# TIMEBRIDGE: NON-STATIONARITY MATTERS FOR LONG-TERM TIME SERIES FORECASTING

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

030

Paper under double-blind review

## Abstract

Non-stationarity poses significant challenges for multivariate time series forecasting due to the inherent short-term fluctuations and long-term trends that can lead to spurious regressions or obscure essential long-term relationships. Most existing methods either eliminate or retain non-stationarity without adequately addressing its distinct impacts on short-term and long-term modeling. Eliminating non-stationarity is essential for avoiding spurious regressions and capturing local dependencies in short-term modeling, while preserving it is crucial for revealing long-term cointegration across variates. In this paper, we propose TimeBridge, a novel framework designed to bridge the gap between non-stationarity and dependency modeling in long-term time series forecasting. By segmenting input series into smaller patches, TimeBridge applies Integrated Attention to mitigate shortterm non-stationarity and capture stable dependencies within each variate, while Cointegrated Attention preserves non-stationarity to model long-term cointegration across variates. Extensive experiments show that TimeBridge consistently achieves state-of-the-art performance in both short-term and long-term forecasting. Additionally, TimeBridge demonstrates exceptional performance in financial forecasting on the CSI 500 and S&P 500 indices, further validating its robustness and effectiveness. The code is available in the supplementary material.

## 1 INTRODUCTION

Multivariate time series forecasting aims to predict future changes based on historical observations of time series data, which holds significant applications in fields such as financial investment (Sezer et al., 2020), weather forecasting (Karevan & Suykens, 2020), and traffic flow prediction (Shu et al., 2021). However, the inherent non-stationarity of time series (Kim et al., 2022), characterized by short-term fluctuations and long-term trends, introduces challenges such as spurious regressions, making time series forecasting a particularly complex task.

Recently, many methods have emerged to utilize a normalization-and-denormalization paradigm to address non-stationarity in time series (Kim et al., 2022; Fan et al., 2023; Liu et al., 2023; 2024b). For instance, RevIN (Kim et al., 2022) normalizes the input data and subsequently applies its dis-040 tributional characteristics to denormalize the output predictions. Building on this approach, other 041 methods have designed distributional prediction networks (Fan et al., 2023) and more refined nor-042 malization techniques (Liu et al., 2023) to further mitigate non-stationarity. On the other hand, some 043 studies (Liu et al., 2022b; Ma et al., 2024; Fan et al., 2024) argue that over-stabilizing time series 044 may actually reduce the richness of embedded features, leading to a decline in model performance. Existing methods for addressing non-stationarity in time series face a dilemma: some prioritize eliminating non-stationary factors to reduce overfitting, while others attempt to incorporate these factors 046 but lack comprehensive theoretical frameworks. Furthermore, they do not adequately explain the 047 trade-offs between removing non-stationarity and leveraging it for modeling. 048

It is notable that non-stationary characteristics have distinct impacts on modeling short-term and
 long-term dependencies. Non-stationarity can lead to spurious regressions in short-term modeling
 due to the high randomness and unpredictability of short-term fluctuations (Noriega & Ventosa Santaulària, 2007) (see Fig. 1b). For example, a sudden drop in temperature could be caused by
 a typhoon or cold front, both of which have no intrinsic connection. Retaining non-stationarity is crucial



Figure 1: Visualization of the impact of non-stationarity on short-term and long-term modeling. (a) 069 Comparison of two cointegrated series  $X_t$  and  $Y_t$  before and after de-trending, with  $X_t$  showing two 070 random fluctuations in Phase A. (b) Short-term patch attention map of  $X_t$  in Phase A, showing that 071 non-stationarity leads to spurious regression in short-term modeling. (c) Long-term relationship be-072 tween  $X_t$  and  $Y_t$ , where removing non-stationarity eliminates the cointegration. (d) Non-stationarity 073 enables modeling long-term cointegration between variates but causes spurious regressions in short-074 term modeling (orange line), while de-trending benefits short-term modeling but disrupts long-term 075 relationships (green line).

for capturing long-term cointegration relationships among variates, reflecting their co-movement or 078 synchronized changes over time (Fanchon & Wendel, 1992). Removing non-stationarity may also 079 eliminate these essential long-term dependencies (see Fig. 1c). Therefore, while non-stationarity can cause spurious regressions in short-term modeling, it is essential for modeling long-term depen-081 dencies between variates. Conversely, eliminating non-stationarity benefits short-term modeling but 082 erases long-term cointegration (see Fig. 1d). 083

In this paper, to address the dual challenges posed by non-stationarity in short-term and long-term 084 modeling, we propose distinct strategies tailored to each scenario. For short-term modeling, we 085 eliminate non-stationarity to capture the strong temporal dependencies within each variate, as shortterm causal relationships exist mainly between consecutive time points within a single variate rather 087 than across variates. This strategy reduces the risk of spurious regressions from non-stationary fluc-088 tuations, enabling the model to better capture the local causal dynamics. For long-term modeling, 089 we utilize the preserved non-stationarity to uncover long-term cointegration relationships between 090 different variates, thereby enabling more accurate and reliable long-term forecasting. 091

Technically, based on the above motivations, we propose TimeBridge as a novel framework to *bridge* 092 the gap between non-stationarity and dependency modeling in long-term time series forecasting. 093 TimeBridge first captures short-term fluctuations by partitioning the input sequence into small-094 length patches, followed by utilizing Integrated Attention to model these stabilized sub-sequences 095 within each variate. Here, "Integrated" reflects the non-stationary nature of the short-term series, 096 also referred to as integrated series (Park & Phillips, 2001). Subsequently, we downsample the patches to reduce their quantity, thereby enriching each patch with more long-term information. 098 Cointegrated Attention retains the non-stationary characteristics of the sequences to effectively cap-099 ture the long-term cointegration relationships among variates. Experiments across multiple datasets 100 demonstrate that TimeBridge achieves consistent state-of-the-art performance in both long-term and short-term forecasting. Furthermore, we validate the effectiveness of TimeBridge on two financial 101 datasets, the CSI 500 and S&P 500, which exhibit significant short-term volatility and strong long-102 term cointegration relationships among sectors. 103

- 104 In a nutshell, our contributions are summarized in three folds:
- 105

1. Going beyond previous methods, we establish a novel connection between non-stationarity and dependency modeling, highlighting the importance of eliminating non-stationarity in 107 short-term variations while preserving it for long-term cointegration.

- 2. We propose TimeBridge, a novel framework that employs Integrated Attention to model temporal dependencies by mitigating short-term non-stationarity, and Cointegrated Attention to capture long-term cointegration across variates while retaining non-stationarity.
- 3. Comprehensive experiments demonstrate that TimeBridge achieves state-of-the-art performance in both long-term and short-term forecasting across various datasets. Moreover, we further validate its robustness and effectiveness on the CSI 500 and S&P 500 indices, which pose additional challenges due to their complex volatility and cointegration characteristics.
- 114 115 116 117

109

110

111

112

113

## 2 RELATED WORKS

As shown in Fig. 2, recent advancements in multivariate time series forecasting have predominantly focused on two core directions: Normalization and Dependecy Modeling.

121 **Normalization** can be divided into stationary and non-stationary methods. Stationary methods (Kim et al., 2022; Fan et al., 2023; Liu et al., 2024b; 2023) aim to eliminate non-stationarity through 122 model-agnostic normalization techniques, thereby preventing spurious regressions and enhancing 123 model performance. For example, RevIN (Kim et al., 2022) applies Z-normalization to the input 124 sequence and then reverses the normalization on the output using the distributional characteristics of 125 the input, assuming that both share similar distributional properties. Dish-TS (Fan et al., 2023) takes 126 this further by predicting the statistical characteristics of the output with a distribution prediction 127 model. Additionally, SAN (Liu et al., 2023) offers a more granular patch-level prediction method. 128 Conversely, some approaches (Liu et al., 2022b; Ma et al., 2024; Fan et al., 2024) advocate preserv-129 ing non-stationarity, as excessive normalization can eliminate diverse sequence characteristics and 130 limit predictive accuracy.

131 **Dependency Modeling** focuses on designing methods to capture the relationships within multivari-132 ate time series, which can be classified into Channel Independent (CI) and Channel Dependent (CD) 133 methods. CI methods (Zeng et al., 2023; Das et al., 2023; Nie et al., 2023; Dai et al., 2024; Lin et al., 134 2024) rely exclusively on the historical values of each individual channel for prediction, deliberately 135 avoiding cross-channel interactions. This strategy not only stabilizes the training process but also 136 excels at capturing rapid temporal dynamics unique to each variate. In contrast, CD methods (Wu et al., 2021; Zhou et al., 2022; Wu et al., 2023; Zhang & Yan, 2023; Liu et al., 2024a) leverage the 137 138 interrelationships between variates for modeling. While these methods utilize more information, they struggle with spurious regressions when modeling short-term dependencies, failing to capture 139 rapid changes effectively. 140

The challenge with previous methods lies in their isolated treatment of non-stationarity and dependency modeling, overlooking their intrinsic connection. Due to non-stationarity, time series often exhibit significant short-term fluctuations, leading to severe spurious regressions when modeling short-term dependencies. However, capturing long-term cointegration requires preserving this underlying variability. Therefore, short-term random fluctuations need to be addressed by elimi-



161 Eigur

146 147 148

149

150

151

152

153 154

156

157

Figure 2: Time series forecasting methods categorized by normalization and dependency modeling.

162 nating non-stationarity and modeling intra-variate temporal dependencies, while long-term cointe-163 gration demands preserving non-stationarity for inter-variate modeling. Our proposed TimeBridge 164 addresses these issues by employing Integrated Attention and Cointegrated Attention, respectively. 165

#### 3 METHOD

166

167 168

170

171

172

173

174

191

192

193

194

195 196 197

198

In the task of multivariate time series forecasting, the objective is to predict future sequences  $\mathbf{Y} = [\mathbf{x}_{I+1}, \cdots, \mathbf{x}_{I+O}] \in \mathbb{R}^{C \times O}$  given historical input sequences  $\mathbf{X} = [\mathbf{x}_1, \cdots, \mathbf{x}_I] \in \mathbb{R}^{C \times I}$ . Here, I and O denote the lengths of the input and output sequences, respectively, and C represents the number of time variates. It is important to recognize that real-world time series data often exhibit high short-term uncertainty, while long-term dynamics may reveal cointegration relationships among different time variates.



Figure 3: Overall architecture of TimeBridge: (a) Patch Embedding divides the input sequence into non-overlapping patches and embeds each as a token; (b) Integrated Attention models temporal dependencies within each variate by mitigating short-term non-stationarity; (c) Patch Downsampling reduce patches to aggregates long-term information and lower complexity; (d) Cointegrated Attention captures long-term relationships across variates while keeping non-stationarity.

#### 3.1 STRUCTURE OVERVIEW

199 As illustrated in Fig. 3, our proposed TimeBridge consists of four key components: (a) Patch Embedding segments the input sequence into non-overlapping patches and transforms each patch into a 200 patch token; (b) Integrated Attention models the dependencies among all patch tokens of the same 201 variates. By eliminating non-stationarity within each patch token, it mitigates the risk of spurious 202 regressions that could arise from abrupt short-term changes; (c) Patch Downsampling aggregates 203 global information and reduces the number of patches to encapsulate richer long-term features within 204 each patch, while simultaneously lowering computational complexity; (d) Cointegrated Attention 205 preserves the non-stationary characteristics of the sequence and models the long-term cointegration 206 relationships across different variates within the same temporal window. 207

208 3.2 PATCH EMBEDDING 209

210 In this stage, each variate of the input sequence X is first divided into non-overlapping patches, 211 and each patch is then mapped to an embedded patch token. Since the process is identical for each 212 variate, we use  $\mathbf{X}$  to represent a single variate and later restore the dimensionality of the variates. 213 Formally, this process can formulated as follows:

$$\{\mathbf{p}_1, \cdots, \mathbf{p}_N\} = \text{Patching}(\mathbf{X}), \quad \mathbf{P} = \text{Embedding}(\mathbf{p}_1, \cdots, \mathbf{p}_N)$$
 (1)

Here, each patch  $\mathbf{p}_i$  has a length of S, and the number of patches  $N = \lfloor \frac{I}{S} \rfloor$ . The Embedding(·) operation transforms each patch from its original length S to a hidden dimension D through a trainable linear layer. This results in embedded patch tokens  $\mathbf{P} \in \mathbb{R}^{C \times N \times D}$ , where each of the C variates contains N patches, capturing local information that is typically subject to rapid shortterm fluctuations. For convenience, we denote  $\mathbf{P}_{c,:}$  as the set of all patches within a single variate and  $\mathbf{P}_{:,n}$  as the patches across all variates at the same time position in the following sections.

## 223 3.3 INTEGRATED ATTENTION

The embedded patch tokens **P** represent short-term non-stationary sequences, also referred to as integrated series of order k (k > 0) (Park & Phillips, 2001; Mushtaq, 2011). This non-stationarity makes it challenging to model dependencies across different variates, as short-term fluctuations are highly susceptible to external shocks. Furthermore, modeling temporal dependencies within the same variate can lead to spurious regression due to the inherent non-stationarity of the patches. To address this, we first apply a patch-wise normalization to all patches within a variate:

$$\mathbf{p}_{i}^{\text{Trend}} = \text{AvgPool}(\text{Padding}(\mathbf{p}_{i})), \quad \mathbf{p}_{i}' = \mathbf{p}_{i} - \mathbf{p}_{i}^{\text{Trend}}, \quad \mathbf{P}_{c,:}' = \{\mathbf{p}_{1}', \cdots, \mathbf{p}_{N}'\},$$
(2)

where the AvgPool( $\cdot$ ) operation is moving average with the Padding( $\cdot$ ) operation to keep the series length unchanged. We then employ the proposed Integrated Attention mechanism to capture temporal dependencies within the same variate:

 $\hat{\mathbf{P}}_{c,:} = \text{LayerNorm} \left( \mathbf{P}_{c,:} + \text{Attention} (\mathbf{P}'_{c,:}, \mathbf{P}'_{c,:}, \mathbf{P}_{c,:}) \right),$ (3)

$$\mathbf{P}_{c,:} = \text{LayerNorm}\left(\hat{\mathbf{P}}_{c,:} + \text{MLP}(\hat{\mathbf{P}}_{c,:})\right),\tag{4}$$

where MLP( $\cdot$ ) represents a multi-layer feedforward network, and LayerNorm( $\cdot$ ) denotes layer normalization. The attention mechanism uses the normalized P'<sub>c,:</sub> as both Query and Key, while the original P<sub>c,:</sub> serves as the Value. This design generates a stationary attention map, which is then directly multiplied by the Value, removing the need for subsequent denormalization. By leveraging Integrated Attention in this way, we effectively model the temporal dependencies without being affected by the short-term non-stationary nature of the sequences.

## 249 250 3.4 PATCH DOWNSAMPLING

Long-term equilibrium relationships between sequences, or cointegration among different variates,
 often require sequences to contain sufficient long-term information to emerge. Therefore, before
 modeling the cointegration between variates, it is crucial to increase the amount of global informa tion represented by each patch. This is achieved by reducing the number of patches and aggregating
 global information through the attention mechanism:

256

222

224

231 232 233

237 238

257 258

$$\mathbf{P}_{c,:}' = \text{Downsample}(\mathbf{P}_{c,:}), \quad \mathbf{P}_{c,:} = \text{Attention}(\mathbf{P}_{c,:}', \mathbf{P}_{c,:}, \mathbf{P}_{c,:}). \tag{5}$$

Here, Downsample(·) reduces the N patches in  $\mathbf{P}_{c,:}$  to M patches (M < N) using an MLP. By employing the downsampled  $\mathbf{P}'_{c,:} \in \mathbb{R}^{M \times D}$  as the Query and the original  $\mathbf{P}_{c,:} \in \mathbb{R}^{N \times D}$  as the Key and Value in the attention mechanism, we leverage the long-range modeling capability of attention to dynamically aggregate global information. This allows each patch to encapsulate richer long-term information, making it possible to capture the intricate cointegration relationships that emerge only over sufficiently extended temporal horizons.

265

266 3.5 COINTEGRATED ATTENTION

267

Although short-term relationships between integrated series are susceptible to spurious regressions,
 accurately modeling long-term cointegration between sequences necessitates retaining their inher ent non-stationary characteristics. Since each downsampled patch now encapsulates more extensive

long-term information, we leverage Cointegrated Attention to directly model the cointegration relationships among all variates at the same time interval  $\mathbf{P}_{:,n} \in \mathbb{R}^{C \times D}$ :

$$\hat{\mathbf{P}}_{:,n} = \text{LayerNorm}(\mathbf{P}_{:,n} + \text{Attention}(\mathbf{P}_{:,n}, \mathbf{P}_{:,n}, \mathbf{P}_{:,n})),$$
(6)

$$\mathbf{P}_{:,n} = \text{LayerNorm}(\hat{\mathbf{P}}_{:,n} + \text{MLP}(\hat{\mathbf{P}}_{:,n})).$$
(7)

This attention mechanism not only captures the global cointegration relationships across variates but also adaptively assesses the strength of these relationships: stronger cointegration is reflected by higher attention weights, while weaker connections receive lower weights. Finally, the embedded patch tokens  $\mathbf{P} \in \mathbb{R}^{C \times M \times D}$  are unpatched and projected to the final output  $\mathbf{Y} \in \mathbb{R}^{C \times O}$ .

## 4 EXPERIMENT

285 286

273 274

275

276

277 278

279

280

To validate the effectiveness of the proposed TimeBridge, we conduct extensive experiments on a
 variety of time series forecasting tasks, including both long-term and short-term forecasting. Addi tionally, we evaluate TimeBridge on financial forecasting tasks characterized by significant short term volatility and strong long-term cointegration relationships among sectors.

291 **Baselines.** For long-term forecasting, we select a diverse set of state-of-the-art baselines represen-292 tative of recent advancements in time series forecasting, including the Transformer-based PDF (Dai 293 et al., 2024), the CNN-based ModernTCN (Donghao & Xue, 2024), the MLP-based TimeMixer 294 (Wang et al., 2024), as well as other competitive methods such as iTransformer (Liu et al., 2024a), PatchTST (Nie et al., 2023), Crossformer (Zhang & Yan, 2023), FEDformer (Zhou et al., 2022), 295 MICN (Wang et al., 2022), TimesNet (Wu et al., 2023), and DLinear (Zeng et al., 2023). For 296 short-term forecasting, we add a well-performing baseline SCINet (Liu et al., 2022a). For financial 297 forecasting, we also incorporate the momentum strategy CSM (Jegadeesh & Titman, 1993) and the 298 reversal strategy BLSW (Poterba & Summers, 1988), along with two classic deep learning models, 299 LSTM (Hochreiter & Schmidhuber, 1997) and Transformer (Vaswani et al., 2017a), to provide a 300 comprehensive evaluation. 301

**Implementation Details.** All experiments are implemented in PyTorch (Paszke et al., 2019) and conducted on two NVIDIA RTX 3090 24GB GPUs. We use the Adam optimizer (Kingma, 2014) with a learning rate selected from  $\{1e-3, 1e-4, 5e-4\}$ . The number of patches N is set to 30. For additional details on hyperparameters and settings, please refer to the Appendix D.

306 307

308

## 4.1 LONG-TERM FORECASTING

309 Setups. We conduct long-term forecasting experiments on eight widely-used real-world datasets, 310 including the Electricity Transformer Temperature (ETT) dataset with its four subsets (ETTh1, 311 ETTh2, ETTm1, ETTm2), as well as Weather, Electricity, Traffic, and Solar (Liu et al., 2024a). 312 These datasets exhibit strong non-stationary characteristics, detailed in Appendix C. Following pre-313 vious works (Zhou et al., 2021; Wu et al., 2021), we use Mean Square Error (MSE) and Mean 314 Absolute Error (MAE) as evaluation metrics. We set the input length I to 720 for our method. For 315 other baselines, we adopt two settings: one uses the original results from the baseline papers, and the other involves searching for the optimal input length I and other hyperparameters. Details of the 316 search process can be found in Appendix E.1. 317

Results. As shown in Tab. 1, TimeBridge consistently outperforms other methods on these non-stationary datasets, with an average improvement of over 10% compared to the baselines. Moreover, in experiments with different hyperparameter search settings in Tab. 2, TimeBridge continues to achieve the best overall results. Specifically, compared to the current state-of-the-art methods—Transformer-based PDF (Dai et al., 2024), CNN-based ModernTCN (Donghao & Xue, 2024), and MLP-based TimeMixer (Wang et al., 2024)—TimeBridge reduces MSE and MAE by 3.10%/1.64%, 3.55%/0.81%, and 6.92%/4.54%, respectively.

Models	TimeI (Ou	Bridge urs)	iTrans (20	former (24a)	PI (20	DF (24)	Time (20	Mixer (24)	Patcl (20	nTST (23)	Cross	former (23)	FEDf (20	ormer	Moder (20	mTCN 24)	MI (20	CN (22)	Time (20	esNet 23)	DL (2
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE
ETTm1	0.353	0.388	0.407	0.410	0.339	0.372	0.381	0.395	0.353	0.382	0.431	0.443	0.448	0.452	0.351	0.381	0.392	0.414	0.400	0.406	0.357
ETTm2	0.246	0.310	0.288	0.332	0.252	0.313	0.275	0.323	0.256	0.317	0.621	0.510	0.305	0.349	0.253	0.314	0.290	0.343	0.291	0.333	0.267
ETTh1	0.399	<u>0.420</u>	0.454	0.447	0.399	0.419	0.447	0.440	0.413	0.434	0.446	0.464	0.440	0.460	0.404	0.420	0.440	0.462	0.458	0.450	0.423
ETTh2	0.346	0.394	0.383	0.407	0.327	0.376	0.364	0.395	0.324	0.381	0.835	0.675	0.434	0.447	0.322	<u>0.379</u>	0.411	0.440	0.414	0.427	0.431
Weather	0.218	0.259	0.258	0.279	0.225	0.261	0.240	0.271	0.226	0.264	0.343	0.386	0.309	0.360	0.224	0.264	0.243	0.299	0.259	0.287	0.246
Electricity	0.149	0.246	0.178	0.270	0.159	0.250	0.182	0.272	0.159	0.253	0.293	0.351	0.205	0.315	<u>0.156</u>	0.253	0.187	0.295	0.192	0.295	0.166
Traffic	0.360	0.255	0.428	0.282	0.383	0.254	0.484	0.297	0.391	0.264	0.535	0.300	0.573	0.347	0.396	0.270	0.542	0.316	0.620	0.336	0.434
Solar <sup>†</sup>	0.181	0.239	0 233	0.262	0 205	0.265	0.216	0.280	0 207	0 294	0 204	0 248	0 296	0 407	0.228	0.282	0.247	0.296	0.319	0.348	0.329

Table 1: Long-term forecasting results from the original papers. All results are averaged across four different prediction lengths:  $O \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively. See Tab. 10 for full results.

Models	TimeI (Or	Bridge Irs)	iTrans (20	former 24a)	PI (20	DF 24)	Time (20	Mixer (24)	Pate (20	hTST )23)	Cross	former 123)	FEDf	ormer 122)	Moder (20	rnTCN 124)	MI (20	CN (22)	Time (20	esNet )23)	DLi (20	inear 023)
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	0.353	0.388	0.362	0.391	0.335	0.373	0.348	<u>0.375</u>	0.353	0.382	0.420	0.435	0.382	0.422	0.351	0.381	0.383	0.406	0.400	0.406	0.357	0.379
ETTm2	0.246	0.310	0.269	0.329	0.247	<u>0.311</u>	0.256	0.315	0.256	0.317	0.518	0.501	0.292	0.343	0.253	0.314	0.277	0.336	0.291	0.333	0.267	0.332
ETTh1	<u>0.399</u>	0.420	0.439	0.448	0.395	0.420	0.411	<u>0.423</u>	0.413	0.434	0.440	0.463	0.428	0.454	0.404	0.420	0.433	0.462	0.458	0.450	0.423	0.437
ETTh2	0.346	0.394	0.374	0.406	0.326	0.376	0.316	0.384	0.324	0.381	0.809	0.658	0.388	0.434	0.322	<u>0.379</u>	0.385	0.430	0.414	0.427	0.431	0.447
Weather	0.218	0.259	0.233	0.271	0.220	0.259	0.222	<u>0.262</u>	0.226	0.264	0.228	0.287	0.305	0.287	0.224	0.264	0.242	0.298	0.259	0.287	0.240	0.300
Electricity	0.149	0.246	0.164	0.261	0.156	<u>0.250</u>	0.156	0.246	0.159	0.253	0.181	0.279	0.205	0.315	0.156	0.253	0.182	0.292	0.192	0.295	0.166	0.264
Traffic	0.360	0.255	0.397	0.282	0.377	<u>0.256</u>	0.387	0.262	0.391	0.264	0.523	0.284	0.573	0.347	0.396	0.270	0.535	0.312	0.620	0.336	0.434	0.295
Solar	0.181	0.239	0.200	0.260	0.205	0.265	0.192	0.244	0.194	0.245	<u>0.191</u>	0.242	0.243	0.350	0.228	0.282	0.213	0.266	0.244	0.334	0.247	0.309

Table 2: Long-term forecasting hyperparameter search results. All results are averaged across four different prediction lengths:  $O \in \{96, 192, 336, 720\}$ . See Tab. 11 for full results.

## 4.2 SHORT-TERM FORECASTING

Setups. For short-term forecasting, we conduct experiments on the PeMS datasets (Wang et al., 2024), which capture complex spatiotemporal correlations among multiple variates across city-wide traffic networks. We use mean absolute error (MAE), mean absolute percentage error (MAPE), and root mean squared error (RMSE) as evaluation metrics. The input length I is set to 96 and the output length O to 12 for all baselines. Details of datasets and metrics are in Appendix C and Appendix B.2.

**Results.** As shown in Tab. 3, methods that perform well in long-term forecasting with channel-independent approaches, such as PatchTST (Nie et al., 2023) and DLinear (Zeng et al., 2023), suffer from significant performance degradation on the PeMS dataset due to its strong inter-variable de-pendencies. In contrast, TimeBridge demonstrates robust performance on this challenging task, outperforming even the recent state-of-the-art method TimeMixer (Wang et al., 2024), which highlights its effectiveness in capturing complex spatiotemporal relationships.

Models	TimeBridge (Ours)	TimeMixe (2024)	r SCINet (2022a)	Crossforme (2023)	r PatchTST (2023)	TimesNe (2023)	t MICN I (2022)	DLinear (2023)	FEDformer (2022)	r Stationary (2022b)	Autoforme (2021)	r Informe (2021)
PeMS03 MAE RMSE	14.52 14.21 23.10	$\frac{\underline{14.63}}{\underline{14.54}}$	15.97 15.89 25.20	15.64 15.74 25.56	18.95 17.29 30.15	16.41 15.17 26.72	15.71 15.67 24.55	19.70 18.35 32.35	19.00 18.57 30.05	17.64 17.56 28.37	18.08 18.75 27.82	19.19 19.58 32.70
PeMS04 MAE RMSE	<u>19.24</u> <b>12.42</b> <u>31.12</u>	<b>19.21</b> <u>12.53</u> <b>30.92</b>	20.35 12.84 32.31	20.38 12.84 32.41	24.86 16.65 40.46	21.63 13.15 34.90	21.62 13.53 34.39	24.62 16.12 39.51	26.51 16.76 41.81	22.34 14.85 35.47	25.00 16.70 38.02	22.05 14.88 36.20
PeMS07 MAPE RMSE	20.43 8.42 33.44	$\frac{20.57}{8.62}$ 33.59	22.79 9.41 35.61	22.54 9.38 35.49	27.87 12.69 42.56	25.12 10.60 40.71	22.28 9.57 35.40	28.65 12.15 45.02	27.92 12.29 42.29	26.02 11.75 42.34	26.92 11.83 40.60	27.26 11.63 45.81
PeMS08 MAPE RMSE	14.98 9.56 23.77	$\frac{15.22}{9.67}$ 24.26	17.38 10.80 27.34	17.56 10.92 27.21	20.35 13.15 31.04	19.01 11.83 30.65	17.76 10.76 27.26	20.26 12.09 32.38	20.56 12.41 32.97	19.29 12.21 38.62	20.47 12.27 31.52	20.96 13.20 30.61

Table 3: Short-term forecasting results in the PeMS datasets.

## 3783794.3 FINANCIAL FORECASTING

Setups. We conduct experiments on both the U.S. and Chinese stock markets, including the S&P 500 and CSI 500 indices. Stock price movements are influenced by various factors such as economic indicators, market sentiment, geopolitical events, and company-specific news, leading to high non-stationarity. We predict next-day returns using historical data and generate investment portfolios with a buy-hold-sell strategy (Sanderson & Lumpkin-Sowers, 2018). At day t + 1 open, traders sell day t stocks and buy top-ranked ones based on predicted returns. Following previous work (Lin et al., 2021; Li et al., 2024), we evaluate performance using Annual Return Ratio (ARR), Annual Volatility (AVol), Maximum Drawdown (MDD), Annual Sharpe Ratio (ASR), Calmar Ratio (CR), and Information Ratio (IR). Details of datasets and metrics are in Appendix C and Appendix B.3. 

Results. As shown in Tab. 4, due to the non-stationary dynamics and complex market dependencies, other baseline methods struggle to consistently and accurately identify broadly optimal portfolios across different markets. In contrast, TimeBridge consistently performs best in both financial markets, demonstrating its ability to capture short-term fluctuations within financial time series and long-term cointegration between sectors.

Models			CSI	500					S&P	500		
1100015	ARR↑	AVol↓	MDD↓	ASR↑	CR↑	IR↑	<b>ARR</b> ↑	AVol↓	$\text{MDD}{\downarrow}$	ASR↑	CR↑	IR↑
BLSW (1988)	0.110	0.227	-0.155	0.485	0.710	0.446	0.199	0.318	-0.223	0.626	0.892	0.774
LSTM (1997)	-0.008	0.159	-0.172	-0.047	-0.044	-0.128	0.077	0.162	-0.139	0.370	0.712	0.929
Transformer (2017a)	0.154	0.156	-0.135	0.986	1.143	0.867	0.135	0.159	-0.140	0.852	0.968	0.908
PatchTST (2023)	0.118	0.152	<u>-0.127</u>	0.776	0.923	0.735	0.146	0.167	-0.140	0.877	1.042	0.949
Crossformer (2023)	-0.039	0.163	-0.217	-0.238	-0.179	-0.350	<u>0.284</u>	0.159	-0.114	<u>1.786</u>	2.491	<u>1.646</u>
iTransformer (2024a)	0.214	0.168	-0.164	1.276	1.309	1.173	0.159	0.170	-0.139	0.941	1.150	0.955
TimeMixer (2024)	0.078	<u>0.153</u>	-0.114	0.511	0.685	0.385	0.254	<u>0.162</u>	<u>-0.131</u>	1.568	1.938	1.448
TimeBridge (Ours)	0.285	0.203	-0.196	1.405	1.453	1.317	0.326	0.169	-0.142	1.927	<u>2.298</u>	1.842

Table 4: Results for financial time series forecasting in CSI 500 and S&P 500 datasets. See Tab. 12 for full results.

## 5 ABLATION STUDIES

To validate the effectiveness of the proposed TimeBridge, we conduct a comprehensive ablation
study on its architectural design. In Tab. 5, Tab. 6, and Tab. 7, the rows highlighted in gray correspond to the original TimeBridge configuration, serving as a baseline for comparison with various modified versions.

Ablation on removing or keeping non-stationarity. We conduct the following experiments: ① Non-stationarity retained in both Integrated and Cointegrated Attention. 2 Retained in Integrated, removed from Cointegrated. 3 Removed from Integrated, retained in Cointegrated. 4 Removed from both. Results in Tab. 5 show that the best performance is achieved when non-stationarity is removed in Integrated Attention, which models short-term intra-variate fluctuations, and retained in Cointegrated Attention, which captures long-term inter-variate dependencies. Conversely, retaining non-stationarity in Integrated Attention while removing it from Cointegrated Attention yields the worst results. 

Case	Integrated Attention	Cointegrated Attention	Wea	ather	So	lar	Elect	ricity	Tra	ffic
cuse	+ Norm?	+ Norm?	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
1	×	×	0.220	<u>0.260</u>	<u>0.183</u>	<u>0.242</u>	<u>0.153</u>	<u>0.249</u>	<u>0.371</u>	<u>0.260</u>
2	×	√	0.220	<u>0.260</u>	<u>0.183</u>	0.252	0.155	0.251	0.381	0.263
3	$\checkmark$	×	0.218	0.259	0.181	0.239	0.149	0.246	0.360	0.255
4	$\checkmark$	√	<u>0.219</u>	0.259	<u>0.183</u>	<u>0.242</u>	<u>0.153</u>	0.250	0.374	0.289

Table 5: Ablation on the effect of removing non-stationarity in Integrated Attention and Cointegrated Attention.  $\checkmark$  indicates the use of patch normalization to eliminate non-stationarity, while  $\times$  means non-stationarity is retained.

432 Ablation on Integrated and Cointegrated Attention impact and order. We conduct the following 433 experiments: 1 Integrated Attention only, 2 Cointegrated Attention only, 3 Integrated Attention 434 followed by Cointegrated Attention, and ( Cointegrated Attention followed by Integrated Attention, 435 with patch downsampling replaced by upsampling in this case. The results in Tab. 6 show that both 436 1 and 2 underperform compared to 3, indicating that both components are beneficial. Additionally, Ishows the weakest performance, possibly because modeling long-term cointegrated relationships 437 first leads to a loss of important short-term temporal features. 438

439											
440	Case	Integrated Attention	Cointegrated Attention	Wea	ther	So	lar	Elect	ricity	Tra	ffic
441	cuse	Order	Order	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
442	1	1	×	<u>0.220</u>	<u>0.262</u>	<u>0.184</u>	<u>0.244</u>	<u>0.158</u>	<u>0.252</u>	0.388	<u>0.264</u>
443	2	×	1	0.222	0.264	0.191	0.260	0.165	0.263	<u>0.369</u>	0.265
444	3	1	2	0.218	0.259	0.181	0.239	0.149	0.246	0.360	0.255
445	4	2	1	0.227	0.266	0.190	0.252	0.160	0.255	0.396	0.281

Table 6: Ablation on the impact and order of Integrated Attention and Cointegrated Attention. "Order" specifies the sequence, with lower numbers indicating earlier placement.  $\times$  indicates the component is removed.

450 Ablation on modeling approaches for Integrated Attention and Cointegrated Attention. We 451 conduct the following experiments: <sup>①</sup> both Integrated and Cointegrated Attention use channel-452 independent (CI) modeling, <sup>(2)</sup> Integrated Attention uses channel-dependent (CD) modeling while 453 Cointegrated Attention uses CI, 3 Integrated Attention uses CI while Cointegrated Attention uses 454 CD, and ④ both use CD modeling. The results in Tab. 7 show that modeling short-term inter-455 variates relationships can lead to severe spurious regression. CI modeling generally outperforms CD in scenarios with fewer channels (e.g., Weather), while CD excels when the number of channels 456 is large (e.g., Traffic). This aligns with recent findings that inter-channel dependencies become 457 increasingly important as the number of channels grows. We attribute this to the model's ability to 458 extract potential long-term stable relationships from non-stationary sequences when more channels 459 are present, thereby improving both forecasting accuracy and robustness. 460

Case	Integrated Attention	Cointegrated Attention	Wea	ther	So	lar	Elect	ricity	Tra	ffic
	CI or CD	CI or CD	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
1	CI	CI	0.218	0.259	<u>0.183</u>	<u>0.243</u>	0.157	<u>0.252</u>	0.387	0.276
2	CD	CI	<u>0.222</u>	<u>0.262</u>	0.184	0.247	0.160	0.255	0.387	0.280
3	CI	CD	0.218	0.259	0.181	0.239	0.149	0.246	0.360	0.255
4	CD	CD	<u>0.222</u>	0.263	<u>0.183</u>	0.247	<u>0.156</u>	0.254	<u>0.376</u>	0.269

Table 7: Ablation on modeling approaches for Integrated Attention and Cointegrated Attention. "CI" denotes channel independent and "CD" denotes channel-dependent modeling.

## 470 471 472

473

469

446

447

448

449

#### 6 NON-STATIONARITY AND DEPENDENCY MODELING ANALYSIS

474 Intra-variate Modeling. As shown in Fig. 4a, when non-stationarity is retained, the attention map 475 in the temporal dimension diverges, with the model focusing on multiple patches across a broader 476 time span. However, after removing non-stationarity, the attention map becomes more concentrated 477 on adjacent time steps, aligning with the causal nature of time series, where closer time steps are usually more correlated. Non-stationarity may cause the model to mistake distant similarities for 478 causality. By eliminating non-stationarity, the model better captures short-term variations and local 479 dependencies, enhancing its robustness and interpretability in handling complex time series data. 480

481 Inter-variate Modeling. Fig. 4b shows that removing non-stationarity narrows the model's atten-482 tion to a few inter-variate dependencies, while retaining non-stationarity enables the capture of more diverse and richer relationships. Non-stationary sequences help the model identify cointegration, 483 revealing hidden equilibrium mechanisms in multivariate time series. Preserving non-stationarity 484 enhances the model's ability to express complex inter-variate dependencies. Additionally, Fig. 4c 485 shows the impact of different patch downsampling rates on performance. For datasets with more



Figure 4: (a) Comparison of intra-variate attention maps under stationary and non-stationary conditions for different patches in the Electricity dataset. (b) Comparison of inter-variate attention maps between different variates under stationary and non-stationary conditions in the Solar dataset. (c) Impact of varying the number of downsampled patches M on forecasting performance across different datasets. See Tab. 16 for full results.



Figure 5: Visualization of the effect of retaining or removing non-stationarity in Integrated Attention and Cointegrated Attention on the Weather dataset for temperature (T) and dew point temperature  $(T_{dew})$ . (a) Both Integrated and Cointegrated Attention retain non-stationarity. (b) Both remove nonstationarity. (c) Only Integrated Attention retains non-stationarity. (d) Only Cointegrated Attention retains non-stationarity.

515

516

517

518

509

510

511

494

495

496

497

498

channels and stronger cointegration (e.g., Solar and Traffic), increasing downsampled patches initially improves predictions by preserving long-term features. However, too much downsampling adds computational cost and negatively affects smaller-channel datasets (e.g., Weather), so we carefully balanced downsampling rates based on dataset characteristics, as detailed in Tab. 9.

Real Case of Weather Forecast. Given the strong interrelationships between weather variables, 519 we analyze temperature T and dew point temperature  $T_{dew}$  from the Weather dataset. Dew point 520 measures atmospheric moisture and is typically closely linked to temperature. Without external 521 influences, such as water vapor or heat sources, the difference between temperature and dew point 522 is minimal, showing long-term cointegration. However, temperature tends to exhibit more short-523 term fluctuations due to external factors like sunlight and weather systems. As shown in Fig. 5, the 524 results demonstrate that spurious regressions can only be avoided by eliminating non-stationarity 525 during short-term modeling, while preserving it during long-term dependency modeling to capture 526 the underlying cointegration between variables.

### 527 528

529

## 7 CONCLUSION

530 In this paper, we address the dual challenges of non-stationarity in multivariate time series forecast-531 ing, specifically focusing on its distinct impacts on short-term and long-term modeling. To this end, 532 we propose TimeBridge, a novel framework that bridges the gap between non-stationarity and de-533 pendency modeling. By employing Integrated Attention to mitigate short-term non-stationarity and 534 Cointegrated Attention to preserve long-term dependencies, TimeBridge effectively captures both local dynamics and long-term cointegration. Comprehensive experiments across diverse datasets demonstrate that TimeBridge consistently achieves state-of-the-art performance in both short-term 537 and long-term forecasting tasks. Moreover, its exceptional performance on the CSI 500 and S&P 500 indices underscores its robustness and adaptability to complex real-world financial scenarios. 538 Our work paves the way for further exploration of models that balance the nuanced effects of nonstationarity, offering a promising direction for advancing multivariate time series forecasting.

#### 540 REFERENCES 541

547

551

553

559

565

566

569

581

582

583

- Faik Bilgili. Stationarity and cointegration tests: Comparison of engle-granger and johansen 542 methodologies. Erciyes Universitesi İktisadi ve İdari Bilimler Fakültesi Dergisi, pp. 131–141, 543 1998. 544
- Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of 546 gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555, 2014.
- Tao Dai, Beiliang Wu, Peiyuan Liu, Naiqi Li, Jigang Bao, Yong Jiang, and Shu-Tao Xia. Periodicity 548 decoupling framework for long-term series forecasting. International Conference on Learning 549 Representations, 2024. 550
- Abhimanyu Das, Weihao Kong, Andrew Leach, Rajat Sen, and Rose Yu. Long-term forecasting 552 with tiDE: Time-series dense encoder. arXiv preprint arXiv:2304.08424, 2023.
- Luo Donghao and Wang Xue. ModernTCN: A modern pure convolution structure for general time 554 series analysis. International Conference on Learning Representations, 2024. 555
- 556 Wei Fan, Pengyang Wang, Dongkun Wang, Dongjie Wang, Yuanchun Zhou, and Yanjie Fu. Dish-TS: a general paradigm for alleviating distribution shift in time series forecasting. In *Proceedings* 558 of the AAAI Conference on Artificial Intelligence, volume 37, pp. 7522–7529, 2023.
- Wei Fan, Kun Yi, Hangting Ye, Zhiyuan Ning, Qi Zhang, and Ning An. Deep frequency derivative 560 learning for non-stationary time series forecasting. International Joint Conference on Artificial 561 Intelligence, pp. 3944-3952, 2024. 562
- 563 Phillip Fanchon and Jeanne Wendel. Estimating var models under non-stationarity and cointegration: 564 alternative approaches for forecasting cattle prices. Applied Economics, 24(2):207–217, 1992.
  - S. Hochreiter and J. Schmidhuber. Long short-term memory. Neural Comput., 1997.
- 567 Narasimhan Jegadeesh and Sheridan Titman. Returns to buying winners and selling losers: Impli-568 cations for stock market efficiency. Journal of Finance, 48:65-91, 1993.
- Zahra Karevan and Johan AK Suykens. Transductive lstm for time-series prediction: An application 570 to weather forecasting. Neural Networks, 125:1-9, 2020. 571
- 572 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-573 versible instance normalization for accurate time-series forecasting against distribution shift. In 574 International Conference on Learning Representations, 2022.
- 575 Diederik P Kingma. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 576 2014. 577
- 578 Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term 579 temporal patterns with deep neural networks. In International ACM SIGIR conference on research 580 & development in information retrieval, pp. 95–104, 2018.
  - Naiqi Li, Zhikang Xia, Yiming Li, Ercan E. Kuruoğlu, Yong Jiang, and Shu-Tao Xia. Portfolio selection via graph-aware gaussian processes with generalized gaussian likelihood. IEEE Transactions on Artificial Intelligence, 5(2):505-515, 2024. doi: 10.1109/TAI.2023.3262456.
- 585 Hengxu Lin, Dong Zhou, Weiqing Liu, and Jiang Bian. Learning multiple stock trading patterns with temporal routing adaptor and optimal transport. In Proceedings of the 27th ACM SIGKDD 586 Conference on Knowledge Discovery & Data Mining, 2021.
- 588 Shengsheng Lin, Weiwei Lin, Wentai Wu, Haojun Chen, and Junjie Yang. Sparsetsf: Modeling long-589 term time series forecasting with 1k parameters. International Conference on Machine Learning, 590 2024. 591
- Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. SCInet: 592 Time series modeling and forecasting with sample convolution and interaction. Advances in Neural Information Processing Systems, 35:5816–5828, 2022a.

594 595 596	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</i> , 35:9881–9893, 2022b.
597 598 599 600	Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. iTransformer: Inverted transformers are effective for time series forecasting. <i>International Con-</i> <i>ference on Learning Representations</i> , 2024a.
601 602 603	Yong Liu, Chenyu Li, Jianmin Wang, and Mingsheng Long. Koopa: Learning non-stationary time series dynamics with koopman predictors. <i>Advances in Neural Information Processing Systems</i> , 36, 2024b.
605 606 607	Zhiding Liu, Mingyue Cheng, Zhi Li, Zhenya Huang, Qi Liu, Yanhu Xie, and Enhong Chen. Adap- tive normalization for non-stationary time series forecasting: A temporal slice perspective. <i>Ad-</i> <i>vances in Neural Information Processing Systems</i> , 36, 2023.
608 609 610 611	Xiang Ma, Xuemei Li, Lexin Fang, Tianlong Zhao, and Caiming Zhang. U-mixer: An unet-mixer architecture with stationarity correction for time series forecasting. <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pp. 14255–14262, 2024.
612	Rizwan Mushtaq. Augmented dickey fuller test. 2011.
613 614 615 616	Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. In <i>International Conference on Learning Representations</i> , 2023.
617 618	Antonio E Noriega and Daniel Ventosa-Santaulària. Spurious regression and trending variables. Oxford Bulletin of Economics and Statistics, 69(3):439–444, 2007.
619 620 621	Joon Y Park and Peter CB Phillips. Nonlinear regressions with integrated time series. <i>Econometrica</i> , 69(1):117–161, 2001.
622 623 624 625	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
626 627 628	James M. Poterba and Lawrence H. Summers. Mean reversion in stock prices: Evidence and impli- cations. <i>Journal of Financial Economics</i> , 1988.
629 630 631	Yao Qin, Dongjin Song, Haifeng Chen, Wei Cheng, Guofei Jiang, and Garrison W. Cottrell. A dual- stage attention-based recurrent neural network for time series prediction. In <i>International Joint</i> <i>Conference on Artificial Intelligence</i> , pp. 2627–2633, 2017. doi: 10.24963/ijcai.2017/366.
632 633 634	Rohnn Sanderson and Nancy L Lumpkin-Sowers. Buy and hold in the new age of stock market volatility: A story about etfs. <i>International Journal of Financial Studies</i> , 6(3):79, 2018.
635 636 637 638	Omer Berat Sezer, Mehmet Ugur Gudelek, and Ahmet Murat Ozbayoglu. Financial time series fore- casting with deep learning: A systematic literature review: 2005–2019. <i>Applied soft computing</i> , 90:106181, 2020.
639 640 641	Wanneng Shu, Ken Cai, and Neal Naixue Xiong. A short-term traffic flow prediction model based on an improved gate recurrent unit neural network. <i>IEEE Transactions on Intelligent Transportation Systems</i> , 23(9):16654–16665, 2021.
642 643 644 645	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in Neural Information Processing Systems</i> , 30, 2017a.
646	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,

Kaiser, and Illia Polosukhin. Attention is all you need. Advances in Neural Information Processing Systems, 30, 2017b.

648 649 650	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>International Confer-</i> <i>ence on Learning Representations</i> , 2022.
651 652 653 654	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. <i>International Conference on Learning Representations</i> , 2024.
655 656 657	Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans- formers with auto-correlation for long-term series forecasting. <i>Advances in Neural Information</i> <i>Processing Systems</i> , 34:22419–22430, 2021.
658 659 660 661	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>International Conference on Learning Representations</i> , 2023.
662 663 664	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.
665 666 667	Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In <i>International Conference on Learning Representations</i> , 2023.
669 670 671	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.
672 673 674 675 676	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</i> on Machine Learning, pp. 27268–27286. PMLR, 2022.
677 678 679	
680 681 682	
683 684 685	
686 687 688	
689 690 691	
692 693 694	
695 696	
698 699	

#### TIME SERIES INTEGRATION AND COINTEGRATION ANALYSIS А

## A.1 INTEGRATION AND ADF TEST

A time series is said to be integrated of order k, denoted as I(k), if it becomes stationary after differencing k times. For instance, a series  $X_t$  is I(1) if its first difference  $\Delta X_t = X_t - X_{t-1}$  is stationary. To test for non-stationarity, the Augmented Dickey-Fuller (ADF) test (Mushtaq, 2011) is commonly used. It examines the null hypothesis that a unit root is present, indicating non-stationarity: 

$$\Delta X_t = \alpha + \beta t + \gamma X_{t-1} + \sum_{i=1}^p \delta_i \Delta X_{t-i} + \epsilon_t$$

Here,  $\Delta X_t$  is the differenced series,  $\gamma$  is the coefficient on the lagged series, and  $\epsilon_t$  is the error term. Rejecting the null hypothesis ( $\gamma = 0$ ) indicates stationarity, while failing to reject it implies non-stationarity. Non-stationary data can lead to spurious regressions, where unrelated temporal intervals appear to be correlated due to common trends. We report the ADF test results in Tab. 8.

#### A.2 COINTEGRATION AND EG TEST

Cointegration occurs when two or more non-stationary series move together over time, maintaining a stable, long-term relationship. For example, if  $X_t$  and  $Y_t$  are both I(1), they are cointegrated if there exists a stationary linear combination,  $Z_t = X_t - \beta Y_t$ . This indicates a shared stochastic trend. The Engle-Granger (EG) test (Bilgili, 1998) for cointegration involves two steps:

1. Estimate Long-term Relationship. Regress  $X_t$  on  $Y_t$  using Ordinary Least Squares (OLS):

$$X_t = \alpha + \beta Y_t + \epsilon_t,$$

where  $\epsilon_t$  are the residuals.

2. ADF Test on Residuals. Apply the ADF test to  $\epsilon_t$ :

$$\Delta \epsilon_t = \gamma \epsilon_{t-1} + \sum_{i=1}^p \delta_i \Delta \epsilon_{t-i} + \nu_t$$

If the residuals are stationary,  $X_t$  and  $Y_t$  are cointegrated.

Cointegration is vital for capturing long-term relationships between variables, providing a robust foundation for multivariate time series modeling. Ignoring cointegration can result in models that miss significant underlying connections, reducing forecasting accuracy and reliability. We report the EG test results in Tab. 8. 

#### В METRICS

## **B.1** LONG-TERM FORECASTING

We use Mean Squared Error (MSE) and Mean Absolute Error (MAE) as evaluation metrics. Given the ground truth values  $X_i$  and the predicted values  $X_i$ , these metrics are defined as follows:

$$\mathsf{MSE} = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{X}_i - \hat{\mathbf{X}}_i)^2, \quad \mathsf{MAE} = \frac{1}{N} \sum_{i=1}^{N} |\mathbf{X}_i - \hat{\mathbf{X}}_i|,$$

where N is the total number of predictions.

## B.2 SHORT-TERM FORECASTING

We use MAE (the same as defined above), Mean Absolute Percentage Error (MAPE), and Root
Mean Squared Error (RMSE) to evaluate the performance. These metrics are defined as follows:

$$MAPE = \frac{1}{N} \sum_{i=1}^{N} \left| \frac{\mathbf{X}_i - \hat{\mathbf{X}}_i}{\mathbf{X}_i} \right| \times 100, \quad RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (\mathbf{X}_i - \hat{\mathbf{X}}_i)^2}.$$

## 765 B.3 FINANCIAL FORECASTING

We use six widely recognized metrics to assess the overall performance of each method: Annual
Return Ratio (ARR), Annual Volatility (AVol), Maximum Drawdown (MDD), Annual Sharpe Ratio
(ASR), Calmar Ratio (CR), and Information Ratio (IR). Lower absolute values of AVol and MDD,
coupled with higher values of ARR, ASR, CR, and IR, indicate better performance.

• **ARR** quantifies the percentage increase or decrease in the value of an investment over a year.

$$ARR = (1 + \text{Total Return})^{\frac{1}{n}} - 1$$

• AVol measures the volatility of an investment's returns over the course of a year.  $R_p$  denotes the daily return of the portfolio.

$$\text{AVol} = \sqrt{\text{Var}(R_p)}.$$

• MDD indicates the maximum decline from a peak to a trough in the value of an investment.

$$ext{MDD} = - ext{max}igg(rac{p_{peak} - p_{trough}}{p_{peak}}igg).$$

• ASR reflects the risk-adjusted return of an investment over a year.

$$ASR = \frac{ARR}{AVol}$$

• **CR** compares the average annual return of an investment to its maximum drawdown.

$$CR = \frac{ARR}{|MDD|}$$

• IR evaluates the excess return of an investment relative to a benchmark, adjusted for its volatility.  $R_b$  is the daily return of the market index.

$$IR = \frac{\text{mean}(R_p - R_b)}{\text{std}(R_p - R_b)}$$

## C DATASETS

We conduct extensive experiments on eight widely-used time series datasets for long-term forecasting. Additionally, we use the PeMS datasets for short-term forecasting and the CSI 500 and S&P 500 indices for financial forecasting. We report the statistics in Tab. 8. Detailed descriptions of these datasets are as follows:

- ETT (Electricity Transformer Temperature) dataset (Zhou et al., 2021) encompasses temperature and power load data from electricity transformers in two regions of China, spanning from 2016 to 2018. This dataset has two granularity levels: ETTh (hourly) and ETTm (15 minutes).
- 807 (2) Weather dataset (Wu et al., 2023) captures 21 distinct meteorological indicators in Germany, meticulously recorded at 10-minute intervals throughout 2020. Key indicators in this dataset include air temperature, visibility, among others, offering a comprehensive view of the weather dynamics.

(3) Electricity dataset (Wu et al., 2023) features hourly electricity consumption records in kilowatt-hours (kWh) for 321 clients. Sourced from the UCL Machine Learning Repository, this dataset covers the period from 2012 to 2014, providing valuable insights into consumer electricity usage patterns.

- (4) **Traffic** dataset (Wu et al., 2023) includes data on hourly road occupancy rates, gathered by 862 detectors across the freeways of the San Francisco Bay area. This dataset, covering the years 2015 to 2016, offers a detailed snapshot of traffic flow and congestion.
- (5) **Solar-Energy** dataset (Lai et al., 2018) contains solar power production data recorded every 10 minutes throughout 2006 from 137 photovoltaic (PV) plants in Alabama.
- (6) PeMS dataset (Liu et al., 2022a) comprises four public traffic network datasets (PeMS03, PeMS04, PeMS07, and PeMS08), constructed from the Caltrans Performance Measurement System (PeMS) across four districts in California. The data is aggregated into 5-minute intervals, resulting in 12 data points per hour and 288 data points per day.
- (7) **CSI 500**<sup>1</sup> contains 502 stocks listed on the Shanghai and Shenzhen stock exchanges in China from 2018 to 2023, including close, open, high, low, volume and turnover data.
- (8) **S&P 500**<sup>2</sup> contains 487 stocks representing diverse sectors within the U.S. economy from 2018 to 2023, including close, open, high, low and volume data.

Tasks	Dataset	Dim	Prediction Length	Dataset Size	Frequency	ADF <sup>†</sup>	EG <sup>‡</sup>
	ETTm1	7	$\{96, 192, 336, 720\}$	(34465, 11521, 11521)	15 min	-14.98	20
	ETTm2	7	$\{96, 192, 336, 720\}$	(34465, 11521, 11521)	$15 \min$	-5.66	17
	ETTh1	7	$\{96, 192, 336, 720\}$	(8545, 2881, 2881)	$15 \min$	-5.91	11
Long-term	ETTh2	7	$\{96, 192, 336, 720\}$	(8545, 2881, 2881)	$15 \min$	-4.13	10
Forecasting	Electricity	321	$\{96, 192, 336, 720\}$	(18317, 2633, 5261)	1 hour	-8.44	39567
	Traffic	862	$\{96, 192, 336, 720\}$	(12185, 1757, 3509)	1 hour	-15.02	354627
	Weather	21	$\{96, 192, 336, 720\}$	(36792, 5271, 10540)	$10 \min$	-26.68	77
	Solar-Energy	137	$\{96, 192, 336, 720\}$	(36601, 5161, 10417)	$10 \min$	-37.23	8373
	PeMS03	358	12	(15617, 5135, 5135)	$5 \min$	-19.05	-
Short-term	PeMS04	307	12	(10172, 3375, 3375)	$5 \min$	-15.66	-
Forecasting	PeMS07	883	12	(16911, 5622, 5622)	$5 \min$	-20.60	-
	PeMS08	170	12	(10690, 3548, 265)	$5 \min$	-16.04	-
Financial	CSI 500	502	1	(943, 242, 242)	1 day	-3.06	-
Forecasting	S&P 500	487	1	(1008, 251, 249)	1 day	-2.80	-

† Augmented Dickey-Fuller (ADF) Test: A smaller ADF test result indicates a more stationary time series data.

‡ Engle-Granger (EG) Test: A bigger EG test result indicates the data contains more cointegration relationships.

Table 8: Dataset detailed descriptions. "Dataset Size" denotes the total number of time points in (Train, Validation, Test) split respectively. "Prediction Length" denotes the future time points to be predicted. "Frequency" denotes the sampling interval of time points.

To further illustrate the degree of non-stationarity in the datasets, we conduct additional experiments using a Random Walk series (representing maximum non-stationarity) and Gaussian white noise (representing near-stationarity). The Random Walk series is generated using the formula  $X_t$  =  $X_{t-1} + \epsilon_t$  with  $\epsilon_t \sim \mathcal{N}(0,1)$ , where we set t = 10,000 and simulate 100 iterations. The average ADF value for the Random Walk series is -1.53, indicating a high degree of non-stationarity. In contrast, for the Gaussian white noise series, generated as  $X_t \sim \mathcal{N}(0,1)$  with the same settings, the average ADF value is -97.54, indicating strong stationarity. Comparing these results with those in Tab. 8, we can see that most datasets exhibit significant non-stationarity, especially the ETT, CSI 500, and S&P 500 datasets. 

<sup>&</sup>lt;sup>1</sup>https://cn.investing.com/indices/china-securities-500

<sup>&</sup>lt;sup>2</sup>https://hk.finance.yahoo.com/quote/%5EGSPC/history/

To analyze cointegration, we conducted the Engle-Granger (EG) test on all eight datasets of longterm forecasting. The results indicate that datasets with more channels tend to exhibit more extensive cointegration relationships. This is particularly evident in datasets like Electricity and Traffic, which show significantly higher EG test values, reflecting a greater abundance of long-term equilibrium relationships among variates. For these high-dimensional datasets, effectively modeling the intricate cointegration structures is crucial, as neglecting these long-term dependencies can result in suboptimal predictions.

## D IMPLEMENTATION DETAILS

All experiments are implemented in PyTorch (Paszke et al., 2019) and conducted on two NVIDIA RTX 3090 24GB GPUs. We use the Adam optimizer (Kingma, 2014). The batch size is set to 16 for the Electricity and Traffic datasets, and 32 for all other datasets. All models are trained for 10 epochs. Tab. 9 provides detailed hyperparameter settings for each dataset. For the four ETT datasets, the relatively small number of channels results in less pronounced long-term cointegration relationships, as evidenced by the low EG test results in Tab. 8. Therefore, we focus on modeling short-term intra-variate variations only.

	Num. of Integrated	Num. of Cointegrated	$\mid N$	$\mid M$	lr	d_model	d_ff
ETTh1	2	0	30	_	1 <i>e</i> -4	128	256
ETTh2	2	0	30	-	1e-4	64	128
ETTm1	2	0	30	-	1e-4	32	128
ETTm2	2	0	30	-	1e-4	32	64
Weather	1	1	30	12	1e-4	128	128
Solar	1	1	30	12	5e-4	128	128
Electricity	2	1	30	4	5e-4	512	512
Traffic	2	2	30	4	1e-3	512	512

Table 9: Hyperparameter settings for different datasets. "N" denotes the number of patches. "M" denotes the number of patches after the patch downsampling block. "Ir" denotes the learning rate. "d\_model" and "d\_ff" denote the model dimension of attention layers and feed-forward layers, respectively.

## E FULL RESULTS

## E.1 MAIN EXPERIMENTS

Tab. 10 and Tab. 11 present the full results for long-term forecasting, including both the original results from their respective papers and those obtained through hyperparameter search. The hyper-parameter search process involved exploring input lengths  $I \in \{96, 192, 336, 512, 720\}$ , learning rates from  $10^{-5}$  to 0.05, encoder layers from 1 to 5,  $d_{model}$  values from 16 to 512, and training epochs from 10 to 100. In both settings, TimeBridge consistently achieved the best performance, demonstrating its effectiveness and robustness. Additionally, for financial forecasting, we included three additional strong baselines: ALSTM (Qin et al., 2017), GRU (Chung et al., 2014), and TRA (Lin et al., 2021). The results in Tab. 12 show that TimeBridge continues to outperform these meth-ods, further validating its superiority.

## 913 E.2 ABLATION STUDIES

We present the full results of the ablation studies discussed in the main text. Tab. 13 provides the complete results of the ablation on removing non-stationarity in both Integrated and Cointegrated Attention. Tab. 14 reports the full results on the impact and order of Integrated and Cointegrated Attention, with an illustrative visualization in Fig. 6. Additionally, Tab. 15 shows the results of abla-

918 010	M	odels	Timel	Bridge	iTrans	former 24a)	PE (20)	0F 24)	Timel	Mixer	Patcl	nTST 23)	Crossf	former 23)	FEDf	ormer	Moder (20	nTCN 24)	MI (20	CN 22)	Time (20)	sNet	DLi (20	near (23)
919	М	letric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
921	_	96	0.297	0.353	0.334	0.368	0.280	0.335	0.320	0.357	0.293	0.346	0.316	0.373	0.379	0.419	$\frac{0.292}{0.332}$	0.346	0.316	0.362	0.338	0.375	0.299	$\frac{0.343}{0.365}$
922	TTm	336	0.366	0.375	0.426	0.391	0.354	0.339	0.390	0.381	0.355	0.392	0.431	0.411	0.420	0.459	0.352	0.308	0.303	0.390	0.374	0.387	0.369	0.386
923	ш	/20	0.414	0.423	0.491	0.459	0.405	0.413	0.454	0.441	0.416	0.420	0.600	0.547	0.543	0.490	0.416	0.381	0.481	0.476	0.478	0.450	0.425	0.421
924		96	0.355	0.388	0.180	0.264	0.162	0.253	0.175	0.393	0.355	0.382	0.431	0.358	0.203	0.287	0.166	0.256	0.392	0.287	0.400	0.400	0.167	0.260
925	m2	192	0.215	0.291	0.250	0.309	$\frac{0.219}{0.270}$	0.291	0.237	0.299	0.223	0.296	0.345	0.400	0.269	0.328	0.222	$\frac{0.293}{0.324}$	0.262	0.326	0.249	0.309	0.224	0.303
926	ΕIJ	720	0.348	0.325	0.412	0.407	0.358	0.320 0.380	0.391	0.396	0.362	0.325	1.208	0.753	0.421	0.415	<u>0.351</u>	0.381	0.389	0.333	0.408	0.403	0.397	0.421
927		Avg.	0.246	0.310	0.288	0.332	<u>0.252</u>	<u>0.313</u>	0.275	0.323	0.256	0.317	0.621	0.510	0.305	0.349	0.253	0.314	0.290	0.343	0.291	0.333	0.267	0.332
928	_	96 192	0.358 0.388	0.392 0.411	0.386	0.405	0.357 0.397	<b>0.388</b> 0.412	0.375 0.429	0.400	0.370	0.400	0.405	0.426 0.444	0.376	0.419 0.448	0.368	0.394	0.398	0.427	0.384	0.402	0.375	0.399 0.416
929	ETTh	336 720	$\frac{0.401}{0.447}$	$\frac{0.419}{0.458}$	0.487	0.458	0.409	0.422	0.484	0.458	0.422	0.440	0.440	0.461	0.459	0.465	0.391 0.450	<b>0.412</b> 0.461	0.440	0.460	0.491	0.469	0.439	0.443
930	Π	Avg.	0.399	0.420	0.454	0.447	0.399	0.455	0.447	0.440	0.413	0.434	0.446	0.464	0.440	0.460	0.404	0.420	0.440	0.462	0.458	0.450	0.423	0.437
931		96	0.295	0.354	0.297	0.349	0.272	0.333	0.289	0.341	0.274	0.337	0.628	0.563	0.346	0.388	0.263	0.332	0.332	0.377	0.340	0.374	0.289	0.353
932	Th2	192 336	0.351	0.389 0.397	0.380	0.400 0.432	0.335	$\frac{0.375}{0.377}$	0.372 0.386	0.392 0.414	0.314 0.329	0.382	0.703 0.827	0.624 0.675	0.429	0.439 0.487	0.320 0.313	0.374 0.376	0.422 0.447	0.441 0.474	0.402 0.452	0.414 0.452	0.383	0.418 0.465
933	ΕT	720	0.388	0.436	0.427	0.445	0.375	0.417	0.412	0.434	0.379	<u>0.422</u>	1.181	0.840	0.463	0.474	0.392	0.433	0.442	0.467	0.462	0.468	0.605	0.551
934		Avg.	0.346	0.394	0.383	0.407	0.327	0.376	0.364	0.395	<u>0.324</u>	0.381	0.835	0.675	0.434	0.447	0.322	<u>0.379</u>	0.411	0.440	0.414	0.427	0.431	0.447
935	er	96 192	0.143 0.185	0.192 0.235	0.174 0.221	0.214 0.254	$\frac{0.147}{0.192}$	<b>0.194</b> 0.239	0.163 0.208	0.209 0.250	0.149 0.194	$\frac{0.198}{0.241}$	0.153 0.197	0.217 0.269	0.217 0.276	0.296 0.336	0.149 0.196	0.200 0.245	0.161 0.220	0.229 0.281	0.172 0.219	0.220 0.261	0.176 0.220	0.237 0.282
936	Veath	336 720	0.237 0.307	0.277 0.330	0.278	0.296	0.244	0.279	0.251	0.287	0.245	0.282	0.495	0.515 0.542	0.339	0.380 0.428	<b>0.238</b> 0.314	<b>0.277</b> 0.334	0.278	0.331	0.280	0.306	0.265	0.319
937	-	Avg.	0.218	0.259	0.258	0.279	0.225	0.261	0.240	0.271	0.226	0.264	0.343	0.386	0.309	0.360	0.224	0.264	0.243	0.299	0.259	0.287	0.246	0.300
938	~	96	0.118	0.218	0.148	0.240	<u>0.127</u>	0.219	0.153	0.247	0.129	0.222	0.187	0.283	0.183	0.297	0.129	0.226	0.164	0.269	0.168	0.272	0.140	0.237
939	tricit	192 336	0.142 0.156	0.237 0.252	0.162 0.178	0.253 0.269	0.145 0.162	$\frac{0.237}{0.255}$	0.166 0.185	0.256 0.277	0.147 0.163	0.240 0.259	0.258 0.323	0.330 0.369	0.195	0.308 0.313	$\frac{0.143}{0.161}$	0.239 0.259	0.177 0.193	0.285 0.304	0.184 0.198	0.289 0.300	0.153 0.169	0.249 0.267
940	Elec	720	0.179	0.278	0.225	0.317	0.200	0.290	0.225	0.310	0.197	0.290	0.404	0.423	0.231	0.343	0.191	<u>0.286</u>	0.212	0.321	0.220	0.320	0.203	0.301
941		Avg.	0.149	0.246	0.178	0.270	0.159	0.250	0.182	0.272	0.159	0.253	0.293	0.351	0.205	0.315	<u>0.156</u>	0.253	0.187	0.295	0.192	0.295	0.166	0.264
942	<u>.</u> 2	96 192	0.340 0.343	0.240 0.250	0.395	0.268 0.276	$\frac{0.351}{0.374}$	$\frac{0.238}{0.248}$	0.462 0.473	0.285 0.296	0.360 0.379	0.249 0.256	0.512 0.523	0.290 0.297	0.562	0.349 0.346	0.368 0.379	0.253 0.261	0.519 0.537	0.309 0.315	0.593 0.617	0.321 0.336	0.410 0.423	0.282 0.287
943	Traff	336 720	0.363 0.393	0.257 0.271	0.433	0.283	0.386	0.253	0.498	0.296	0.392	0.264	0.530	0.300	0.570	0.323	0.397	0.270	0.534	0.313	0.629	0.336	0.436	0.296
944	-	Avg.	0.360	0.255	0.428	0.282	0.383	0.254	0.484	0.297	0.391	0.264	0.535	0.300	0.573	0.347	0.396	0.270	0.542	0.316	0.620	0.336	0.434	0.295
945		96	0.161	0.224	0.203	0.237	0.179	0.246	0.189	0.259	0.190	0.273	0.181	0.240	0.209	0.330	0.202	0.263	0.190	0.250	0.285	0.330	0.289	0.377
946	olar	192 336	0.177 0.188	0.237 0.244	0.233	0.261	$\frac{0.205}{0.210}$	0.265	0.222	0.283	0.204	0.302	0.196	0.252	0.274	0.400	0.223	0.279	0.225	0.270	0.309	0.342	0.319	0.397 0.415
947	Sc	720	0.197	0.252	0.249	0.275	0.225	0.281	0.223	0.285	0.221	0.310	0.220	0.256	0.365	0.459	0.247	0.292	0.323	0.362	0.346	0.355	0.356	0.412
948		Avg.	0.181	0.239	0.233	0.262	<u>0.205</u>	0.265	0.216	0.280	0.207	0.294	0.204	0.248	0.296	0.407	0.228	0.282	0.247	0.296	0.319	0.348	0.329	0.400

Table 10: Full results of long-term forecasting from the original papers. All results are averaged across four different prediction lengths:  $O \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and underlined, respectively.

tion on different modeling approaches for these attention mechanisms. Finally, Tab. 16 presents the results of varying the number of downsampled patches M and its effect on forecasting performance.

#### F STATISTICAL ANALYSIS

We repeat all experiments three times and report the standard deviations for both our model and the second-best baseline, along with the results of statistical significance tests. Tab. 17, Tab. 18, and Tab. 19 present the results for long-term forecasting, short-term forecasting, and financial forecasting, respectively.

963 964 965

966

949

950

951

952 953 954

955

956 957

958 959

960

961

962

#### VISUALIZATION G

967 Fig. 7 visualizes short-term fluctuations and long-term cointegration across stock sectors. Fig. 8 pro-968 vides additional examples of intra-variate attention maps comparing stationary and non-stationary 969 conditions for different patches in the Electricity dataset. Fig. 9 shows further examples of intervariate attention maps in the Solar dataset under both stationary and non-stationary conditions. 970 Fig. 10, Fig. 11, Fig. 12, and Fig. 13 present long-term forecasting visualizations for Weather, So-971 lar, Electricity, and Traffic datasets, respectively. We display the last 96 input steps based on each model's optimal input length, along with the corresponding 96 predicted steps. Finally, Fig. 14 illustrates short-term forecasting for the PeMS03 dataset, where each model predicts 12 steps from a
96-step input.

Mo	odels	TimeBridge (Ours)	iTransformer (2024a)	PDF (2024)	TimeMixer (2024)	PatchTST (2023)	Crossformer (2023)	FEDformer (2022)	ModernTCN (2024)	MICN (2022)	TimesNet (2023)	DLinear (2023)
М	etric	MSE MAE	MSE MAE	MSE MAE								
ETTm1	96 192 336 720	0.297 0.353 0.333 0.375 0.366 0.399 0.414 0.423	0.300 0.353 0.345 0.382 0.374 0.398 0.429 0.430	0.277 0.337 0.316 0.364 0.346 0.381 0.402 0.409	$\begin{array}{c} 0.291 \\ 0.327 \\ 0.360 \\ 0.415 \\ 0.415 \end{array} \begin{array}{c} 0.340 \\ 0.365 \\ 0.381 \\ 0.417 \end{array}$	0.293 0.346 0.333 0.370 0.369 0.392 0.416 0.420		$ \begin{vmatrix} 0.326 & 0.390 \\ 0.365 & 0.415 \\ 0.392 & 0.425 \\ 0.446 & 0.458 \end{vmatrix} $		0.314 0.360 0.359 0.387 0.398 0.413 0.459 0.464	0.338 0.375 0.371 0.387 0.410 0.411 0.478 0.450	$\begin{array}{c} 0.299 & 0.343 \\ 0.335 & \underline{0.365} \\ 0.369 & \underline{0.386} \\ 0.425 & 0.421 \end{array}$
1	Avg.	0.353 0.388	0.362 0.391	0.335 0.373	0.348 0.375	0.353 0.382	0.420 0.435	0.382 0.422	0.351 0.381	0.383 0.406	0.400 0.406	0.357 0.379
ETTm2	96 192 336 720 <i>Avg.</i>	0.158         0.249           0.215         0.291           0.263         0.323           0.348         0.376           0.246         0.310	0.175 0.266 0.242 0.312 0.282 0.340 0.378 0.398	0.159         0.251           0.217         0.292           0.266         0.325           0.345         0.375	0.164 0.254 0.223 0.295 0.279 0.330 0.359 0.383 0.256 0.315	0.166 0.256 0.223 0.296 0.274 0.329 0.362 0.385 0.256 0.317	0.263 0.359 0.345 0.400 0.469 0.496 0.996 0.750 0.518 0.501	0.180 0.271 0.252 0.318 0.324 0.364 0.410 0.420 0.292 0.343	0.166 0.256 0.222 0.293 0.272 0.324 0.351 0.381 0.253 0.314	0.178 0.273 0.245 0.316 0.295 0.350 0.389 0.406	0.187 0.267 0.249 0.309 0.321 0.351 0.497 0.403 0.291 0.333	0.167 0.260 0.224 0.303 0.281 0.342 0.397 0.421
ETTh1	96 192 336 720 <i>Avg.</i>	0.358         0.392           0.388         0.411           0.401         0.419           0.447         0.458           0.399         0.420	0.386 0.405 0.424 0.440 0.449 0.460 0.495 0.487	0.356         0.391           0.390         0.413           0.402         0.421           0.432         0.455           0.395         0.420	0.361 0.390 0.409 0.414 0.430 0.429 0.445 0.460	0.370 0.400 0.413 0.429 0.422 0.440 0.447 0.468	0.386 0.426 0.413 0.442 0.440 0.461 0.519 0.524	0.376 0.415 0.423 0.446 0.444 0.462 0.469 0.492	0.368 0.394 0.405 0.413 0.391 0.412 0.450 0.461	0.396 0.427 0.430 0.453 0.433 0.458 0.474 0.508	0.384 0.402 0.557 0.436 0.491 0.469 0.521 0.500 0.458 0.450	0.375 0.399 0.405 0.416 0.439 0.443 0.472 0.490 0.423 0.437
ETTh2	96 192 336 720	0.295 0.354 0.351 0.389 0.351 0.397 0.388 0.436	0.297 0.348 0.371 0.403 0.404 0.428 0.424 0.444	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	0.271 0.330 0.317 0.402 0.332 0.396 0.342 0.408	0.274 0.337 0.314 0.382 0.329 0.384 0.379 0.422	0.611 0.557 0.703 0.624 0.827 0.675 1.094 0.775	0.332 0.374 0.407 0.446 0.400 0.447 0.412 0.469	<b>0.263</b> 0.332 0.320 0.374 0.313 0.376 0.392 0.433	0.289 0.357 0.409 0.438 0.417 0.452 0.426 0.473	0.340 0.374 0.402 0.414 0.452 0.452 0.462 0.468	0.289 0.353 0.383 0.418 0.448 0.465 0.605 0.551
Weather	96 192 336 720	0.143 0.192 0.185 0.235 0.237 0.277 0.307 0.330	0.159 0.208 0.200 0.248 0.253 0.289 0.321 0.338	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	0.147 0.197 0.189 0.239 0.241 0.280 0.310 <u>0.330</u>	0.149 0.198 0.194 0.241 0.245 0.282 0.314 0.334	$\begin{vmatrix} 0.809 & 0.038 \\ 0.146 & 0.212 \\ 0.195 & 0.261 \\ 0.252 & 0.311 \\ 0.318 & 0.363 \end{vmatrix}$	0.217 0.296 0.275 0.329 0.339 0.377 0.389 0.409	0.149 0.200 0.196 0.245 0.238 0.277 0.314 0.334	0.161 0.226 0.220 0.283 0.275 0.328 0.311 0.356	0.172 0.220 0.219 0.261 0.280 0.306 0.365 0.359	0.431 0.447 0.152 0.237 0.220 0.282 0.265 0.319 0.323 0.362
Electricity	Avg. 96 192 336 720 Avg.	0.218         0.259           0.118         0.218           0.142         0.237           0.156         0.252           0.179         0.278           0.149         0.246	0.233         0.271           0.138         0.237           0.157         0.256           0.167         0.264           0.194         0.286           0.164         0.261	0.220         0.259           0.126         0.220           0.145         0.237           0.159         0.255           0.194         0.287           0.156         0.250	0.222         0.262           0.129         0.224           0.140         0.220           0.161         0.255           0.194         0.287           0.156         0.246	0.226         0.264           0.129         0.222           0.147         0.240           0.163         0.259           0.197         0.290           0.159         0.253	0.228         0.287           0.135         0.237           0.160         0.262           0.182         0.282           0.246         0.337           0.181         0.279	0.305 0.287 0.183 0.297 0.195 0.308 0.212 0.313 0.231 0.343 0.205 0.315	0.224         0.264           0.129         0.226           0.143         0.239           0.161         0.259           0.191         0.286           0.156         0.253	0.242         0.298           0.159         0.267           0.168         0.279           0.196         0.308           0.203         0.312           0.182         0.292	0.259 0.287 0.168 0.272 0.184 0.289 0.198 0.300 0.220 0.320 0.192 0.295	0.240 0.300 0.140 0.237 0.152 0.249 0.169 0.267 0.203 0.301 0.166 0.264
Traffic	96 192 336 720	0.340 0.240 0.343 0.250 0.363 0.257 0.393 0.271	0.363 0.265 0.385 0.273 0.396 0.277 0.445 0.312	0.350         0.239           0.363         0.247           0.376         0.258           0.419         0.279	0.360 0.249 0.375 0.250 0.385 0.270 0.430 0.281	0.360 0.249 0.379 0.256 0.392 0.264 0.432 0.286	0.512 0.282 0.501 0.273 0.507 0.279 0.571 0.301	0.562 0.349 0.562 0.346 0.570 0.323 0.596 0.368	0.368 0.253 0.379 0.261 0.397 0.270 0.440 0.296	0.508 0.301 0.536 0.315 0.525 0.310 0.571 0.323	0.593 0.321 0.617 0.336 0.629 0.336 0.640 0.350	0.410 0.282 0.423 0.287 0.436 0.296 0.466 0.315
Solar	Avg. 96 192 336 720 Avg.	0.360 0.255 0.161 0.224 0.177 0.237 0.188 0.244 0.197 0.252 0.181 0.239	0.397         0.282           0.188         0.242           0.193         0.258           0.195         0.259           0.223         0.281           0.200         0.260	0.377         0.256           0.179         0.246           0.205         0.265           0.210         0.270           0.225         0.281           0.205         0.265	0.387         0.262           0.167         0.220           0.187         0.249           0.200         0.258           0.215         0.250           0.192         0.244	0.391         0.264           0.178         0.229           0.189         0.246           0.198         0.249           0.209         0.256           0.194         0.245	0.523         0.284           0.166         0.230           0.186         0.237           0.203         0.243           0.210         0.256           0.191         0.242	0.573         0.347           0.201         0.304           0.237         0.337           0.254         0.362           0.280         0.397           0.243         0.350	0.396         0.270           0.202         0.263           0.223         0.279           0.241         0.292           0.247         0.292           0.248         0.282	0.535         0.312           0.188         0.252           0.215         0.280           0.222         0.267           0.226         0.264           0.213         0.266	0.620 0.336 0.219 0.314 0.231 0.322 0.246 0.337 0.280 0.363 0.244 0.334	0.434 0.295 0.216 0.287 0.244 0.305 0.263 0.319 0.264 0.324 0.247 0.309

Table 11: Full results of long-term forecasting of hyperparameter searching. All results are averaged across four different prediction lengths:  $O \in \{96, 192, 336, 720\}$ . The best and second-best results are highlighted in **bold** and <u>underlined</u>, respectively.

Models			CSI	500					S&P	500	
Models BLSW (1988) CSM (1993) LSTM (1997) ALSTM (2017) GRU (2014) Transformer (2017b) TRA (2021) PatchTST (2023) iTransformer (2024a) TimeMixer (2024) Crossformer (2023)	ARR↑	AVol↓	$\text{MDD}{\downarrow}$	ASR↑	CR↑	IR↑	ARR↑	AVol↓	$\text{MDD}{\downarrow}$	ASR↑	CR↑
BLSW (1988)	0.110	0.227	-0.155	0.485	0.710	0.446	0.199	0.318	-0.223	0.626	0.892
CSM (1993)	0.015	0.229	-0.179	0.066	0.084	0.001	0.099	0.250	-0.139	0.396	0.712
LSTM (1997)	-0.008	0.159	-0.172	-0.047	-0.044	-0.128	0.142	0.162	-0.178	0.877	0.798
ALSTM (2017)	0.016	0.162	-0.192	0.101	0.086	0.014	0.191	0.161	-0.150	1.186	1.273
GRU (2014)	-0.004	0.159	-0.193	-0.028	-0.023	-0.118	0.124	0.169	-0.139	0.734	0.829
Transformer (2017b)	0.154	0.156	-0.135	0.986	1.143	0.867	0.135	0.159	-0.140	0.852	0.968
TRA (2021)	0.125	0.162	-0.145	0.776	0.866	0.657	0.184	0.166	-0.158	1.114	1.172
PatchTST (2023)	0.118	0.152	-0.127	0.776	0.923	0.735	0.146	0.167	-0.140	0.877	1.042
iTransformer (2024a)	0.214	0.168	-0.164	1.276	1.309	1.173	0.159	0.170	-0.139	0.941	1.150
TimeMixer (2024)	0.078	0.153	-0.114	0.511	0.685	0.385	0.254	0.162	-0.131	1.568	1.938
Crossformer (2023)	-0.039	0.163	-0.217	-0.238	-0.179	-0.350	<u>0.284</u>	0.159	-0.114	<u>1.786</u>	2.491
TimeBridge	0.285	0.203	-0.196	1.405	1.453	1.317	0.326	0.169	-0.142	1.927	2.298

Table 12: Full results for financial time series forecasting in CSI 500 and S&P 500 datasets.

1	0	26
1	0	27

Integrated Attention	Cointegrated Attention		Wea	ather	So	lar	Elect	ricity	Tra	aff
+ Norm?	+ Norm?	Length	MSE	MAE	MSE	MAE	MSE	MAE	MSE	
		96	0.144	0.193	0.163	0.227	0.124	0.221	0.342	
		192	0.186	0.235	0.180	0.240	0.144	0.240	0.351	
×	×	336	0.239	0.279	0.191	0.248	0.158	0.254	0.374	
		720	0.311	0.333	0.197	0.253	0.184	0.282	0.418	
		Avg.	0.220	<u>0.260</u>	<u>0.183</u>	<u>0.242</u>	<u>0.153</u>	<u>0.249</u>	<u>0.371</u>	
		96	0.144	0.193	0.164	0.227	0.124	0.220	0.348	
×		192	0.188	0.237	0.180	0.240	0.146	0.240	0.370	
	$\checkmark$	336	0.239	0.279	0.191	0.248	0.161	0.258	0.382	
		720	0.308	0.331	0.197	0.293	0.189	0.285	0.422	
		Avg.	0.220	0.260	<u>0.183</u>	0.252	0.155	0.251	0.381	
		96	0.143	0.192	0.161	0.224	0.118	0.218	0.340	
		192	0.185	0.235	0.177	0.237	0.142	0.237	0.343	
$\checkmark$	×	336	0.237	0.277	0.188	0.244	0.156	0.252	0.363	
		720	0.307	0.330	0.197	0.252	0.179	0.278	0.393	
		Avg.	0.218	0.259	0.181	0.239	0.149	0.246	0.360	
		96	0.143	0.193	0.163	0.227	0.123	0.219	0.343	
		192	0.186	0.236	0.180	0.239	0.144	0.239	0.367	
$\checkmark$	$\checkmark$	336	0.238	0.278	0.191	0.247	0.159	0.256	0.379	
		720	0.307	0.330	0.197	0.253	0.185	0.284	0.405	
		Avg.	<u>0.219</u>	0.259	0.183	0.242	<u>0.153</u>	0.250	0.374	

1047Table 13: Full results of ablation on the effect of removing non-stationarity in Integrated Attention1048and Cointegrated Attention.  $\checkmark$  indicates the use of patch normalization to eliminate non-stationarity,1049while  $\times$  means non-stationarity is retained.

Integrated Attention	Cointegrated Attention		Wea	ther	So	lar	Elect	tricity	Tr
Order	Order	Length	MSE	MAE	MSE	MAE	MSE	MAE	MSE
		96	0.144	0.196	0.163	0.227	0.127	0.221	0.356
		192	0.186	0.238	0.182	0.242	0.145	0.239	0.377
1	×	336	0.241	0.283	0.192	0.246	0.162	0.257	0.390
		720	0.310	0.332	0.199	0.260	0.197	0.289	0.427
		Avg.	<u>0.220</u>	<u>0.262</u>	<u>0.184</u>	<u>0.244</u>	<u>0.158</u>	<u>0.252</u>	0.388
	1	96	0.147	0.200	0.161	0.240	0.127	0.227	0.348
		192	0.191	0.242	0.195	0.259	0.155	0.254	0.358
×	1	336	0.242	0.283	0.198	0.268	0.173	0.273	0.367
		720	0.308	0.331	0.209	0.273	0.203	0.299	0.401
		Avg.	0.222	0.264	0.191	0.260	0.165	0.263	<u>0.369</u>
		96	0.143	0.193	0.161	0.224	0.118	0.218	0.340
		192	0.185	0.235	0.177	0.237	0.142	0.237	0.343
1	2	336	0.237	0.277	0.188	0.244	0.156	0.252	0.363
		720	0.307	0.330	0.197	0.252	0.179	0.278	0.393
		Avg.	0.218	0.259	0.181	0.239	0.149	0.246	0.360
		96	0.148	0.199	0.174	0.237	0.130	0.225	0.370
		192	0.193	0.243	0.187	0.251	0.147	0.240	0.386
2	1	336	0.245	0.284	0.195	0.258	0.165	0.262	0.394
		720	0.320	0.336	0.203	0.262	0.199	0.291	0.432
		Avg.	0.227	0.266	0.190	0.252	0.160	0.255	0.396





Figure 6: Illustration of the impact and order of Integrated Attention and Cointegrated Attention in Tab. 14: 1 Integrated Attention only, 2 Cointegrated Attention only, 3 Integrated Attention followed by Cointegrated Attention, and (1) Cointegrated Attention followed by Integrated Attention, with patch downsampling replaced by upsampling. 

Integrated Attention	Cointegrated Attention		Wea	ather	So	lar	Elect	ricity	Tra	ffic
CI or CD	CI or CD	Length	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
		96	0.144	0.192	0.163	0.227	0.125	0.221	0.358	0.260
		192	0.185	0.236	0.180	0.240	0.144	0.242	0.373	0.271
CI	CI	336	0.237	0.278	0.191	0.247	0.161	0.255	0.391	0.282
		720	0.307	0.330	0.197	0.258	0.196	0.288	0.425	0.292
		Avg.	0.218	0.259	<u>0.183</u>	<u>0.243</u>	0.157	0.252	0.387	0.276
		96	0.146	0.195	0.165	0.229	0.125	0.222	0.362	0.266
		192	0.188	0.239	0.178	0.245	0.145	0.241	0.374	0.274
CD	CI	336	0.242	0.284	0.191	0.252	0.166	0.263	0.388	0.282
		720	0.310	0.331	0.201	0.261	0.205	0.293	0.423	0.296
		Avg.	<u>0.222</u>	<u>0.262</u>	0.184	0.247	0.160	0.255	0.387	0.280
		96	0.143	0.193	0.161	0.224	0.118	0.218	0.340	0.240
		192	0.185	0.235	0.177	0.237	0.142	0.237	0.343	0.250
CI	CD	336	0.237	0.277	0.188	0.244	0.156	0.252	0.363	0.257
		720	0.307	0.330	0.197	0.252	0.179	0.278	0.393	0.271
		Avg.	0.218	0.259	0.181	0.239	0.149	0.246	0.360	0.255
		96	0.146	0.197	0.162	0.229	0.125	0.221	0.352	0.254
		192	0.188	0.238	0.178	0.245	0.148	0.246	0.361	0.266
CD	CD	336	0.241	0.282	0.191	0.254	0.161	0.261	0.377	0.267
		720	0.313	0.333	0.199	0.260	0.189	0.289	0.412	0.288
		Avg.	0.222	0.263	0.183	0.247	0.156	0.254	0.376	0.269

Table 15: Full results of ablation on modeling approaches for Integrated Attention and Cointegrated Attention. "CI" denotes channel independent and "CD" denotes channel-dependent modeling.

Downsampled		Wea	ther	So	lar	Elect	ricity	Tra	ffic
Patch Number ${f M}$	Length	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	96	0.171	0.231	0.172	0.234	0.135	0.231	0.349	0.252
	192	0.207	0.261	0.187	0.252	0.158	0.255	0.358	0.259
1	336	0.257	0.298	0.196	0.257	0.182	0.278	0.382	0.269
	720	0.323	0.344	0.203	0.258	0.191	0.290	0.414	0.282
	Avg.	0.239	0.284	0.189	0.250	0.166	0.264	0.375	0.266
	96	0.147	0.201	0.168	0.231	0.118	0.218	0.340	0.240
	192	0.189	0.241	0.183	0.245	0.142	0.237	0.343	0.250
4	336	0.240	0.282	0.195	0.251	0.156	0.252	0.363	0.257
	720	0.309	0.332	0.200	0.255	0.179	0.278	0.393	0.271
	Avg.	0.221	0.264	0.186	0.246	0.149	0.246	0.360	0.255
	96	0.144	0.195	0.163	0.225	0.119	0.219	0.338	0.240
	192	0.186	0.237	0.183	0.243	0.146	0.244	0.341	0.249
8	336	0.238	0.280	0.195	0.251	0.161	0.260	0.379	0.264
	720	0.308	0.330	0.196	0.250	0.177	0.277	0.400	0.280
	Avg.	<u>0.219</u>	<u>0.261</u>	0.184	0.242	<u>0.151</u>	<u>0.250</u>	0.365	0.258
	96	0.143	0.192	0.161	0.224	0.120	0.220	0.334	0.238
	192	0.185	0.235	0.177	0.237	0.148	0.247	0.337	0.250
12	336	0.237	0.277	0.188	0.244	0.163	0.264	0.363	0.256
	720	0.307	0.330	0.197	0.252	0.176	0.286	0.387	0.268
	Avg.	0.218	0.259	<u>0.181</u>	<u>0.239</u>	<u>0.151</u>	0.254	<u>0.355</u>	0.253
	96	0.145	0.195	0.149	0.223	0.122	0.223	0.333	0.235
	192	0.187	0.238	0.175	0.236	0.149	0.249	0.343	0.254
16	336	0.241	0.282	0.187	0.244	0.165	0.265	0.354	0.256
	720	0.312	0.328	0.196	0.250	0.179	0.276	0.386	0.269
	Avg.	0.222	0.261	0.177	0.238	0.152	0.253	0.354	0.254

1161 Table 16: Full results of varying the number of downsampled patches M on forecasting perfor-1162 mance. 1163

Model	Timel	Bridge	PDF (	(2024)	Confidence
Dataset	MSE	MAE	MSE	MAE	Interval
ETTm1	$0.353 \pm 0.014$	$0.388 \pm 0.010$	$0.339 \pm 0.008$	$0.372\pm0.006$	99%
ETTm2	$0.246 \pm 0.004$	$0.310\pm0.012$	$0.252 \pm 0.003$	$0.313 \pm 0.003$	99%
ETTh1	$0.399 \pm 0.010$	$0.420 \pm 0.008$	$0.399 \pm 0.015$	$0.419 \pm 0.006$	99%
ETTh2	$0.346 \pm 0.018$	$0.394 \pm 0.015$	$0.327 \pm 0.009$	$0.376 \pm 0.010$	99%
Weather	$0.218 \pm 0.006$	$0.259 \pm 0.004$	$0.225\pm0.009$	$0.261 \pm 0.006$	99%
Electricity	$0.149 \pm 0.011$	$0.246 \pm 0.007$	$0.159\pm0.010$	$0.250 \pm 0.015$	99%
Traffic	$0.360 \pm 0.008$	$0.255 \pm 0.013$	$0.383 \pm 0.016$	$0.254 \pm 0.010$	99%
Solar	$0.181 \pm 0.002$	$0.239 \pm 0.003$	$0.205 \pm 0.008$	$0.265 \pm 0.005$	99%

Table 17: Standard deviation and statistical tests for TimeBridge and second-best method (PDF) on 1176 ETT, Weather, Electricity, Traffic, and Solar datasets. 1177

Model		TimeBridge		Г	imeMixer (2024	4)	Confidenc
Dataset	MAE	MAPE	RMSE	MAE	MAPE	RMSE	Interval
PeMS03	$14.63 \pm 0.164$	$14.21\pm0.133$	$23.10\pm0.186$	$14.63 \pm 0.112$	$14.54\pm0.105$	$23.28\pm0.128$	99%
PeMS04	$19.24\pm0.131$	$12.42\pm0.108$	$31.12\pm0.112$	$19.21 \pm 0.217$	$12.53\pm0.154$	$30.92\pm0.143$	99%
PeMS07	$20.43 \pm 0.173$	$8.42 \pm 0.155$	$33.44\pm0.190$	$20.57 \pm 0.158$	$8.62 \pm 0.112$	$33.59 \pm 0.273$	99%
PeMS08	$14.98\pm0.278$	$9.56 \pm 0.126$	$23.77\pm0.142$	$15.22 \pm 0.311$	$9.67 \pm 0.101$	$24.26\pm0.212$	99%

1178

1186 Table 18: Standard deviation and statistical tests for TimeBridge and second-best method 1187 (TimeMixer) on the PeMS dataset.

1105														
1190	Model			TimeB	ridge					Crossfor	mer (2023)			Confidence
	Dataset	ARR	AVol	MDD	ASR	CR	IR	ARR	AVol	MDD	ASR	CR	IR	Interval
1191	CSI 500 S&P 500	$\begin{array}{c} 0.285 \pm 0.033 \\ 0.326 \pm 0.022 \end{array}$	$\begin{array}{c} 0.203 \pm 0.012 \\ 0.169 \pm 0.009 \end{array}$	$\begin{array}{c} -0.196\pm 0.010 \\ -0.142\pm 0.010 \end{array}$	$\begin{array}{c} 1.405 \pm 0.016 \\ 1.927 \pm 0.017 \end{array}$	$\begin{array}{c} 1.453 \pm 0.021 \\ 2.298 \pm 0.019 \end{array}$	$\begin{array}{c} 1.317 \pm 0.019 \\ 1.842 \pm 0.016 \end{array}$	$\begin{array}{c} -0.039 \pm 0.027 \\ 0.284 \pm 0.024 \end{array}$	$\begin{array}{c} 0.163 \pm 0.009 \\ 0.159 \pm 0.011 \end{array}$	$\begin{array}{c} -0.217 \pm 0.011 \\ -0.114 \pm 0.014 \end{array}$	$\begin{array}{c} -0.238 \pm 0.021 \\ 1.786 \pm 0.014 \end{array}$	$\begin{array}{c} -0.179 \pm 0.014 \\ 2.491 \pm 0.018 \end{array}$	$\begin{array}{c} -0.350\pm 0.014 \\ 1.842\pm 0.012 \end{array}$	95% 95%
1192													-	

Table 19: Standard deviation and statistical tests for TimeBridge and second-best method (Crossformer) on the CSI 500 and S&P 500 dataset.



Figure 7: Visualization of short-term fluctuations and long-term cointegration across stock sectors. "NBFIs" represents Non-Bank Financial Institutions. The figure highlights how sectors experience short-term price volatility while maintaining long-term cointegration. 



Figure 8: Additional examples comparing intra-variate attention maps under stationary and nonstationary conditions for different patches in the Electricity dataset.



Figure 9: Additional examples comparing inter-variate attention maps between different variates under stationary and non-stationary conditions in the Solar dataset.



Figure 12: Visualization of predictions from different models on the Electricity dataset.



Figure 14: Visualization of predictions from different models on the PeMS03 dataset.