TWINSFORMER: REVISITING INHERENT DEPENDEN CIES VIA <u>TWO</u> <u>INTERACTIVE COMPONENTS</u> FOR TIME SERIES FORECASTING

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032 033 034 Paper under double-blind review

ABSTRACT

Due to the remarkable ability to capture long-term dependencies, Transformerbased models have shown great potential in time series forecasting. However, real-world time series usually present intricate temporal patterns, making forecasting still challenging in many practical applications. To better grasp inherent dependencies, in this paper, we propose TwinsFormer, a Transformer-based framework utilizing two interactive components for time series forecasting. Unlike the mainstream paradigms of plain decomposition that train the model with two independent branches, we design an interactive strategy around the attention module and the feed-forward network to strengthen the dependencies via decomposed components. Specifically, we adopt dual streams to facilitate progressive and implicit information interactions for trend and seasonal components. For the seasonal stream, we feed the seasonal component to the attention module and feed-forward network with a subtraction mechanism. Meanwhile, we construct an auxiliary highway (without the attention module) for the trend stream by the supervision of seasonal signals. Finally, we incorporate the dual-stream outputs into a linear layer leading to the ultimate prediction. In this way, we can avoid the model overlooking inherent dependencies between different components for accurate forecasting. Our interactive strategy, albeit simple, can be adapted as a plug-and-play module to existing Transformer-based methods with negligible extra computational overhead. Extensive experiments on various real-world datasets show the superiority of TwinsFormer, which can outperform previous state-of-theart methods in terms of both long-term and short-term forecasting performance.

1 INTRODUCTION

As a ubiquitous and paramount task in many real-world scenarios (e.g., weather (Wu et al., 2023b), energy (Yin et al., 2021), mar-037 ket (Granger & Newbold, 2014), and transportation (Yin et al., 2021)), time series forecasting has been explored with ongoing passion. Generally, time series forecasting aims to predict future 040 temporal variations based on historical observations of time se-041 ries, where the primary challenge is how to effectively capture 042 temporal patterns from observed data. Benefiting from the ad-043 vancements in deep learning, various representative models with 044 well-designed architectures, such as MLP-based (Wang et al., 2024; Zeng et al., 2023; Li et al., 2023), CNN-based (Wang et al., 2023; Wu et al., 2023a; Liu et al., 2022a), and Transformer-046 based (Liu et al., 2024; Zhang & Yan, 2023; Zhou et al., 2022) 047 methods, have been proposed to tackle time series forecasting 048 tasks and demonstrate impressive performance. Since the complex and non-stationary nature of the real world or systems, the





observed series usually involves multitudinous variations, such as increasing, decreasing, and fluctuating, making it still hard to grasp reliable and stable temporal dependencies.

To tackle intricate temporal patterns, series decomposition (Robert et al., 1990), utilizing the moving average kernel to smooth out short-term fluctuations or noise in the time series, has been involved



Figure 2: Trend-seasonal Decomposition. The top two subplots showcase the observed values of Electricity and Traffic from two different channels. The lower four subplots present the decomposed trend and seasonal components of the two datasets. Please zoom in for more details.

in deep models as a basic module. Empowered with various decomposition designs, existing methods (Wu et al., 2021; Zhou et al., 2022; Wang et al., 2023; Zeng et al., 2023) generally utilize two 071 independent branches to highlight seasonal and trend properties separately, then combine the sea-072 sonal and trend representations for the final prediction. As seen in Figure 2, we decompose the 073 raw signals (i.e., observed values) of two different channels on Electricity and Traffic datasets into 074 trend and seasonal components for better understanding. Comparatively, the trend and seasonal 075 components exhibit different but indispensable characteristics for the observed time series. More 076 specifically, the former shows the overall variation while the latter presents the cyclical fluctuation. 077 Since trend-seasonal decomposition is an untrainable linear transformation, the decomposed components obtained by the moving average kernel cannot directly reflect precise and intricate patterns 079 for the observed time series. Taking channel 11 on the Electricity dataset as an example, significant 080 variations and fluctuations of the observed series lag behind those of the trend and seasonal compo-081 nents. Such inconsistencies lead to the learned trend and seasonal representations by independent branches may not satisfy the temporal patterns of the observed series. Therefore, a more rational decomposition design should consider the interactions between decomposed components to precisely 083 unravel inherent dependencies for observed values. 084

085 To fill this gap, we propose TwinsFormer, a Transformer-based framework that explicitly explores inherent dependencies via two interactive components for time series forecasting. First, we de-087 compose the observed time series rather than the time series embeddings into trend and seasonal components better to capture the characteristics of the time series itself. Second, since the trend 088 components reflect the long-term fluctuations of the time series, we only feed the seasonal compo-089 nents to the attention and feed-forward modules with a subtraction mechanism to alleviate redundant 090 coding. Most importantly, we regard the outputs of the attention and feed-forward modules as su-091 pervision information to guide the model to capture the representation of the trend components. With 092 our interactive design, TwinsFormer can successfully aggregate seasonal and trend information to learn inherent dependencies between different components for accurate forecasting. Experimentally, 094 our proposed TwinsFormer achieves state-of-the-art performance on seven real-world forecasting 095 scenarios shown in Figure 1 and effectively provides an interactive learning scheme for time series 096 forecasting. The primary contributions are summarized as follows:

097 098

099

100

102

103

067

068

- We delve into the existing decomposition designs for time series forecasting and figure out that the interaction between different components is not explored: these designs simply learn separate representations for trend and seasonal components and overlook non-linear dependencies or significant noise levels among time series.
- We propose TwinsFormer, a Transformer-based framework (perhaps the first to our best knowledge) that explicitly explores inherent dependencies by learning implicit and progressive interactions between different components for time series forecasting.
- Extensive experimental results on 13 real-world benchmarks show the superiority of our TwinsFormer against previous state-of-the-art methods. Specifically, TwinsFormer ranks in the top 1 among 11 models on 18 out of the 22 average settings including various prediction lengths and metrics over long-term and short-term forecasting tasks.

108 2 RELATED WORK

110 2.1 DECOMPOSITIONS FOR TIME SERIES FORECASTING

112 Due to the capacity of the moving average kernel to smooth out short-term fluctuations or noise in the time series, Autoformer (Wu et al., 2021) initially proposed using the moving average ker-113 nel to decompose complex temporal variations into seasonal and trend components. Since then, 114 trend-seasonal decomposition designs based on the moving average kernel have been frequently 115 introduced in time series forecasting works. For instance, SCINet (Liu et al., 2022a) devises a 116 downsample-convolve-interact architecture to extract dynamic temporal features at multiple resolu-117 tions with two sub-sequences. DLinear (Zeng et al., 2023) utilizes the series decomposition as the 118 pre-processing before linear regression. MICN (Wang et al., 2023) adopts multi-scale branches to 119 model the local and global context by decomposing input series into seasonal and trend terms while 120 TimesNet (Wu et al., 2023a) designs a modular architecture to obtain decomposed components with 121 Fourier Transform. More recently, TimeMixer (Wang et al., 2024) mixes multiscale decomposable 122 components for time series forecasting. Due to the non-linear or non-stationary properties of time 123 series, however, a rudimentary moving averaging kernel may inadequately capture precise trends, which impedes the model from learning inherent dependencies through two independent branches. 124

125 126

2.1.1 TRANSFORMERS FOR TIME SERIES FORECASTING

127 Benefiting from the ability to model long-term temporal patterns, Transformer-based methods (Li 128 et al., 2019; Zhou et al., 2021; Liu et al., 2022b) have shown significant success in time series fore-129 casting. Since the quadratic complexity and redundant coding for the self-attention mechanism, 130 most existing Transformer-based approaches modify the attention module to reduce computational 131 overhead. Representative works include Informer (Zhou et al., 2021) introducing ProbSparse self-132 attention and distilling techniques, Autoformer (Wu et al., 2021) incorporating series decomposition with an auto-correlation mechanism, FEDformer (Zhou et al., 2022) implementing the attention 133 module with a Fourier-enhanced structure, etc. Without modifications to the Transformer, some 134 other attempts pay attention to the inherent processing of time series, such as stationarity (Liu et al., 135 2022c; 2023), patching (Du et al., 2023), channel independence (Nie et al., 2023), and inverting 136 operations (Liu et al., 2024), bringing consistently improved performance for time series forecast-137 ing. Besides, refurbishing Transformer in both aspects mentioned above, Crossformer Zhang & 138 Yan (2023) introduces a two-stage-attention mechanism and dimension-segment-wise embedding 139 strategy to capture cross-time and cross-variate dependencies. 140

Going beyond the designs in previous works, TwinsFormer devises an interactive dual-stream architecture without modifying the native components of Transformer. Moreover, we replace the observed series with trend and seasonal components, where the model can better learn inherent dependencies with their interactions. To the best of our knowledge, TwinsFormer is the first attempt to consider interactions between decomposed components on Transformer for time series forecasting.

3 TWINSFORMER

Preliminary Given the historical observation data $X = \{x_1, x_2, \dots, x_M\} \in \mathbb{R}^{M \times N}$ with Mlength look-back window and N variates, the goal of multivariate time series forecasting is to predict the future time series $Y = \{x_{M+1}, \dots, x_{M+\tau}\} \in \mathbb{R}^{\tau \times N}$ at next τ time steps ($\tau > 1$). Following the idea of decomposition (Robert et al., 1990; Wu et al., 2021), time series can be divided into trend and seasonal parts by the moving average kernel. For length-M input series $X \in \mathbb{R}^{M \times N}$, the decomposition process can be formulated as

$$X_T = AvgPool(Padding(X)),$$

$$X_S = X - X_T,$$
(1)

155 156 157

158

160

146

147

where X_T and X_S denote the trend and seasonal components, respectively.

- 159 3.1 STRUCTURE OVERVIEW
- 161 Our proposed TwinsFormer illustrated on the left of Figure 3 adopts the encoder-only architecture, renovating Transformer to a dual-stream structure with two decomposed components. Before



Figure 3: Overall framework of TwinsFormer.

embedding the time series, we decompose the observed series into trend (T) and seasonal (S) components in the channel dimension. Then, we feed seasonal embeddings E_S to the attention module and feed-forward network (FFN) with a subtraction mechanism, while feeding trend embeddings E_T to the interactive module with the supervision of seasonal information (i.e., A_S and F_s). Finally, we aggregate seasonal and trend representations for the ultimate prediction.

Embedding the decomposed series as tokens For convenience, we denote $X_{m,:}$ as the simultaneously recorded values for all the variates at the *m* time point, while $X_{:,n}$ as the whole time series of each variate indexed by *n*. Based on Equation 1, the trend and seasonal components of the time series can be formulated as $X_T = AvgPool(Padding(X_{:,n}))$ and $X_S = X - X_T$, respectively. Then, we utilize straightforward linear and dropout layers to create trend and seasonal embeddings with global covariates X_{mark} as follows:

171

172

185 186 $E_T = Dropout(Linear(Concat(X_T, X_{mark}))),$ $E_S = Dropout(Linear(Concat(X_S, X_{mark}))).$ (2)

Note that, $Concat(\cdot)$ is used on the dataset containing X_{mark} information and different components (i.e., X_T and X_S) through separate liner layers in our experiments. In this way, we map decomposed series data $X_T, X_S \in \mathbb{R}^{N \times M}$ from the original space into a new space, where $E_T, E_S \in \mathbb{R}^{N \times D}$ and D is the dimension of the linear layer.

Learning interactions with our TwinsBlock Unlike the existing Transformer variants that strug gle to come up with efficient attention mechanisms to learn multivariate correlation, our Twins Former incorporates interactive learning into Transformer block to explore the interactions between
 decomposed components for better inherent dependencies. Thus, a bundle of efficient attention
 mechanisms can be the plugins and our interactive strategy can promote the performance of existing
 Transformer variants on time series forecasting, which are evaluated in Table 4.

197

3.2 DUAL-STREAM DESIGN WITH INTERACTIVE MODULE

Keeping the original modules (i.e., the self-attention and feed-forward network (FFN) modules)
of Transformer unchanged, our key design lies in the computationally efficient interactive module,
which can guide the model to learn and aggregate the effective representations of trend and seasonal
information. We provide the pseudo-code of our framework in Algorithm 1 of the Appendix.

Seasonal Branch Since the seasonal components obtained without the moving average kernel can better highlight the intrinsic characteristics of the time series data, we feed the seasonal embeddings to the attention and FFN modules to effectively capture the dependencies among the multivariates. Following the attention process of iTransformer (Liu et al., 2024), we regard $E_S \in \mathbb{R}^{N \times D}$ as N D-dimension tokens and utilize queries, keys, and values $Q, K, V \in \mathbb{R}^{N \times d_k}$ to obtain the attentionweighted seasonal representations $A_S \in \mathbb{R}^{N \times D}$, where d_k is the projected dimension

$$Q = E_s W_1 + b_1, K = E_s W_2 + b_2, V = E_s W_3 + b_3, \quad W_i \in \mathcal{R}^{d_k \times d_k}, b_i \in \mathcal{R}^{1 \times d_k},$$

$$A_S = Softmax(\frac{QK^T}{\sqrt{d_k}})V.$$
(3)

According to Equation 1, the seasonal components can be regarded as the residual part of the observed time series data. Intuitively, we adopt the idea of residual learning to implement a corrective strategy by subtracting the outputs of the Attention and FFN modules from the corresponding inputs. The learning process can be formulated as follows:

$$H_1 = LayerNorm(E_S - A_S),$$

$$H_2 = H_1 - FFN(H_1).$$
(4)

218 219 220

221

222

224

225

226

227 228

229

230

237 238

248

251

253 254 255

256

257 258

259 260 **Trend Branch** Considering that untrainable moving average kernels lead to unreliable trend patterns, we fuse seasonal information to assist the learning of the trend branch through our interactive module (IM). On the one hand, attention-weighted A_S well reflects the dependencies among multivariate, which can be converted into a coefficient matrix to update the trend embeddings. On the other hand, the signals discarded by the seasonal branch can be regarded as meaningful information to guide the representation of the trend embeddings. Our interactive module is illustrated on the right of Figure 3, we only use simple structures to train and update the trend branch network:

$$E'_{T} = Sigmoid(Conv_{1\times 1}(\sum Multiconv(A_{S}))) * E_{T},$$

$$F_{T} = Concate(A_{S}, F_{S}, E'_{T}),$$
(5)

where we concatenate A_S , F_S , and E'_T in the embedding dimension and the kernel sizes of multiscale convolutions are 1×1 , 3×3 , and 5×5 , respectively.

Gate Mechanism Inspired by the inherent control of cells in RNNs (Zhao et al., 2017; Dey & Salem, 2017), we devise a gate mechanism σ at the end of each block for both streams to autonomously regulate the pace of information transmission. The gate mechanism for both seasonal and trend streams can be formulated as

$$O_S = Sigmoid(Conv_1(H_2)) \cdot Conv_2(H_2),$$

$$O_T = Sigmoid(Conv_3(F_T)) \cdot Conv_4(F_T),$$
(6)

where $Conv_1, Conv_2, Conv_3$ and $Conv_4$ are four 1×1 convolution operations with different parameters. Taking the output of the former TwinsBlock as the input of the latter TwinsBlock, we stack L TwinsBlocks to learn seasonal and trend representations, and then add them together through a linear projection for the ultimate predictive outcomes, i.e., $\{Y = Projection(O_S + O_T)\} \in \mathbb{R}^{\tau \times N}$.

Rationality Analysis Given historical time series data X, we can obtain its trend and seasonal components (i.e., X_t and X_s) by the moving average kernel. For existing time series forecasting methods, we regard the models as $F(\cdot)$, while regarding the independent branches with decomposition designs as $f_t(\cdot)$ and $f_s(\cdot)$, then we can formulate the predictive outputs Y as

$$Y = F(f_t(X_t) + f_s(X_s)), \quad where \ X = X_t + X_s.$$
(7)

Similarly, we definite the attention module, FFN, interactive module, and gate mechanism of Twins-Former as $g(\cdot)$, $h(\cdot)$, $\phi(\cdot)$, and σ respectively. Then, the outcomes of TwinsFormer are:

$$Y = F(\overbrace{\sigma_t(\phi([X_t, g(X_s), h(X_s - g(X_s))]))}^{\text{interactive learning}}) + \overbrace{\sigma_s(X_s - g(X_s) - h(X_s - g(X_s))))}^{\text{residual learning}}), \quad (8)$$

where $[\cdot]$ indicates concatenation operation. By omitting the constraints from various functions on variables and replacing ' $[\cdot]$ ' with '+' operation, our interactive components can be simplified as

$$X'_{s} = X_{s} - X_{1} - X_{2}, \quad X'_{t} = X_{t} + X_{1} + X_{2}.$$
(9)

Then, the sum of our two interactive components can be expressed as

$$X = X'_{s} + X'_{t} = X_{s} - X_{1} - X_{2} + X_{t} + X_{1} + X_{2} = X_{s} + X_{t}.$$
 (10)

261 Based on Equation (10), we can find that our interaction strategy perfectly fits the requirements of 262 the decomposition design without bringing in redundant signals. Furthermore, we can elaborate 263 on the practical implications of our TwinsFormer in mitigating the limitations of the trend-seasonal 264 decomposition. Since the untrainable moving average kernel does not accurately capture the trend 265 patterns, the decomposed trend and seasonal components are not completely disentangled. Twins-266 Former adopts a dual-stream interaction strategy to implicitly and progressively promote the decoupling of both components by using residual learning and interactive learning. Concretely, we filter 267 out the coupled information (i.e., X_1 and X_2) from the seasonal components and compensate them 268 to the trend components by some transformation mechanisms, so that we can learn more robust and 269 reliable decomposed representations for accurate time series forecasting.

270 4 **EXPERIMENTS** 271

272

We conduct extensive experiments to evaluate the performance of TwinsFormer, covering long-273 term and short-term time series forecasting, including 13 real-world benchmarks and 10 well-274 acknowledged baselines. Moreover, we dive into the effectiveness and generality of the proposed 275 framework to existing Transformer-based methods with indispensable ablation studies.

276 **Benchmarks** For long-term forecasting, we experiment on 9 well-established benchmarks: 277 ETT (Zhou et al., 2021) (Electricity Transformer Temperature) datasets including ETTm1, ETTm2, 278 ETTh1, and ETTh2, ECL (Wu et al., 2021) (Electricity Consuming Load), Exchange (Wu et al., 279 2021), Traffic (Wu et al., 2021), Weather (Wu et al., 2021) and Solar-energy (Lai et al., 2018) 280 datasets. For short-term forecasting, we use 4 public traffic network PeMS (Liu et al., 2022a) 281 datasets, namely PEMS03, PEMS04, PEMS07 and PEMS08. We follow standard protocols like 282 (Wu et al., 2021; Liu et al., 2024) and split all datasets into training, validation and test sets in chronological order by the ratio of 6:2:2 for the ETT dataset and 7:1:2 for the other datasets. De-283 tailed dataset descriptions are provided in Table 5 of the Appendix. 284

285 **Baselines** We compare TwinsFormer with 10 representative baselines, including 1) Transformer-286 based methods: iTransformer (Liu et al., 2024), PatchTST(Nie et al., 2023), Crossformer (Zhang 287 & Yan, 2023), FEDformer (Zhou et al., 2022); 2)Linear-based methods: TimeMixer (Wang et al., 288 2024), DLinear (Zeng et al., 2023), TiDE (Das et al., 2023), and RLinear (Li et al., 2023); and 3) TCN-based methods: TimesNet Wu et al. (2023a), SCINet (Liu et al., 2022a). 289

290 **Implementation details** All the experiments are implemented in PyTorch (Paszke et al., 2019) 291 and conducted on one NVIDIA 4090 24GB GPU. We use the L2 loss to train the model with the 292 Adam (Kingma & Ba, 2015) optimizer, where the training process is early stopped within 10 epochs. 293 Our TwinsBlock is applicable to Transformer-based architectures without introducing any additional hyperparameters. Following iTransformer (Liu et al., 2024), we use the Mean Square Error (MSE) 294 and Mean Absolute Error (MAE) as the core metrics for the evaluation. 295

4.1 MAIN RESULTS 297

296

306

298 **Long-term forecasting** As shown in Table 1, TwinsFormer achieves leading performance on most 299 benchmarks, covering various time series with different frequencies, variate numbers and real-world 300 scenarios. For example, TwinsFormer outperforms iTransformer by a considerable margin, with a 301 6.2% MSE reduction in ECL and a 5.1% MSE reduction in Traffic. On Weather and Solar-energy 302 datasets, although TimeMixer has a subtle advantage over TwinsFormer of 0.4% and 4.8% in MSE 303 reduction, Twinsformer achieves lower MAE scores than TimeMixer by 1.1% and 9.3% reduc-304 tion, respectively. It is worth noting that TwinsFormer exhibits better performance than TimeMixer 305 among other datasets, which further highlights the superiority and robustness of TwinsFormer.

Table 1: Long-term forecasting results. The lookback length is set to T = 96 and all the results are 307 averaged from all predictions $S \in \{96, 192, 336, 720\}$. Avg means further averaged by subsets. A 308 lower MSE or MAE indicates a better prediction. See Table 6 in the Appendix for the full results. 309

Models	Twins (O	Former urs)	iTrans (20	former 24)	Time (20	Mixer 24)	Patcl (20	nTST 23)	RLi (20	near 23)	Cross (20	former 23)	TiE (<mark>202</mark>	DE 23)	Time (20	esNet 23a)	DL: (<mark>20</mark>	inear)23)	SC (20	INet 22a)	FEDf (20	'orm)22)
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	М
ETT (Avg)	0.372	0.392	0.383	0.399	0.381	<u>0.396</u>	0.381	0.397	<u>0.380</u>	0.392	0.685	0.578	0.482	0.470	0.391	0.404	0.442	0.444	0.689	0.597	0.408	0.4
ECL	0.167	0.262	0.178	<u>0.270</u>	0.183	0.272	0.205	0.290	0.219	0.298	0.244	0.334	0.251	0.344	0.192	0.295	0.212	0.300	0.268	0.365	0.214	0.3
Exchange	0.346	0.395	0.360	<u>0.403</u>	0.380	0.417	0.367	0.404	0.378	0.417	0.940	0.707	0.370	0.413	0.416	0.443	0.354	0.414	0.750	0.626	0.519	0.4
Traffic	0.406	0.273	0.428	<u>0.282</u>	0.496	0.298	0.481	0.304	0.626	0.378	0.550	0.304	0.760	0.473	0.620	0.336	0.625	0.383	0.804	0.509	0.610	0.3
Weather	0.246	0.271	0.258	0.278	0.245	<u>0.274</u>	0.259	0.281	0.272	0.291	0.259	0.315	0.271	0.320	0.259	0.287	0.265	0.317	0.292	0.363	0.309	0.3
Solar-energy	y <u>0.227</u>	0.254	0.233	<u>0.262</u>	0.216	0.280	0.270	0.307	0.369	0.356	0.641	0.639	0.347	0.417	0.301	0.319	0.330	0.401	0.282	0.375	0.291	0.3
																						_

319 Short-term forecasting TwinsFormer also performs well in short-term forecasting on PeMS 320 datasets. Due to the complex spatiotemporal dependencies among citywide traffic networks in PeMS 321 benchmarks, many advanced models degenerate a lot in this task. For instance, TimeMixer adopts the multiscale mixing architecture to model complex temporal variations, but its performance is not 322 as good as iTransformer which simply tokenizes the embedding of time series in the variate di-323 mension. By contrast, TwinsFormer learns the inherent dependencies from the interactions between

decomposed components, which can better capture accurate patterns for multivariate time series.
 Although SCINet obtains the best performance on the PEMS04 dataset using hierarchical sample convolution and interaction, it is inferior to TwinsFormer on other datasets. Remarkably, Twins Former achieves leading performance when averaging all the subsets, affirming the capacity of our interactive strategy in modeling complex temporal dynamics.

Table 2: Short-term forecasting results on PEMS datasets. The lookback length is set to T = 96and all the results are averaged from all predictions $S \in \{12, 24, 48, 96\}$. A lower MSE or MAE indicates a better prediction. See Table 7 in the Appendix for the full results.

Models	Twins (O	Former urs)	iTrans (<mark>20</mark>	former)24)	Time (20	Mixer)24)	Patc (20	hTST)23)	RLi (<mark>2(</mark>	inear)23)	Cross (2	former	Ti (<mark>20</mark>	DE)23)	Time (20	esNet 23a)	DLi (20	near 23)	SC (20	INet 22a)	FEDf	former
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
PEMS03	3 0.109	0.219	0.116	0.226	0.145	0.253	0.180	0.291	0.495	0.472	0.169	0.281	0.326	0.419	0.147	0.248	0.278	0.375	<u>0.114</u>	<u>0.224</u>	0.213	0.327
PEMS04	4 <u>0.111</u>	<u>0.219</u>	0.121	0.232	0.162	0.268	0.195	0.307	0.526	0.491	0.209	0.314	0.353	0.437	0.129	0.241	0.295	0.388	0.092	0.202	0.231	0.337
PEMS0	7 0.094	0.196	<u>0.100</u>	<u>0.204</u>	0.152	0.248	0.211	0.303	0.504	0.478	0.235	0.315	0.380	0.440	0.124	0.225	0.329	0.395	0.119	0.234	0.165	0.283
PEMS08	8 0.133	0.222	<u>0.151</u>	<u>0.234</u>	0.209	0.296	0.280	0.321	0.529	0.487	0.268	0.307	0.441	0.464	0.193	0.271	0.379	0.416	0.158	0.244	0.286	0.358
Avg	0.112	0.214	0.122	0.224	0.167	0.266	0.217	0.305	0.514	0.482	0.220	0.304	0.375	0.440	0.148	0.246	0.320	0.394	<u>0.121</u>	<u>0.222</u>	0.224	0.327

4.2 Ablation studies

To verify the effectiveness of each main component of TwinsFormer, we provide indispensable ablation studies for every possible design on decomposition and interactions. To be concrete, we disable or replace certain designs as model variants and experiment on two long-term (i.e., ECL and Traffic) and two short-term forecasting (i.e., PEMS03 and PEMS07) datasets. As seen in Table 3, we conduct an insightful analysis of decomposition and interactions through the following observation.

Table 3: Ablation studies for TwinsFormer. We disable or replace each component of both decomposition and interactions over four datasets. ✓ and ✗ indicate with and without certain components, respectively. The average results of all predicted lengths are listed here. See Table 9 in the Appendix for complete ablation results.

Design	Decomposition		Int	eractio	ns		E0	CL	Tra	uffic	PEN	1S03	PEN	4S07
8		-	E'_T	A_S	F_S	σ	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
TwinsFormer	 ✓ 	1	1	1	1	1	0.167	0.262	0.406	0.273	0.109	0.219	0.094	0.196
1	×	11	1	1	1	1	0.176	0.272	0.417	0.282	0.116	0.226	0.102	0.204
2	swap	1	1	1	1	1	0.172	0.265	0.413	0.277	0.114	0.224	0.101	0.204
3	1	+	1	1	1	1	0.180	0.275	0.416	0.283	0.118	0.228	0.102	0.207
4	1	1	X	1	1	1	0.185	0.278	0.418	0.283	0.122	0.232	0.105	0.210
5		1	1	X	1	1	0.183	0.277	0.413	0.278	0.118	0.229	0.103	0.208
6		1	1	1	X	1	0.176	0.271	0.412	0.281	0.121	0.228	0.104	0.210
\overline{O}		1	1	1	1	X	0.176	0.268	0.413	0.278	0.118	0.228	0.107	0.215

363

364

365

366

367

343

344

345

346

347

348

Ablation on decomposition Considering that the trend and seasonal components in the decomposition design are fed to different network branches, we disable the decomposition by using two original observed series as inputs (i.e., ①) and swap trend and seasonal components (i.e., ②) for ablation analysis. In ablation ① and ②, we can find significant decreases in forecasting performance for both long and short-term predictions, which demonstrates that our integration of the decomposition into Transformer architecture is reasonable and effective.

368 **Ablation on interactions** For the interactions, we verify the effectiveness by removing or replacing 369 components gradually. In ablation (3), we replace the subtraction mechanism (i.e., -) with original 370 addition skip connections (i.e., +), and the results on 3 show a decline in forecasting accuracy. 371 This illustrates that decomposed components can better satisfy the requirements of Transformer ar-372 chitecture by using the subtraction mechanism. Meanwhile, the results in 3 further highlight the 373 rationality of the decomposition design, which is consistent with the rationality analysis in Section 374 3.2. In ablations (4), (5), (6), and (7), we eliminate the impact of E'_T , A_S , F_S , and gate mechanism 375 σ for interactive learning, respectively. These four ablations all cause serious drops in forecasting performance, which indicates that all inputs for interactive learning can effectively boost the perfor-376 mance of TwinsFormer. The above observations highlight the substantial influence of our strategy 377 using residual and interactive learning in Transformer architecture.

Мс	odels	Transt (20	former 17)	Info (20	rmer 21)	Autof (20	ormer 21)	Flowf (20	former (22)	Period (20	former (23)
Me	etrics	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ECL	Original + Twins	0.277 0.167	0.372 0.262	0.311 0.168	0.397 0.261	0.227 0.176	0.338 0.267	0.267 0.168	0.359 0.262	0.219 0.170	0.328 0.262
	Promotion	39.7%	29.6%	46.0%	43.7%	22.5%	21.0%	37.1%	27.0%	22.4%	20.1%
Traffic	Original + Twins	0.665 0.406	0.363 0.273	0.764 0.431	0.416 0.282	0.628 0.433	0.379 0.288	0.750 0.424	0.421 0.280	0.608 0.439	0.373 0.287
	Promotion	38.9%	24.8%	43.6%	32.2%	31.1%	24.0%	43.5%	33.5%	27.8%	23.1%
PEMS03	Original + Twins	0.137 0.109	0.237 0.219	0.193 0.105	0.290 0.214	0.667 0.110	0.601 0.220	0.140 0.109	0.245 0.219	0.265 0.111	0.368 0.220
	Promotion	20.4%	7.6%	45.6%	26.2%	83.5%	63.4%	22.1%	10.6%	58.1%	40.2%
PEMS07	Original + Twins	0.178 0.094	0.243 0.196	0.194 0.092	0.259 0.194	0.367 0.096	0.451 0.200	0.178 0.096	0.240 0.198	0.200 0.098	0.318 0.201
	Promotion	47.2%	19.3%	52.6%	25.1%	73.8%	55.7%	46.1%	17.5%	51.0%	36.8%

378	Table 4: Attention compatibility and performance promotion obtained by applying various efficient
379	attention mechanisms to our interactive framework. The average results of all predicted lengths are
380	listed here. See Table 10 in the Appendix for the full results.

4.3 MODEL ANALYSIS

Compatibility and promotion We evaluate TwinsFormer by applying our interactive strategy to original Transformer (Vaswani et al., 2017) and its variants, which generally address the quadratic complexity of the self-attention mechanism, including Informer (Zhou et al., 2021), Autoformer (Wu et al., 2021), Flowformer (Wu et al., 2022) and Periodformer (Liang et al., 2023). As seen in Table 4, our framework can be adapted to various attention mechanisms with promoted performance for Transformer-based forecasters. On the one hand, the performance under different attention mechanisms illustrates the favorable attention compatibility of TwinsFormer. On the other hand, the performance promotion for different Transformer-based architectures exhibits the superiority of TwinsFormer. Overall, it achieves averaged 28.4% promotion on Transformer, 39.4% on Informer, **46.9%** on Autoformer, **29.7%** on Flowformer and **34.9%** on Periodformer.



418 419

396 397

398 399

400

401

402

403

404

405

406

407

Figure 4: Forecasting performance with different lookback lengths on three datasets.

420 As argued in (Zeng et al., 2023) and (Nie et al., 2023), most of Lookback length sensitivity 421 the Transformer-based models will not improve the forecasting performance with an increasing 422 lookback length due to the distracted attention on the longer input (Liu et al., 2024). However, 423 our TwinsFormer reduces the MSE scores with enlarged historical information to be utilized, which is consistent with the theoretical analysis by statistical methods (Box & Jenkins, 1968). As seen 424 in Figure 4, the forecasting results keep improving in most cases where the prediction length S425 belongs to $\{96, 192, 336, 720\}$ as the receptive field increases. These improvements confirm that 426 our TwinsFormer can effectively capture inherent dependencies from a longer lookback window. 427

428 **Visualization analysis** To provide an intuitive understanding of the learned representations by our dual-stream framework, we visualize the multivariate correlations, the corresponding representa-429 tions, and prediction results in Figure 5. It can be observed that the multivariate correlations learned 430 by iTransformer are less distinct than those of TwinsFormer in the gold dashed box. Accordingly, the 431 learned representations obtained by TwinsFormer have more abundant variate and temporal infor-



Figure 5: Analysis of multivariate correlations and series representations. Zoom in for more details.

mation than those of iTransformer, which we highlight with the gold and red dashed boxes in (b) of Figure 5. Consequently, TwinsFormer has better accuracy than iTransformer in forecasting performance, which is circled with a red line in the (c) part. Those observations indicate that TwinsFormer can better learn inherent dependencies among time series to achieve more accurate forecasting.

454 Efficiency analysis Our TwinsFormer is 455 a Transformer-based architecture with dual-456 stream interactions, where the trend branch is 457 composed of linear layers and sigmoid acti-458 vation functions. Therefore, like other Trans-459 former models, the main complexity of Twins-460 Former is $O(N^2)$, which comes from the 461 seasonal branch with the attention module. Note that, the N for TwinsFormer is related 462 to the number of variates, while the N for 463 most Tansformer-based models is affected by 464 the lookback length. Going further, the ef-465 ficiency of TwinsFormer exceeds most Trans-



Figure 6: Model efficient comparison on Traffic.

former Variants in datasets with a relatively small number of variates (i.e., N < 96). Although the memory cost of TwinsFormer is not dominant when $N \gg 96$, we can choose only a part of the variates based on the correlations among variates to improve the training efficiency. As shown in Figure 6, TwinsFormer can obtain the best performance compared with 8 baselines on Traffic dataset, since the multivariate correlations can be well explored. Meanwhile, our efficient TwinsFormer-E (trained with 20% variates and prediction for all variates) can achieve comparable forecasting performance while substantially reducing the memory footprint.

473 474

475

449 450

451

452

453

5 CONCLUSION AND FUTURE WORK

476 Considering that the decomposition design can mine temporal patterns and the attention mechanism 477 can capture multivariate correlations, we propose TwinsFormer, a Transformer-based framework 478 revisiting inherent dependencies via two interactive branches for time series forecasting. Empow-479 ered with our interactive module, TwinsFormer handily incorporates the idea of decomposition into 480 Transformer architectures and effectively learns time series representations. Experimentally, Twins-Former achieves state-of-the-art performances in both long-term and short-term forecasting tasks. 481 Furthermore, the detailed visualization, ablations, and analysis illustrate the effectiveness and gen-482 erality of our framework. Benefiting from the multivariate correlations learned with variate tokens, 483 TwinsFormer demonstrates favorable run-time efficiency for high-dimension channel datasets. In 484 the future, we will explore more efficient interaction design with decomposed components in the 485 MLP architectures and analyze the performance of the interaction design in more time series tasks.

486 REFERENCES

488 489	George EP Box and Gwilym M Jenkins. Some recent advances in forecasting and control. <i>Journal</i> of the Royal Statistical Society. Series C (Applied Statistics), 17(2):91–109, 1968.
490 491	Abhimanyu Das, Weihao Kong, Andrew Leach, Shaan Mathur, Rajat Sen, and Rose Yu. Long-term
492 493	Torecasting with tide: Time-series dense encoder. Trans. Mach. Learn. Res., 2025, 2025.
494 495	60th International Midwest Symposium on circuits and systems (MWSCAS), pp. 1597–1600, 2017.
496 497 498	Dazhao Du, Bing Su, and Zhewei Wei. Preformer: predictive transformer with multi-scale segment- wise correlations for long-term time series forecasting. In <i>IEEE International Conference on</i> <i>Acoustics, Speech and Signal Processing (ICASSP)</i> , pp. 1–5, 2023.
499 500 501	Clive William John Granger and Paul Newbold. <i>Forecasting economic time series</i> . Academic press, 2014.
502 503 504	Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In 3rd Interna- tional Conference on Learning Representations, ICLR, 2015.
505 506 507	Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term temporal patterns with deep neural networks. In <i>The 41st international ACM SIGIR conference on research & development in information retrieval</i> , pp. 95–104, 2018.
508 509 510 511 512	Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyou Zhou, Wenhu Chen, Yu-Xiang Wang, and Xifeng Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. In <i>Advances in Neural Information Processing Systems, NeurIPS</i> , pp. 5244–5254, 2019.
513 514 515	Zhe Li, Shiyi Qi, Yiduo Li, and Zenglin Xu. Revisiting long-term time series forecasting: An investigation on linear mapping. <i>arXiv preprint arXiv:2305.10721</i> , 2023.
516 517 518	Daojun Liang, Haixia Zhang, Dongfeng Yuan, Xiaoyan Ma, Dongyang Li, and Minggao Zhang. Does long-term series forecasting need complex attention and extra long inputs? <i>arXiv preprint arXiv:2306.05035</i> , 2023.
519 520 521 522	Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. Scinet: Time series modeling and forecasting with sample convolution and interaction. In <i>Advances in</i> <i>Neural Information Processing Systems, NeurIPS</i> , 2022a.
523 524 525 526	Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar. Pyraformer: Low-complexity pyramidal attention for long-range time series modeling and fore- casting. In <i>The Tenth International Conference on Learning Representations, ICLR</i> , 2022b.
527 528 529	Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Exploring the stationarity in time series forecasting. <i>Advances in Neural Information Processing Systems, NeurIPS</i> , 35:9881–9893, 2022c.
530 531 532 533	Yong Liu, Chenyu Li, Jianmin Wang, and Mingsheng Long. Koopa: Learning non-stationary time series dynamics with koopman predictors. In <i>Advances in Neural Information Processing Systems, NeurIPS</i> , 2023.
534 535 536 537	Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. itransformer: Inverted transformers are effective for time series forecasting. In <i>The Twelfth Inter-</i> <i>national Conference on Learning Representations, ICLR</i> , 2024.
538 539	Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In <i>The Eleventh International Conference</i> on Learning Representations, ICLR, 2023.

540 541 542 543	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Ed- ward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
544 545	performance deep learning library. In Advances in Neural Information Processing Systems, NeurIPS, pp. 8024–8035, 2019.
546 547	Cleveland Robert, C William, and Terpenning Irma. Stl: A seasonal-trend decomposition procedure based on loess. J Off Stat, 6:3–73, 1990.
548 549 550 551	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>Advances in Neural Information Processing Systems, NeurIPS</i> , pp. 5998–6008, 2017.
552 553 554 555	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 <i>The Twelfth Interna-</i> <i>tional Conference on Learning Representations, ICLR</i> , 2023.
556 557 558	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 <i>The Twelfth International Conference on Learning Representations, ICLR</i> , 2024.
559 560 561	Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans- formers with auto-correlation for long-term series forecasting. In <i>Advances in Neural Information</i> <i>Processing Systems, NeurIPS</i> , pp. 22419–22430, 2021.
563 564 565	Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Flowformer: Linearizing transformers with conservation flows. In <i>International Conference on Machine Learning, ICML</i> , volume 162, pp. 24226–24242, 2022.
566 567 568	Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In <i>The Eleventh International</i> <i>Conference on Learning Representations, ICLR</i> , 2023a.
569 570 571 572	Haixu Wu, Hang Zhou, Mingsheng Long, and Jianmin Wang. Interpretable weather forecasting for worldwide stations with a unified deep model. <i>Nature Machine Intelligence</i> , 5(6):602–611, 2023b.
573 574 575	Xueyan Yin, Genze Wu, Jinze Wei, Yanming Shen, Heng Qi, and Baocai Yin. Deep learning on traffic prediction: Methods, analysis, and future directions. <i>IEEE Transactions on Intelligent Transportation Systems</i> , 23(6):4927–4943, 2021.
576 577 578	Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 37, pp. 11121–11128, 2023.
580 581 582	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, ICLR</i> , 2023.
583 584 585	Zheng Zhao, Weihai Chen, Xingming Wu, Peter CY Chen, and Jingmeng Liu. Lstm network: a deep learning approach for short-term traffic forecast. <i>IET intelligent transport systems</i> , 11(2): 68–75, 2017.
586 587 588 589	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, volume 35, pp. 11106–11115, 2021.
590 591 592 593	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, ICML</i> , volume 162, pp. 27268–27286, 2022.

594 A IMPLEMENTATION DETAILS

596

597 598

600

613

614 615

618

619

Benchmarks details We evaluate the performance of TwinsFormer compared with various baselines on 13 well-established benchmarks¹, which are detailed in Table 5.

Table 5: Detailed descriptions of benchmarks. Channel denotes the number of variates in each dataset. Prediction length points out four prediction settings. The dataset size is split in (Train, Validation, Test). Frequency denotes the sampling interval of time points.

Tasks	Benchmarks	Channels	Prediction Length	Dataset Size	Frequency	Information
	ETTm1	7		(34465, 11521, 11521)	15min	Electricity
	ETTm2	7		(34465, 11521, 11521)	15min	Electricity
	ETTh1	7		(8545, 2881, 2881)	Hourly	Electricity
Long term	ETTh2	7		(8545, 2881, 2881)	Hourly	Electricity
Eoroposting	ECL	321	{96, 192, 336, 720}	(18317, 2633, 5261)	Hourly	Electricity
rorceasting	Traffic	862		(12185, 1757, 3509)	Hourly	Transportation
	Exchange	8		(5120, 665, 1422)	Daily	Economy
	Weather	21		(36792, 5271, 10540)	10min	Weather
	Solar-energy	137		(36601, 5161, 10417)	10min	Electricity
	PEMS03	358		(15617, 5135, 5135)	5min	Transportation
Short-term	PEMS04	307	512 24 48 061	(10172, 3375, 3375)	5min	Transportation
Forecasting	PEMS07	883	12, 24, 40, 90	(16911, 5622, 5622)	5min	Transportation
	PEMS08	170		(10690, 3548, 3548)	5min	Transportation

Metrics details Regarding evaluation metrics, we utilize the mean square error (MSE) and mean absolute error (MAE) for long-term and short-term forecasting:

$$MSE = \frac{1}{L} \sum_{i=1}^{L} (X_i - \hat{X}_i)^2, \qquad MAE = \sum_{i=1}^{L} |X_i - \hat{X}_i|,$$

where $X, \hat{X} \in \mathbb{R}^{L \times N}$ denote the ground truth and prediction results for N variates in the future L time steps. $|\cdot|$ means the absolute value operation.

Algorithm details We provide the pseudo-code of TwinsFormer in Algorithm 1.

Algorithm 1 Workflow of our TwinsFormer. 620 **Input:** Input lookback time series $X \in \mathbb{R}^{T \times N}$; Input length T, prediction length L, and variates 621 622 number N; Token dimension D, TwinsBlock number M, and moving average kernel size k. 623 **Output:** The prediction results $\hat{X} \in \mathbb{R}^{L \times N}$. 624 1: b Using the moving average kernel and padding operations to decompose time series. 625 $\triangleright X_T, X_S \in \mathbb{R}^{T \times N}$ 2: $X_T = AvgPool(Padding(X)), X_S = X - X_T$ 626 3: ▷ Embedding series into variate tokens by Multi-layer Perceptron. 627 $\triangleright E_T^0, E_S^0 \in \mathbb{R}^{N \times D}$ 4: $E_T^0 = Embed_T(X_T.transpose), E_S^0 = Embed_S(X_S.transpose)$ 628 629 5: ▷ Running through TwinsFormer blocks. 630 6: for m in $\{1, \dots, M\}$ do 631 7: ▷ Self-attention mechanism and feed-forward network are applied for the seasonal branch. 632 $E_S^{m-1} = LayerNorm(E_S^{m-1} - \text{Self-Attn}(E_S^{m-1}))$ $\triangleright E_S^{m-1} \in \mathbb{R}^{N \times D}$ 8: 633 $\triangleright E_S^m \in \mathbb{R}^{N \times D}$ $E_S^m = E_S^{m-1} - \text{Feed-Forward}(E_S^{m-1})$ 9: 634 ▷ Interactive module (IM) is utilized for the trend branch. 10: 635
$$\begin{split} E_T^{m-1} &= Sigmoid(Conv(\sum Multiconv(\text{Self-Attn}(E_S^{m-1})))) \cdot E_T^{m-1} \ \triangleright E_T^{m-1} \in \mathbb{R}^{N \times D} \\ E_T^m &= Concate([\text{Self-Attn}(E_S^{m-1}), \text{Feed-Forward}(E_S^{m-1}), E_T^{m-1}]) \ \quad \triangleright E_T^m \in \mathbb{R}^{N \times 3 \times D} \end{split}$$
636 11: 637 12: 638 13. ▷ Adding gate mechanism to seasonal and trend branches. 639 $\triangleright E_S^m \in \mathbb{R}^{N \times D}$ $\triangleright E_T^m \in \mathbb{R}^{N \times D}$ $E_S^m = LayerNorm(Sigmoid(Conv(E_S^m)) * E_S^m)$ 14: 640 15: $E_T^m = Sigmoid(Conv(E_T^m)) * E_T^m$ 641 16: end for 642 $\triangleright \, \hat{X} \in \mathbb{R}^{N \times L}$ 17: $\hat{X} = Projector(E_S^m + E_T^m)$ 643 $\triangleright \hat{X} \in \mathbb{R}^{L \times N}$ 644 18: $\hat{X} = \hat{X}.transpose$ 645 19: return \hat{X} 646

⁶⁴⁷

¹All the datasets are publicly available at https://github.com/thuml/iTransformer

648 B BASELINE METHODS 649

650 651	We provide brief descriptions for the selected baselines as follows:
652 653	• iTransformer (Liu et al., 2024) is a Transformer-based model that captures multivari- ate correlations on the variate dimension for forecasting. The source code is available
654	athttps://github.com/thuml/iTransformer.
655	• TimeMixer (Wang et al., 2024) is an MLP-based predictor that adopts Past-Decomposable-
656	Mixing and Future-Multipredictor-Mixing blocks to aggregate disentangled information.
657	The source code is available at https://github.com/kwuking/TimeMixer.
658	• PatchTST (Nie et al., 2023) is a Transformer-based model that utilizes patching design
659 660	and channel-independence to learn temporal patterns for forecasting. The source code is available at https://github.com/yuqinie98/PatchTST.
661	• RLinear (Li et al., 2023) is an MLP-based model that employs linear mapping with re-
663	versible normalization and independent channel operations for forecasting. The source code is available at https://github.com/plumprc/RTSF.
664	• Crossformer (Zhang & Yan 2023) is a Transformer-based predictor capturing cross-time
666	and cross-dimension dependencies with dimension-segment-wise embedding. The source
667	code is available at https://github.com/Thinklab-SJTU/Crossformer.
668	• TiDE (Das et al., 2023) is an MLP-based predictor that handles covariates and non-
669	linear dependencies with the dense encoder. The source code is available at https:
670	//github.com/ZihangHLiu/TiDE.
671	• TimesNet (Wu et al., 2023a) is a CNN-based model that ravels out the complex tempo-
672	ral variations into multiple intraperiod- and interperiod-variations for general time series
673	analysis. The source code is available at https://github.com/thuml/TimesNet.
674	• DLinear (Zeng et al., 2023) is an MLP-based model that combines a decomposition
675 676	scheme with two one-layer linear layers for forecasting. The source code is available at https://github.com/cure-lab/LTSFLinear.
677	• SCINet (Line et al. 2022a) is a CNN based model that conducts sample convolution and
678	interaction for temporal modeling and forecasting. The source code is available at https:
679	//github.com/cure-lab/SCINet.
680	• FEDformer (Zhou et al. 2022) is a Transformer-based model that utilizes the seasonal-
681	trend decomposition with frequency-enhanced blocks to learn temporal dependency for
682	forecasting. The source code of FEDformer is available at https://github.com/
683	MAZiqing/FEDformer.
684	• Transformer (Vaswani et al., 2017) utilizes self-attention mechanism to capture cross-
696	time dependency for forecasting. The source code of Transformer is available at https:
687	//github.com/zhouhaoyi/Informer2020.
688	• Informer (Zhou et al., 2021) is a Transformer-based model using the ProbSparse self-
689	attention to learn cross-time correlations for forecasting. The source code of Informer is
690	available at https://github.com/zhouhaoyi/Informer2020.
691	• Autoformer (Wu et al., 2021) is a Transformer-based predictor introducing decomposition
692	design and Auto-Correlation mechanism to capture temporal dependency. The source code
693	of Autoformer is available at https://github.com/thuml/Autoformer.
694	• Flowformer (Wu et al., 2022) is a Transformer-based model that replaces the self-attention
695	mechanism with the flow-attention. The source code is available at https://github.
696	com/thuml/Flowformer.
697	• Periodformer (Liang et al., 2023) is a Transformer-based predictor capturing the periodic-
698	ity of time series based on the relationship between different moments. The source code is
699	available at https://github.com/Anoise/Minusformer.
700	To ensure the fairness of comparison, all the baselines are reproduced with the same repository
/01	which is available at https://github.com/thuml/Time-Series-Library.

702 C FULL MAIN RESULTS

Due to the space limitation, we provide the full multivariate forecasting results here. Specifically, Table 6 contains the detailed results of all prediction lengths on 9 well-acknowledged benchmarks for long-term forecasting, while Table 7 includes the full short-term forecasting results on 4 challenging citywide traffic datasets.

709Table 6: Full results of the long-term forecasting task. We compare extensive competitive models710under different prediction lengths $S \in \{96, 192, 336, 720\}$. The input sequence length is set to 96711for all baselines. Avg means the average results from all four prediction lengths.

М	odels	Twins (O	Former urs)	r iTrans (<mark>20</mark>	former 24)	Time (20	Mixer 24)	Patch (20	nTST 23)	RLi (20	near 23)	Cross (20	former 23)	Ti (20	DE 123)	Time (20)	esNet 23a)	DLi (<mark>20</mark>	near 23)	SCI (202	Net 22a)	FED	former
Μ	letric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	96 192 336 720	$\begin{array}{c} \underline{0.325} \\ 0.372 \\ 0.406 \\ 0.467 \end{array}$	0.364 0.390 0.412 0.448	0.334 0.377 0.426 0.491	0.368 0.391 0.420 0.459	0.320 0.362 0.396 0.458	0.355 0.382 0.406 0.445	0.329 0.367 0.399 0.454	0.367 <u>0.385</u> <u>0.410</u> 0.439	0.355 0.391 0.424 0.487	0.376 0.392 0.415 0.450	0.404 0.540 0.532 0.666	0.426 0.451 0.515 0.589	0.364 0.398 0.428 0.487	0.387 0.404 0.425 0.461	0.338 0.374 0.410 0.478	0.375 0.387 0.411 0.450	0.345 0.380 0.413 0.474	0.372 0.389 0.413 0.453	0.418 0.439 0.490 0.595	0.438 0.450 0.485 0.550	0.379 0.426 0.445 0.543	0.419 0.441 0.459 0.490
	Avg	0.393	0.404	0.407	0.410	0.384	0.397	<u>0.387</u>	<u>0.400</u>	0.414	0.407	0.513	0.496	0.419	0.419	0.400	0.406	0.403	0.407	0.485	0.481	0.448	0.452
ETTm2	96 192 336 720	0.173 0.239 0.298 0.397	0.256 0.300 0.339 0.397	0.180 0.250 0.311 0.412	0.264 0.309 0.348 0.407	0.176 0.242 0.303 0.396	0.259 0.303 0.339 0.399	0.175 0.241 0.305 0.402	0.259 0.302 0.343 0.400	0.182 0.246 0.307 0.407	0.265 0.304 <u>0.342</u> 0.398	0.287 0.414 0.597 1.730	0.366 0.492 0.542 1.042	0.207 0.290 0.377 0.558	0.305 0.364 0.422 0.524	0.187 0.249 0.321 0.408	0.267 0.309 0.351 0.403	0.193 0.284 0.369 0.554	0.292 0.362 0.427 0.522	0.286 0.399 0.637 0.960	0.377 0.445 0.591 0.735	0.203 0.269 0.325 0.421	0.287 0.328 0.366 0.415
_	Avg	0.277	0.323	0.288	0.332	0.279	0.325	0.281	0.326	0.286	0.327	0.757	0.610	0.358	0.404	0.291	0.333	0.350	0.401	0.571	0.537	0.305	0.349
ETTh1	96 192 336 720	0.385 0.439 0.480 0.480	$0.401 \\ 0.431 \\ \underline{0.452} \\ \underline{0.474}$	0.386 0.441 0.487 0.503	0.405 0.436 0.458 0.491	0.384 0.437 0.472 0.586	0.400 0.429 0.446 0.531	0.414 0.460 0.501 0.500	0.419 0.445 0.466 0.488	0.386 0.437 <u>0.479</u> <u>0.481</u>	0.395 0.424 0.446 0.470	0.423 0.471 0.570 0.653	0.448 0.474 0.546 0.621	0.479 0.525 0.565 0.594	0.464 0.492 0.515 0.558	0.384 0.436 0.491 0.521	0.402 0.429 0.469 0.500	0.386 0.437 0.481 0.519	0.400 0.432 0.459 0.516	0.654 0.719 0.778 0.836	0.599 0.631 0.659 0.699	0.376 0.420 0.459 0.506	0.419 0.448 0.465 0.507
	Avg	0.446	<u>0.440</u>	0.454	0.447	0.470	0.451	0.469	0.454	0.446	0.434	0.529	0.522	0.541	0.507	0.458	0.450	0.456	0.452	0.747	0.647	0.440	0.460
ETTh2	96 192 336 720	0.292 0.375 0.417 0.406	0.345 0.395 0.429 0.430	0.297 0.380 0.428 0.427	0.349 0.400 0.432 0.445	0.297 0.369 0.427 0.462	0.348 0.392 0.435 0.463	0.302 0.388 0.426 0.431	0.348 0.400 0.433 0.446	0.288 0.374 0.415 0.420	0.338 0.390 0.426 <u>0.440</u>	0.745 0.877 1.043 1.104	0.584 0.656 0.731 0.763	0.400 0.528 0.643 0.874	0.440 0.509 0.571 0.679	0.340 0.402 0.452 0.462	$\begin{array}{c} 0.374 \\ 0.414 \\ 0.452 \\ 0.468 \end{array}$	0.333 0.477 0.594 0.831	0.387 0.476 0.541 0.657	0.707 0.860 1.000 1.249	0.621 0.689 0.744 0.838	0.358 0.429 0.496 0.463	0.397 0.439 0.487 0.474
	Avg	0.373	<u>0.400</u>	0.383	0.407	0.389	0.409	0.387	0.407	<u>0.374</u>	0.398	0.942	0.684	0.611	0.550	0.414	0.427	0.559	0.515	0.954	0.723	0.437	0.449
ECL	96 192 336 720	0.139 0.158 0.172 0.200	0.233 0.252 0.267 0.293 0.262	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{r} 0.240 \\ 0.253 \\ 0.269 \\ 0.317 \\ 0.270 \\ \end{array} $	0.153 0.168 0.185 0.227	0.244 0.259 0.275 0.312 0.272	0.181 0.188 0.204 0.246	0.281 0.274 0.293 0.324	0.201 0.201 0.215 0.257	0.281 0.283 0.298 0.331 0.298	0.219 0.231 0.246 0.280	0.314 0.322 0.337 0.363 0.334	0.237 0.236 0.249 0.284	0.329 0.330 0.344 0.373	0.168 0.184 0.198 0.220	0.272 0.289 0.300 0.320	0.197 0.196 0.209 0.245	0.282 0.285 0.301 0.333 0.300	0.247 0.257 0.269 0.299	0.345 0.355 0.369 0.390 0.365	0.193 0.201 0.214 0.246	0.308 0.315 0.329 0.355
	96	0.081	0.200	10.086	0.206	0.099	0.218	0.088	0.205	0.093	0.217	0 256	0.367	0 094	0.218	0 107	0.234	0.088	0.218	0.267	0.396	0 148	0.278
Exchange	192 336 720	0.172 0.320 0.812	0.295 0.409 0.677	0.177 0.331 0.847	0.200 0.299 0.417 0.691	0.196 0.359 0.864	0.313 0.432 0.703	0.000 0.176 0.301 0.901	0.299 0.397 0.714	0.184 0.351 0.886	0.307 0.432 0.714	0.250 0.470 1.268 1.767	0.509 0.883 1.068	0.184 0.349 0.852	0.210 0.307 0.431 0.698	0.226 0.367 0.964	0.234 0.344 0.448 0.746	0.176 0.313 0.839	0.315 0.427 0.695	0.351 1.324 1.058	0.459 0.853 0.797	0.271 0.460 1.195	0.276 0.315 0.427 0.695
	Avg	0.346	0.395	<u>0.360</u>	<u>0.403</u>	0.380	0.417	0.367	0.404	0.378	0.417	0.940	0.707	0.370	0.413	0.416	0.443	0.354	0.414	0.750	0.626	0.519	0.429
Traffic	96 192 336 720	0.382 0.392 0.410 0.442	0.260 0.267 0.276 0.292	$ \begin{array}{r} 0.395 \\ 0.417 \\ 0.433 \\ 0.467 \end{array} $	$ \begin{array}{r} $	0.473 0.486 0.488 0.536	0.287 0.294 0.298 0.314	0.462 0.466 0.482 0.514	0.295 0.296 0.304 0.322	0.649 0.601 0.609 0.647	0.389 0.366 0.369 0.387	0.522 0.530 0.558 0.589	0.290 0.293 0.305 0.328	0.805 0.756 0.762 0.719	0.493 0.474 0.477 0.449	0.593 0.617 0.629 0.640	0.321 0.336 0.336 0.350	0.650 0.598 0.605 0.645	0.396 0.370 0.373 0.394	0.788 0.789 0.797 0.841	0.499 0.505 0.508 0.523	0.587 0.604 0.621 0.626	0.366 0.373 0.383 0.382
	Avg	0.406	0.273	0.428	0.282	0.496	0.298	0.481	0.304	0.626	0.378	0.550	0.304	0.760	0.473	0.620	0.336	0.625	0.383	0.804	0.509	0.610	0.376
Weather	96 192 336 720	0.161 <u>0.211</u> <u>0.266</u> <u>0.347</u>	0.201 0.248 0.291 0.343	0.174 0.221 0.278 0.358	0.214 0.254 0.296 0.347	0.163 0.209 0.264 0.345	0.209 0.252 0.293 0.345	0.177 0.225 0.278 0.354	0.218 0.259 0.297 0.348	0.192 0.240 0.292 0.364	0.232 0.271 0.307 0.353	0.158 0.206 0.272 0.398	0.230 0.277 0.335 0.418	0.202 0.242 0.287 0.351	0.261 0.298 0.335 0.386	0.172 0.219 0.280 0.365	0.220 0.261 0.306 0.359	0.196 0.237 0.283 0.345	0.255 0.296 0.335 0.381	0.221 0.261 0.309 0.377	0.306 0.340 0.378 0.427	0.217 0.276 0.339 0.403	0.296 0.336 0.380 0.428
	Avg	<u>0.246</u>	0.271	0.258	0.278	0.245	<u>0.274</u>	0.259	0.281	0.272	0.291	0.259	0.315	0.271	0.320	0.259	0.287	0.265	0.317	0.292	0.363	0.309	0.360
Solar-Energy	96 192 336 720	$ \begin{array}{r} 0.193 \\ 0.223 \\ 0.246 \\ 0.245 \end{array} $	0.224 0.250 0.268 0.272	0.203 0.233 0.248 0.249	$\begin{array}{r} \underline{0.237} \\ \underline{0.261} \\ \underline{0.273} \\ \underline{0.275} \end{array}$	0.189 0.222 0.231 0.223	0.259 0.283 0.292 0.285	0.234 0.267 0.290 0.289	0.286 0.310 0.315 0.317	0.322 0.359 0.397 0.397	0.339 0.356 0.369 0.356	0.310 0.734 0.750 0.769	0.331 0.725 0.735 0.765	0.312 0.339 0.368 0.370	0.399 0.416 0.430 0.425	0.250 0.296 0.319 0.338	0.292 0.318 0.330 0.337	0.290 0.320 0.353 0.356	0.378 0.398 0.415 0.413	0.237 0.280 0.304 0.308	0.344 0.380 0.389 0.388	0.242 0.285 0.282 0.357	0.342 0.380 0.376 0.427
	Avg	<u>0.227</u>	0.254	0.233	<u>0.262</u>	0.216	0.280	0.270	0.307	0.369	0.356	0.641	0.639	0.347	0.417	0.301	0.319	0.330	0.401	0.282	0.375	0.291	0.381
1^{st}	Count	22	30	0	0	<u>18</u>	6	2	2	4	<u>9</u>	0	0	0	0	0	0	0	0	0	0	0	0

TwinsFormer achieves the best forecasting performance among 11 models on various prediction
horizons for both long-term and short-term forecasting tasks. To be concrete, TwinsFormer outperforms all the baselines on 52 out of the 90 settings including different prediction lengths and metrics
over 9 long-term benchmarks. Meanwhile, TwinsFormer beats all the baselines on 29 out of the 40 settings of varying prediction lengths and metrics over 4 short-term datasets.

759																								
760	М	odels	Twins (O	Former	· iTrans	former	Time	Mixer	Patel	hTST	RLi	near	Cross	former	Ti (20	DE	Time (20	esNet	DLi	inear	SCI	Net	FEDf	ormer
761						MAE			MOE	<u></u>	MCE	<u></u>		<u>MAE</u>		<u>MAE</u>		MAE		<u>MAE</u>				
762		letric	INISE	MAE	INISE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
763	\$03	12 24	0.065 0.086	0.169 0.196	0.069	0.175 0.208	0.077 0.112	0.187 0.224	0.099 0.121	0.216 0.240	0.126 0.246	0.236 0.334	0.090	0.203 0.240	0.178	0.305 0.371	0.085	0.192 0.223	0.122 0.201	0.243 0.317	0.066 0.085	0.172 0.198	0.126 0.149	0.251 0.275
764	PEMS	48 96	0.121 0.165	0.234 0.276	0.131	0.243 0.279	0.169 0.22	0.277 0.322	0.202 0.262	0.317 0.367	0.551 1.057	0.529 0.787	0.202	0.317 0.367	0.379	0.463 0.539	0.155 0.228	0.260 0.317	0.333 0.457	0.425 0.515	0.127 0.178	0.238 0.287	0.227 0.348	0.348 0.434
765	Í	Avg	0.109	0.219	0.116	0.226	0.145	0.253	0.180	0.291	0.495	0.472	0.169	0.281	0.326	0.419	0.147	0.248	0.278	0.375	<u>0.114</u>	<u>0.224</u>	0.213	0.327
766		12	0.077	<u>0.181</u>	0.081	0.188	0.092	0.203	0.105	0.224	0.138	0.252	0.098	0.218	0.219	0.340	0.087	0.195	0.148	0.272	0.073	0.177	0.138	0.262
767	MS04	24 48	0.095 0.120	$\frac{0.204}{0.231}$	0.099	0.211 0.247	0.127 0.188	0.239 0.294	0.153 0.229	0.275 0.339	0.258 0.572	0.348 0.544	0.131 0.205	0.256 0.326	0.292	0.398 0.478	0.103	0.215 0.250	0.224	0.340 0.437	0.084 0.099	0.193 0.211	0.177 0.270	0.293 0.368
768	PE	96	0.150	0.261	0.172	0.283	0.240	0.337	0.291	0.389	1.137	0.820	0.402	0.457	0.492	0.532	0.190	0.303	0.452	0.504	0.114	0.227	0.341	0.427
769		Avg	<u>0.111</u>	<u>0.219</u>	0.121	0.232	0.162	0.268	0.195	0.307	0.526	0.491	0.209	0.314	0.353	0.437	0.129	0.241	0.295	0.388	0.092	0.202	0.231	0.337
770	10	12	0.060	0.158	$\left \frac{0.067}{0.086} \right $	$\frac{0.167}{0.189}$	0.069	0.172	0.095	0.207	0.118	0.235	0.094	0.200	0.173	0.304	0.082	0.181	0.115	0.242	0.068	0.171	0.109	0.225
771	MS(48	0.104	0.209	0.000	0.109	0.185	0.212	0.253	0.340	0.562	0.541	0.311	0.369	0.271	0.385	0.134	0.238	0.398	0.458	0.149	0.225	0.125	0.244
772	PE	96	0.132	0.236	0.138	<u>0.244</u>	0.246	0.327	0.346	0.404	1.096	0.795	0.396	0.442	0.628	0.577	0.181	0.279	0.594	0.553	0.141	0.304	0.262	0.376
773		Avg	0.094	0.196	0.100	<u>0.204</u>	0.152	0.248	0.211	0.303	0.504	0.478	0.235	0.315	0.380	0.440	0.124	0.225	0.329	0.395	0.119	0.234	0.165	0.283
774	80	12 24	0.075 0.106	0.174 0.206	$\frac{0.080}{0.118}$	$\frac{0.183}{0.221}$	0.097 0.156	0.205 0.262	0.168 0.224	0.232 0.281	0.133 0.249	0.247 0.343	0.165	0.214 0.260	0.227	0.343 0.409	0.112	0.212 0.238	0.154 0.248	0.276 0.353	0.087	0.184 0.221	0.173 0.210	0.273 0.301
775	PEMS	48 96	0.167 0.184	0.258 0.251	0.186 0.221	$\frac{0.265}{0.267}$	0.269 0.313	0.345 0.373	0.321 0.408	0.354 0.417	0.569 1.166	0.544 0.814	0.315 0.377	0.355 0.397	0.497	0.510 0.592	0.198 0.320	0.283 0.351	0.440 0.674	0.470 0.565	0.189 0.236	0.270 0.300	0.320 0.442	0.394 0.465
776		Ανσ	0.133	0.222	10 151	0 234	0 209	0.296	0 280	0.321	0 529	0.487	 0.268	0 307	0 441	0 464	0 193	0 271	0 379	0.416	0 158	0 244	0 286	0.358
777	1 st	Court	14	15		0	0	0	0	0	0	0	0	0		0		0		0	6	5	0	0
778	1	Count	1.4	15	10	U	U	U	U	U	U	0		U	10	U		U	0	U	<u>v</u>	2	U	0

Table 7: Full results of the short-term forecasting task. We compare extensive competitive models under different prediction lengths $S \in \{12, 24, 48, 96\}$. The input sequence length is set to 96 for all baselines. Avg means the average results from all four prediction lengths.

ERROR BARS D

We obtain the standard deviation of TwinsFormer performance by training the model with 5 different random seeds over 12 datasets. As seen in Table 8, the error bars of all the results are tiny, which exhibits that the performance of TwinsFormer is robust and reliable.

Table 8: Robustness of TwinsFormer performance obtained from 5 random seeds on 12 benchmarks.

Dataset	ETT	ſm1	ETI	ſm2	ET	Th2	E	CL
Metrics	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
96	0.325 ± 0.001	0.364 ± 0.001	0.173 ± 0.001	0.256 ± 0.001	0.292 ± 0.001	0.345 ± 0.000	0.139 ± 0.000	0.233 ± 0.000
192	0.372 ± 0.001	0.390 ± 0.002	0.239 ± 0.001	0.300 ± 0.000	0.375 ± 0.002	0.395 ± 0.001	0.158 ± 0.001	0.252 ± 0.001
336	0.406 ± 0.002	0.412 ± 0.001	0.298 ± 0.000	0.339 ± 0.001	0.417 ± 0.004	0.429 ± 0.002	0.172 ± 0.002	0.267 ± 0.001
720	0.467 ± 0.002	0.448 ± 0.003	0.397 ± 0.002	0.397 ± 0.001	0.406 ± 0.003	0.430 ± 0.001	0.200 ± 0.004	0.293 ± 0.002
Dataset	Tra	ffic	Exch	ange	Solar-	Energy	Wea	ather
Metrics	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
96	$0.382 \pm 0.001 \ 0.260 \pm 0.000$		0.081 ± 0.001	0.200 ± 0.001	0.193 ± 0.002	0.224 ± 0.002	0.161 ± 0.000	0.201 ± 0.000
192	0.392 ± 0.002	0.267 ± 0.001	0.172 ± 0.001	0.295 ± 0.001	0.223 ± 0.002	0.250 ± 0.002	0.211 ± 0.001	0.248 ± 0.001
336	0.410 ± 0.003	0.276 ± 0.003	0.320 ± 0.002	0.409 ± 0.001	0.246 ± 0.000	0.268 ± 0.001	0.266 ± 0.002	0.291 ± 0.001
720	0.442 ± 0.001	0.292 ± 0.001	0.812 ± 0.012	0.677 ± 0.005	0.245 ± 0.001	0.272 ± 0.001	0.347 ± 0.001	0.343 ± 0.001
Dataset	PEM	1803	PEM	1S04	PEN	1S07	PEN	4S08
Metrics	MSE MAE		MSE	MAE	MSE	MAE	MSE	MAE
12	0.065 ± 0.000	0.169 ± 0.000	0.077 ± 0.003	0.181 ± 0.001	0.060 ± 0.000	0.158 ± 0.000	0.075 ± 0.000	0.174 ± 0.000
24	0.086 ± 0.000	0.196 ± 0.000	0.095 ± 0.001	0.204 ± 0.000	0.079 ± 0.000	0.181 ± 0.000	0.106 ± 0.000	0.206 ± 0.000
48	0.121 ± 0.001	0.234 ± 0.001	0.120 ± 0.002	0.231 ± 0.001	0.104 ± 0.001	0.209 ± 0.001	0.167 ± 0.001	0.258 ± 0.001
96	0.165 ± 0.000	0.276 ± 0.001	0.150 ± 0.001	0.261 ± 0.000	0.132 ± 0.001	0.236 ± 0.001	0.184 ± 0.002	0.251 ± 0.001
	=		· · · ± • · • • ·	±	1 ±		1 ±	10000

Ε FULL RESULTS FOR ABLATION STUDIES

To elaborate on the effectiveness of our TwinsFormer, we conduct detailed ablations covering disabling (X), swapping (swap), and replacing components (+). Due to the page limitation, we provide detailed results and analysis here for ablation studies.

As shown in Table 8, we disable the decomposition and use the observed series as input for ①, where the results indicate that decomposed components are more effective for model performance than the

810 observed values as inputs. For 2, we swap the positions of the seasonal and trend components (i.e., 811 feeding the trend components to the attention module while forcing the seasonal components as the 812 input for the interactive module.), and the results show that the multivariate correlations captured by 813 the attention module from the trend components are not effective as those obtained from the seasonal 814 components. For 3, we replace the subtraction mechanism with the addition operation, in which the performance is inferior to TwinsFormer. This highlights that residual learning is more in line with 815 the idea of decomposition design. As for (3), (4), (5), and (6), we disable E'_{T} , A_{S} , F_{S} , and σ variants, 816 respectively. The performance of the four cases is worse than that of the full model, indicating that 817 each component is effective for the interactive module to capture inherent dependencies. 818

Table 9: Full ablation results on four benchmarks. We disable or replace each component of both decomposition and interactions for TwinsFormer. \checkmark and \varkappa indicate with and without certain components, respectively. H1, H2, H3, and H4 denote different prediction lengths, where prediction length belongs to $S \in \{96, 192, 336, 720\}$ for ECL and Traffic, while $S \in \{12, 24, 48, 96\}$ for PEMS03 and PEMS07. Avg means the average results of all predicted lengths.

825	Design	Decomposition	Interactions				Prediction	Prediction ECL		Traffic		PEMS03		PEMS07		
826			-	E_T'	A_S	F_S	σ	lengths	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
827 828 829	TwinsFormer	✓	1	1	1	1	1	H1 H2 H3 H4	0.139 0.158 0.172 0.200	0.233 0.252 0.267 0.293	0.382 0.392 0.410 0.442	0.260 0.267 0.276 0.292	0.065 0.086 0.121 0.165	0.169 0.196 0.234 0.276	0.060 0.079 0.104 0.132	0.158 0.181 0.209 0.236
830								Avg	0.167	0.262	0.406	0.273	0.109	0.219	0.094	0.196
831 832 833	(I)	×	1	1	1	1	1	H1 H2 H3 H4	0.142 0.160 0.178 0.224	0.245 0.258 0.272 0.312	0.388 0.405 0.420 0.456	0.265 0.278 0.284 0.301	0.070 0.093 0.133 0.169	0.176 0.204 0.242 0.283	0.065 0.083 0.112 0.146	0.162 0.185 0.214 0.254
034								Avg	0.176	0.272	0.417	0.282	0.116	0.226	0.102	0.204
535 536 537 538 539 540 541 542 543 544 545 544 545 544 545 545 550 551 552 553 554	2	swap	~	1	1	1	1	H1 H2 H3 H4	0.140 0.161 0.175 0.210	0.238 0.256 0.268 0.298	0.389 0.399 0.414 0.448	0.263 0.271 0.279 0.296	0.069 0.092 0.128 0.167	0.175 0.202 0.239 0.281	0.061 0.082 0.115 0.145	0.160 0.187 0.217 0.251
								Avg	0.172	0.265	0.413	0.277	0.114	0.224	0.101	0.204
	3	✓	+	1	1	1	1	H1 H2 H3 H4	0.145 0.168 0.182 0.223	0.248 0.263 0.276 0.311	0.388 0.402 0.414 0.458	0.268 0.271 0.289 0.304	0.072 0.095 0.135 0.171	0.177 0.204 0.243 0.286	0.063 0.081 0.115 0.147	0.161 0.186 0.224 0.255
								Avg	0.180	0.275	0.416	0.283	0.118	0.228	0.102	0.207
	۹	1	1	×	1	1	1	H1 H2 H3 H4	0.149 0.172 0.188 0.229	0.248 0.265 0.282 0.318	0.392 0.407 0.419 0.452	0.269 0.275 0.284 0.303	0.074 0.099 0.138 0.176	0.178 0.207 0.251 0.290	0.064 0.083 0.124 0.149	0.163 0.189 0.231 0.257
								Avg	0.185	0.278	0.418	0.283	0.122	0.232	0.122	0.210
	5	1	1	1	x	1	1	H1 H2 H3 H4	0.148 0.172 0.188 0.222	0.251 0.263 0.281 0.311	0.389 0.399 0.413 0.449	0.264 0.270 0.278 0.300	0.073 0.095 0.131 0.173	0.179 0.205 0.247 0.284	0.065 0.084 0.119 0.143	0.163 0.192 0.225 0.251
								Avg	0.183	0.277	0.413	0.278	0.118	0.229	0.103	0.208
	6	<i>✓</i>	1	1	1	×	1	H1 H2 H3 H4	0.142 0.163 0.184 0.215	0.243 0.257 0.279 0.303	0.385 0.396 0.417 0.451	0.265 0.273 0.285 0.302	0.072 0.096 0.139 0.177	0.175 0.203 0.244 0.288	0.066 0.084 0.120 0.147	0.165 0.194 0.223 0.259
855								Avg	0.176	0.271	0.412	0.281	0.121	0.228	0.104	0.210
856 857 858	Ø	1	~	1	1	1	x	H1 H2 H3 H4	0.141 0.160 0.178 0.223	0.236 0.255 0.270 0.310	0.388 0.404 0.414 0.445	0.268 0.271 0.279 0.294	0.069 0.094 0.132 0.178	0.174 0.204 0.245 0.289	0.079 0.084 0.116 0.147	0.193 0.190 0.223 0.254
859								Avg	0.176	0.268	0.413	0.278	0.118	0.228	0.107	0.215



66					6			1 4		1 171 (
67		Models				(2021)		Autoformer (2021)		(2022)		(20	former 23)
68		Metrics		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
59			96	0.260	0.358	0.274	0.368	0.201	0.317	0.215	0.320	0.180	0.294
70		Original	192	0.266	0.367	0.296	0.386	0.222	0.334	0.259	0.355	0.193	0.307
ч		Original	720	0.280	0.375	0.300	0.394	0.231	0.338	0.296	0.383	0.208	0.321
1	ECL		Avg	0.277	0.372	0.311	0.397	0.227	0.338	0.267	0.359	0.219	0.328
2		+Twins	96	0 139	0.233	0.138	0.232	0 151	0.242	0 140	0.234	0 142	0 235
3			192	0.158	0.252	0.156	0.248	0.161	0.253	0.159	0.252	0.157	0.249
4			336	0.172	0.267	0.170	0.264	0.176	0.269	0.171	0.267	0.173	0.267
5			Avg	0.167	0.262	0.168	0.261	0.176	0.267	0.168	0.262	0.170	0.262
6			96	0.647	0.357	0.719	0.391	0.613	0.388	0.691	0.393	0.562	0.343
7		Original	192	0.649	0.356	0.696	0.397	0.616	0.382	0.729	0.419	0.587	0.356
8			336	0.667	0.364	0.777	0.420	0.622	0.337	0.756	0.423	0.612	0.370
9	Traffic		1 <u>720</u>	0.697	0.370	0.804	0.472	0.628	0.408	0.825	0.449	0.672	0.423
0	Thund	1	Avg	0.005	0.303	0.764	0.410	0.028	0.379	0.750	0.421	0.008	0.575
			96 192	0.382	0.260	0.397	0.266	0.400	0.268	0.392	0.264	0.408	0.273
1		+Twins	336	0.410	0.276	0.436	0.284	0.439	0.284	0.428	0.281	0.443	0.288
2			720	0.442	0.292	0.471	0.303	0.473	0.305	0.465	0.301	0.477	0.307
3			Avg	0.406	0.273	0.431	0.282	0.433	0.288	0.424	0.280	0.439	0.287
4			96	0.105	0.205	0.202	0.293	0.272	0.385	0.105	0.207	0.128	0.257
5		Original	336	0.121	0.222	0.175	0.277	1.032	0.440	0.127	0.229	0.173	0.306
c			720	0.175	0.274	0.211	0.304	1.031	0.796	0.173	0.287	0.490	0.524
7	PEMS03		Avg	0.137	0.237	0.193	0.290	0.667	0.601	0.140	0.245	0.265	0.368
1		+Twins	96	0.065	0.169	0.064	0.167	0.066	0.171	0.065	0.169	0.065	0.169
8			192	0.086	0.196	0.084	0.193	0.087	0.196	0.085	0.195	0.086	0.195
9			336	0.121	0.234	0.115	0.228	0.122	0.235	0.120	0.233	0.121	0.234
)		l I	<u>720</u>	0.100	0.210	0.105	0.209	0.105	0.277	0.104	0.217	0.110	0.280
1			Avg	0.105	0.219	0.105	0.214	0.110	0.220	0.105	0.219	0.111	0.220
-		Original	96	0.173	0.235	0.189	0.255	0.199	0.336	0.174	0.236	0.117	0.238
2			336	0.181	0.245	0.196	0.261	0.390	0.470	0.183	0.244	0.226	0.353
3			720	0.185	0.252	0.196	0.260	0.554	0.578	0.180	0.245	0.309	0.407
4	PEMS07		Avg	0.178	0.243	0.194	0.259	0.367	0.451	0.178	0.240	0.200	0.318
5			96	0.060	0.158	0.060	0.157	0.061	0.160	0.060	0.158	0.061	0.159
6		+Twins	336	0.079	0.181	0.103	0.181	0.082	0.180	0.079	0.181	0.080	0.184
7			720	0.132	0.236	0.128	0.232	0.134	0.239	0.140	0.245	0.140	0.245
8			Avg	0.094	0.196	0.092	0.194	0.096	0.200	0.096	0.198	0.098	0.201
9													
0	Lear	ning Rat	e		🗔	B	lock Nu	mber	1]		lidden	Dimen
1 0	0.4			•	0.4			-	-	0.4			
2													
3	0.3			H H	ן 0.3					표 0.3			
4 X				2						MS			
5 0					0.2					0.2			
- U					··			-		0.2			
-				-	-		<u> </u>	- i	—		—		
17	0.001 0.000	05 0.000	3 0.0	001	1		2	3	4		256	512	102
18				ETT		ECL	-	Traffic	· –	- Wea	ather		

Table 10: Full results of promotion and attention compatibility on five Transformer-based models
 for both long-term and short-term forecasting tasks.



913 914

915

909

F EXTRA RESULTS FOR MODEL ANALYSIS

We apply the proposed Twinsblocks to Transformer and its variants and compare the performance with the original results. Specifically, we regard Transformer (Vaswani et al., 2017), Informer (Zhou



Figure 8: Performance of generation on unseen variates, comparing with iTransformer on PEMS07 and Traffic datasets. We train the model with $\{20\%, 40\%, 60\%, 80\%, 100\%\}$ variates, and evaluate the performance by forecasting all variates.

et al., 2021), Autoformer (Wu et al., 2021), Flowformer (Wu et al., 2022), and Periodformer (Liang
et al., 2023) as original models, and treat various combinations of attention mechanisms and our
interactive strategy as Twins variants. As shown in Table 8, our framework has robust performance
under various attention mechanisms for the same benchmark and obtains better forecasting accuracy
than original Transformer-based variants when using the same attention mechanism.

We evaluate the hyperparameter sensitivity of TwinsFormer in terms of the learning rate, the number
of Twinsblock, and the hidden dimension of variate tokens. As shown in Figure 7, The performance
fluctuates under different hyperparameter settings. We can observe that the learning rate, as the most
common hyperparameter, should be carefully selected for different datasets. In most cases, increasing the number of the Twinsblock tends to strengthen the model performance, especially in datasets
with numerous varieties. For scenarios involving many attributes, the forecasting performance will
decrease when the hidden dimension of variate tokens is larger than 1024.

As seen in Figure 8, we set 5 variate proportions to explore the model generalization performance
for short-term and long-term forecasting. From the figure, we can observe that the prediction performance is better as the number of variate increases for both forecasting tasks. Comparatively,
the performance of TwinsFormer is better than that of iTransformer under any variate proportion,
which demonstrates that our method has a more powerful generalization ability than iTransformer.
Notably, TwinsFormer can be naturally trained with 20% variates and learn transferable representations to achieve favorable forecasting on all varieties.

958 G SHOWCASES

959 G.1 VISUALIZATION OF REPRESENTATIONS

960

957

935

936

937

To better understand the representations learned by our model, we visualize the learned multivariate 961 correlations by the pre-Softmax scores. Following (Liu et al., 2024), we normalize each variate token 962 on its feature dimension and reveal the variate-wise correlation by the whole score map $\mathbf{A} \in \mathbb{R}^{N \times N}$ 963 among N paired variate tokens. Meanwhile, we visualize the trend, seasonal, and weighted repre-964 sentations by normalized operation. As shown in Figure 9, we visualize the learned representations 965 of iTranformer (Liu et al., 2024) for comparison to illustrate the superiority of our method. Based on 966 Figure 9, we can observe that the multivariate correlations learned by TwinsFormer are clearer than 967 those of iTransformer. Consequently, the learned representation via seasonal and trend branches 968 can capture more effective information than iTranformer without decomposition design. Taking the representation in Figure 9 for example, (b), (c) and (d) contain more obvious information in the gold 969 dashed box, and (a) has no variations within the red dashed boxes in the corresponding positions of 970 (b), (c) and (d). Such differences show that TwinsFormer can better capture inherent dependencies 971 of time series than iTransformer for forecasting tasks.

990

991

992

1008



Figure 9: Analysis of multivariate correlations and series representations, comparing with iTransformer on ECL dataset. The left part of the dashed line denotes the multivariate correlations learned by iTransformer and TwinsFormer. Full series representations are obtained using iTransformer, and the rest are extracted from TwinsFormer. Zoom in for more details.



Figure 10: Analysis of multivariate correlations and series representations in different TwinsBlock
 on ECL dataset. The attention map and decomposed components of Layer 1 and n are obviously different, which indicates that our dual-stream framework can learn the interactions between seasonal and trend components, and thus better capture the multivariate dependencies.

1013 In Figure 10, we provide the visualization of the attention maps and the decomposed components 1014 in different layers. Naturally, the inputs fed to the first TwinsBlock are the initial decomposed 1015 components obtained by the moving average kernel, while the inputs to the last TwinsBlock are 1016 decomposed components learned by our dual-stream interaction structure. Accordingly, the atten-1017 tion map of Layer 1 captures the multivariate correlations among the initial seasonal components, while the attention map of Layer n obtains the multivariate correlations among the learned seasonal 1018 components. The attention map and decomposed components of Layer 1 and n are obviously dif-1019 ferent, which indicates that our dual-stream framework can learn the interactions between seasonal 1020 and trend components, and thus better capture the multivariate dependencies. 1021

In Figure 11-12, we present a detailed decomposed comparison of the moving average kernel and
 our interactions. Note that, the decomposed components by the moving average kernel can be regarded as initial decomposed components, while the decomposed components by our interactions are learned from our TwinsFormer. From these two figures, we can observe that the learned components are different from the initial components, where especially the learned seasonal component





Figure 12: Trend-Seasonal Decomposition Results obtained by the moving average kernel and our interactions on Traffic. To better compare the variations of raw time series with the seasonal com-ponents, we translate the seasonal components on the vertical axis.





Figure 14: Traffic prediction cases among different models under the input-96-predict-96 setting.

1121 1122

1123

G.2 VISUALIZATION OF PREDICTION RESULTS

For clarity and comparison among different models, we present supplementary forecasting showcases on four representative benchmarks in Figure 13, 14, 15, and 16. To be concrete, we provide prediction cases for TwinsFormer, iTransformer (Liu et al., 2024), TimeMixer (Wang et al., 2024), PatchTST (Nie et al., 2023), Crossformer (Zhang & Yan, 2023), and DLinear (Zeng et al., 2023) over ECL, Traffic, Weather, and PEMS07 datasets. Among these models, TwinsFormer exhibits superior forecasting performance with the most precise future series variations. Note that, all the qualitative results of different models are obtained with a fixed input length T = 96 and forecasting horizon S = 96 over four datasets.

- 1131
- 1132



Figure 16: PEMS07 prediction cases among different models under the input-96-predict-96 setting.
1185
1186