.353         0.357         0.356           720         0.387         0.387         0.389           Exchange         96         0.091         0.092         0.091           192         0.183         0.2         0.197           336         0.34         0.345         0.345           720         1.071         1.036

Table 19: Effect of Segment number selection.

segment number of 8 further improves the forecasting performance (compared with default segment number of 32). This suggests that there can not be a perfect fixed value for the choice of the number of segments. The choice of the number of segments may be influenced by the characteristics of different datasets, prediction lengths, as well as the overall increase in the scale of model parameters.

Therefore, for datasets without clear seasonality and with non-stationary characteristics (e.g. financial asset time series, including exchange rates), it is preferable to choose a larger segmentation
number. This allows each segment to capture smaller local spatio-temporal patterns, better addressing unstable variation modes. On the other hand, for datasets with significant seasonality or
relatively stationary characteristics (e.g. electricity and traffic), a relatively smaller segmentation
number often facilitates model training.