VISION-ENHANCED TIME SERIES FORECASTING BY DECOMPOSED FEATURE EXTRACTION AND COMPOSED RECONSTRUCTION

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032

034 035

043

Paper under double-blind review

ABSTRACT

Time series forecasting plays a crucial role in various domains, such as power and weather forecasting. In recent years, different types of models have achieved promising results in long-term time series forecasting. However, these models often produce predictions that lack consistency with the style of the input, resulting in reduced reliability and trust in the forecasts. To address this issue, we propose the Vision-Enhanced Time Series Forecasting by Decomposed Feature Extraction and Composed Reconstruction (VisiTER), which leverages the rich semantic information provided by the image modality to enhance the realism of the predictions. It consists of two main components: the Decomposed Time Series to Image Generation and the Composed Image to Time Series Generation. In the first component, the Decomposed Time Series Feature Extraction Model extracts periodic and trend information, which is then transformed into images using our proposed time series to vision transformation architecture. After converting the input time series into images, the resulting images are used as style features and concatenated with the previously extracted features. In the second component, we use our proposed TimeIR along with the previously obtained feature set to perform image reconstruction for the prediction part. Due to the rich information provided, the reconstructed images exhibit better consistency with the input images, which are then transformed back into time series. Extensive experiments on seven real-world datasets demonstrate that VisiTER achieves state-of-the-art prediction performance on both traditional metrics and new metrics.

1 INTRODUCTION

Time series forecasting has always been a widely studied research area (Lim & Zohren, 2021; Torres et al., 2021), with extensive applications in fields such as economics (Granger & Newbold, 2014), energy (Qian et al., 2019; Martín et al., 2010), and weather forecasting (Wu et al., 2023b). The goal of time series forecasting is to predict future values based on past observed data. Most existing time series prediction methods focus on extracting periodicity, anomalies, random fluctuations, and other features from the time series data. However, relying solely on the information provided by the time series modality is often insufficient, leading to limited prediction accuracy and an inability to efficiently capture complex relationships within the data (Ismail Fawaz et al., 2019).

One promising strategy for addressing this issue is to convert time series into images for prediction (Li 044 et al., 2024; Yang et al., 2024; Hatami et al., 2018). This is because the image modality proposes a new perspective for data modeling compared with time series, providing more rich and diverse 046 information. Nevertheless, there are several key challenges in using images for time series prediction 047 tasks. The first challenge is the difficulty in image transformation and model training. Existing 048 transformation methods directly map time series into scatter plots, which can result in information loss in the images, ultimately leading to poor performance of the subsequent image models. The complexity of image models also results in prolonged training and inference times, consuming 051 substantial GPU memory and increasing the difficulty of training. This is particularly severe for time series with many variables, which correspond to the number of channels in the image version. 052 The second is ineffective image utilization in time series: The advantages and methodologies for effectively utilizing the image modality have not been adequately explored. Some existing prediction

methods that use images rely on image generation for forecasting. However, the generated images often lack fidelity in the context of time series, as they miss crucial temporal feature information.

To address these challenges, we propose **Vision**-Enhanced Time Series Forecasting by Decomposed Feature Extraction and Composed Reconstruction (**VisiTER**). It introduces a novel framework of time series forecasting by transformation and generation between time series and vision domains, which consists of two main components: (1) the Decomposed Time Series to Vision Generation, which extracts decomposed temporal features as images for further prediction; (2) the Composed Vision to Time Series Generation, which utilizes transformed images to generate composed prediction results by image reconstruction.

063 In the first component, we propose two main modules: the Decomposed Time Series Feature 064 Extraction (DTFE) and the Time Series to Vision Transformation (T2V). DTFE decomposes the 065 original time series and leverages the respective strengths of different types of transformer-based 066 models to extract two essential temporal features for forecasting: the periodic and trend features. 067 These decomposed features help the model utilize richer temporal information and improve prediction 068 accuracy, which benefits the following image generation and reconstruction processes. Then, T2V 069 adopts a novel approach to convert time series into images, mapping data points with diminishing pixel values along the y-axis. Utilizing T2V, we transform features from DTFE into image modality, 071 improving the time series data distribution in corresponding images for better reconstruction.

In the second component, we introduce TimeIR, a novel transformer-based image reconstruction 073 model specifically designed for time series data, to generate the forecasting results by image gen-074 eration. By utilizing the periodic features, trend features extracted from DTFE, and style features 075 of the time series as priors, TimeIR can enrich the reconstruction process with valuable temporal 076 information, fully utilizing the potential of vision models. Notably, incorporating style information 077 allows the reconstructed time series to retain the same style as the input time series. To address the challenges of difficult training and inference, we redesign the model architecture and adopt a unique training strategy. Specifically, our model can perform segmented prediction for long sequences, 079 which helps to reduce the computational load. During training, we also employ a channel sampling strategy to further decrease the computational requirements. 081

- 082 In conclusion, our work makes the following key contributions:
 - We propose the VisiTER framework, which enhances time series forecasting using image reconstruction models. By first extracting temporal features as priors, we provide the image model with rich information, thereby fully leveraging the advantages of the image modality.
 - We propose DTFE, designed with various transformer-based models to extract periodic and trend features for more realistic representations. We also introduce T2V, effectively transforming time series into images, and TimeIR, which leverages priors for reconstructing time series images more efficiently.
 - VisiTER achieves state-of-the-art performance on traditional MSE and MAE metrics. Moreover, we introduce the SSIM metric to time series forecasting tasks, enabling a more comprehensive evaluation of the structural integrity of predictions. On this metric, our model also outperforms other state-of-the-art models.

2 REALTED WORK

084

085

090

092

093

095 096

097 098

2.1 TIME SERIES FORECASTING

Time series forecasting methods can be categorized mainly into those based on Recurrent Neural 101 Networks (RNNs)(Tokgöz & Ünal, 2018; Lai et al., 2018; Salinas et al., 2020), Convolutional Neural 102 Networks (CNNs)(Wang et al., 2023; Hewage et al., 2020; Livieris et al., 2020), Transformers, and 103 Multi-Layer Perceptrons (MLPs). In recent years, the Transformer model has emerged as a strong 104 contender in time series forecasting (Liu et al., 2023; Vaswani, 2017; Zhou et al., 2021; Chen et al., 105 2024). Its self-attention mechanisms effectively capture both short-term and long-term dependencies, positioning it as a leading choice for many tasks. Linear models based on MLPs have also shown 106 notable predictive results (Oreshkin et al., 2019; Challu et al., 2023; Wang et al., 2024), particularly 107 for simpler datasets, serving as a useful baseline for comparison with more complex approaches.



Figure 1: The overall architecture of VisiTER. Different colored sequences represent different variables in the time series, while the axes indicate the time series modality and the boxes represent the image modality. In the first part, we predict the periodic and trend features of the time series using the DTFE model, which are then converted into images. In the second part, we utilize TimeIR to reconstruct the images by integrating the periodic features, trend features, and style features, ultimately transforming the results back into time series data.

126 127

2.2 IMAGE RECONSTRUCTION

128 Image reconstruction (Park et al., 2003; Demoment, 1989) is a significant area in computer vision 129 and image processing, aimed at recovering high-quality images from partial or degraded inputs. This 130 task is essential in applications like medical imaging (Zhang & Dong, 2020), remote sensing (He 131 et al., 2019), and video super-resolution (Kappeler et al., 2016; Liu et al., 2022a). Various methods 132 have been proposed to tackle the image reconstruction problem, broadly categorized as follows. 133 Optimization-based methods (Fessler, 2020): These establish mathematical models for reconstruction 134 and use optimization algorithms, such as regularized optimization and dictionary learning. Deep 135 learning-based methods (Liu et al., 2020; Liang et al., 2021): Utilizing deep neural networks for end-136 to-end reconstruction (Ledig et al., 2017; Sajjadi et al., 2017; Goodfellow et al., 2020), these methods, 137 like SwinIR, effectively learn image priors from large datasets for high-quality results. GAN-based methods: Introducing the GAN framework, models like SRGAN and EnhanceNet generate more 138 realistic and natural reconstructions. 139

140 141

2.3 IMAGE TECHNIQUES IN TIME SERIES

142 In existing time series-related tasks, images are primarily used for time series classification (Li 143 et al., 2024; Dosovitskiy et al., 2021). The process typically involves converting time series data 144 into images, followed by the application of traditional image models, such as Vision Transformers, 145 to classify the images, thereby achieving classification of the time series. Additionally, there are 146 extremely few models that utilize image models for time series forecasting (Yang et al., 2024), such 147 as ViTime. These models generally decompose the time series into trend and periodic components 148 and then generate subsequent images based on the provided input. However, this type of model does 149 not fully leverage the rich information that the image modality offers, and they require long input 150 lengths to provide information.

151 152

153

3 VISION-ENHANCED TIME SERIES FORECASTING FRAMEWORK

154 155 3.1 OVERALL ARCHITECTURE

156 We propose the overall framework as shown in the Figure 1. First, we are given the time series 157 $\mathbf{X} \in \mathbb{R}^{L_1 \times N}$, where L_1 and N denote the look back length and the number of variates, which is 158 input to the DTFE in Part I to obtain the trend feature $\mathbf{T} \in \mathbb{R}^{L_2 \times N}$ and periodic feature $\mathbf{P} \in \mathbb{R}^{L_2 \times N}$ 159 of the time series, where L_2 denotes the prediction length. Next, \mathbf{T} , \mathbf{P} and \mathbf{x} are input into the T2V 160 framework to be converted into image formats, and then they are shuffled based on different variables 161 in Part II. For the same variable i, we concatenate the \mathbf{P}_i , \mathbf{T}_i , and \mathbf{X}_i . If L_1 is not equal to L_2 , we also need to align them before concatenation. The subsequent TimeIR model operates on a univariate basis, meaning we perform reconstruction for each variable individually. For a given time series, we
segment it into N images for reconstruction, with each image corresponding to one of the N variables.
Ultimately, we combine the reconstruction results back into a single image, which is then transformed
back into the time series format.

166 167

168

3.2 DECOMPOSED TIME SERIES TO IMAGE GENERATION

169 In this section, we first predict the peri-170 odic information and trend features from 171 the time series. We observe that the trend 172 features of a variable are more susceptible 173 to the influence of other variables, while its 174 periodic information is less affected. There-175 fore, we propose the Decomposed Time 176 Series Feature Extraction Model (DTFE), which employs a decomposed architecture 177 using different Transformers to predict the 178 periodic and trend features, resulting in 179 more accurate outcomes. 180

181 Next, we convert the predicted periodic in-182 formation, trend information, and the input time series into images. The original 183 direct mapping method transforms time se-184 ries into scatter plots, which, while precise, 185 has limited applicability for image reconstruction models due to the reduced number 187 of activated pixels. This results in a sparse 188 representation that does not fully leverage 189 the potential of the image modality, lead-190



Figure 2: Overall architecture of DTFE. Different variables are treated as tokens within the Transformer framework to predict trends, while multiple time steps of a single variable are considered as a token for predicting periodicity. Each row of the time series represents a distinct variable, while each differently colored box signifies a different token within the Transformer backbone. The specific architecture of the Transformer framework is illustrated in Figure 11 of the Appendix B.

ing to insufficient information for the reconstruction models to utilize effectively. To address this,
 we propose a new transformation method, T2V. We diffuse the scatter points, making the image continuous and more suitable for image models.

193 194 195

3.2.1 DECOMPOSED TIME SERIES FEATURE PREDICTION

196 In DTFE, for predicting the trend component, it is necessary to consider the inter-channel interactions 197 more extensively, as the trend of one channel can be influenced by the trends of other channels. In contrast, for the periodic component, the influence of the periodicities in other channels on a 199 given channel is relatively small, and hence, a more channel-independent approach can be adopted. 200 Therefore, for the trend prediction model, we adopt an inverted transformer structure to explicitly 201 capture the cross-channel dependencies. On the other hand, for the periodic component prediction, we use a patch-based transformer architecture, where we treat a few adjacent time steps as a patch 202 and predict the periodicity based on this sequence of patches. By leveraging the strengths of both 203 single-channel and multi-channel modeling approaches, and by tailoring the model structures to the 204 specific characteristics of trend and periodicity, the proposed DTFE succeeds to achieve robust and 205 accurate time series features prediction. The basic framework of DTFE is shown in the Figure 2. 206

When training DTFE, it is essential to compute the loss functions for the periodic Transformer backbone and the trend Transformer backbone separately for predicting the periodic and trend components, and then sum them together. For each periodic Transformer, denoted as P-Trans(·), we use L_{period} to represent its MSE loss function, and for the trend Transformer, denoted as T-Trans(·), we use L_{trend} to represent its MSE loss function. Additionally, we employ the Discrete Fourier Transform (DFT) to decompose the ground truth of forecasting result y into its period and trend components. Thus, we can derive the training losses for period and trend extraction:

214

$$L_{\text{period}} = \sum_{\mathbf{X}_i \in B} \left\| \text{P-Trans}(\mathbf{X}_i) - \text{DFT}(\mathbf{y}_i)_{\text{period}} \right\|_2^2 \tag{1}$$



Figure 3: The top row of the image illustrates the overall framework of T2V, which begins with normalization, followed by direct mapping to images, and then expands on both sides of the y-axis. The second row presents the overall framework of V2T, where the maximum value of each column is selected, mapped to a time series, and finally restored.

$$L_{\text{trend}} = \sum_{\mathbf{X}_i \in B} \left\| \left| \text{T-Trans}(\mathbf{X}_i) - \text{DFT}(\mathbf{y}_i)_{\text{trend}} \right| \right\|_2^2,$$
(2)

where B represents the batch used for training, y_i denotes the ground truth of forecasting result of the input time series \mathbf{X}_i . Thus, we can achieve the overall loss $L_{\text{DTFE}} = L_{\text{period}} + L_{\text{trend}}$.

3.2.2 TIME SERIES TO IMAGE TRANSFORMATION

Given a time series $\mathbf{x} \in R^{L \times N}$, we need to transform it into an image $\mathbf{F} \in R^{H \times L \times N}$, where L 235 corresponds to the length of the time series to be reconstructed, H is a hyperparameter that represents 236 the height of the image after the transformation into a visual format and N denotes the number of 237 variables involved. First, we normalize the time series x to have a variance of 1 and a mean of 0, 238 to facilitate the subsequent operations: $\mathbf{x}' = \frac{\mathbf{x} - \mu(\mathbf{x})}{\sigma(\mathbf{x})}$, where $\mu(\mathbf{x})$ and $\sigma(\mathbf{x})$ represent the mean and 239 standard deviation of x respectively. For each channel, we calculate the maximum value and divide 240 the values in that channel by the maximum value. This maps the entire time series to the range of [-1, 241 1]: $\mathbf{x}'' = \frac{\mathbf{x}'}{\max(\mathbf{x}', axis=0)}$. To control the range of time series within [0, H], we can multiply \mathbf{x}'' by $\frac{H}{2}$ 242 243 and then add $\frac{H}{2}$: $\tilde{\mathbf{x}} = \mathbf{x}'' \times \frac{H}{2} + \frac{H}{2}$, which ensures that each value in the time series corresponds to 244 the vertical coordinate values after conversion into an image format. Based on the normalized values, 245 we can determine which pixels in the image F need to be activated:

223

224

225

226 227 228

229 230

231

232 233

234

248

251

256 257

258 259

260

261 262

263

$$\mathbf{F}_{i,j,c} = \begin{cases} 1, & \text{if } \tilde{\mathbf{x}}_{j,c} = i \\ 0, & \text{otherwise} \end{cases}$$
(3)

249 In this context, $\mathbf{F}_{i,j,c}$ denotes the pixel in \mathbf{F} at the *i*-th row, *j*-th column, and *c*-th channel, while 250 $\tilde{\mathbf{x}}_{j,c}$ represents the value of the c-th variable at the j-th time point in the time series $\tilde{\mathbf{x}}$. To help the reconstruction model perform better, we modify the scatter plot representation by extending the 252 values on the Y-axis. Specifically, we enhance the neighboring pixels on either side of each activated 253 pixel in a given column, with the value decreasing as the distance from the activated pixel increases. The extension range is $[0,\lambda]$, where λ is a hyperparameter to control extension ranges. The specific 254 formula is as follows, where k represents the current extension distance: 255

$$\mathbf{F}_{i,j,c} = \begin{cases} 1 - \frac{k}{\lambda}, & \text{if } \tilde{\mathbf{x}}_{j,c} = i \pm k \quad and \quad k \le \lambda \\ 0, & \text{otherwise} \end{cases}$$
(4)

We transform the periodic features \mathbf{P} , trend features \mathbf{T} and the input \mathbf{x} into images using T2V respectively, resulting in \mathbf{P}_{image} , \mathbf{T}_{image} and \mathbf{X}_{image} .

3.3 **COMPOSED IMAGE TO TIME SERIES GENERATION**

264 Based on the decomposed component images after prediction, we construct an image reconstruction 265 model aimed at generating a complete image of the predicted results. To maintain consistency 266 in style between the reconstructed predicted image and the original sequence, we introduce the image of the original sequence as additional style information. Consequently, we propose TimeIR, 267 a temporal image reconstruction model, which first concatenates these three parts of input for the 268 image reconstruction, then performs the reconstruction, and finally converts the results back into the 269 time series format.



Figure 4: Alignment method of the model under different prediction lengths. When L_1 is greater than L_2 , the length of x is truncated. When L_1 equals L_2 , no special processing is applied. When L_1 is less than L_2 , a sliding window approach is used to segment and truncate P and T for prediction.

3.3.1 COMPOSED IMAGE RECONSTRUCTION

276

277

278 279

280

310

314

321

322

281 TimeIR is a univariate reconstruction model, so the following discussion pertains to the reconstruction 282 of a single variable within a time series. The inputs we have received are $\mathbf{X}_{image} \in R^{H \times L_1 \times N}$, $\mathbf{P}_{image} \in R^{H \times L_2 \times N}$, and $\mathbf{T}_{image} \in R^{H \times L_2 \times N}$, which represent the original time series \mathbf{x} as images, as well as the images depicting the periodicity and trends elucidated by the DTFE. The 283 284 285 length of \mathbf{X}_{image} is L_1 , while the lengths of \mathbf{P}_{image} and \mathbf{T}_{image} are L_2 . Next, we concatenate these three components to obtain the input $\mathbf{I} \in \mathbb{R}^{H \times L \times 3}$, where the three components represent the three 286 channels of I. To perform the concatenation operation, it is necessary to align the lengths of the 287 sequences. We handle different scenarios accordingly: when L_1 equals L_2 , we proceed directly. 288 When L_1 is greater than L_2 , we truncate \mathbf{X}_{image} to the length of L_2 since \mathbf{X}_{image} serves merely as 289 a style feature, and its length does not impact the information. When L_1 is less than L_2 , we maintain 290 a sliding window on P_{image} and T_{image} for learning. During each reconstruction, we select the 291 periodic features and trend features within the sliding window, while keeping the style feature fixed 292 at \mathbf{X}_{image} , as the style of the time series does not change with the sliding window selection. The 293 reconstruction method is illustrated in Figure 4, and additional alignment details can be found in Appendix B.2. Thus, we can obtain style, trend, and periodicity priors, which we then concatenate to 295 form the input to the TimeIR model. 296

The TimeIR model is comprised of three main components. First, a shallow feature extraction module 297 utilizing a CNN network block is employed. This is followed by a deep feature extraction module, 298 which consists of multiple Time Series Swin Transformer Blocks (TSTB). Each TSTB is composed 299 of several Swin Transformer Blocks. The key characteristic of the Swin Transformer used in TimeIR 300 is its utilization of an overlapping patch embedding scheme. This approach allows the model to 301 better focus on the fine-grained details of the input data, which in this case is the time series. Finally, 302 the deep feature representations are passed through a convolutional layer that reduces the channel 303 dimension to 1, producing the final time series reconstruction. The overall architecture of TimeIR 304 illustrated in Figure 10.

Since images and time series are different modalities, we need to convert the ground truth into images for training. We use MSE as the loss function between the predictions from TimeIR and the ground truth converted into images, as shown in the equation below.

$$L_{\text{TimeIR}} = \sum_{\mathbf{X}_i \in B} \left\| \left| \text{TimeIR}(\mathbf{I}_i) - \text{T2V}(\mathbf{y}_i) \right| \right\|_2^2$$
(5)

In this equation, L_{TimeIR} represents the loss function of training TimeIR, B represents the batch used for training, \mathbf{y}_i denotes the ground truth of forecasting result of \mathbf{X}_i , and \mathbf{I}_i refers to the combined information inputted into TimeIR for the time series \mathbf{X}_i .

315 3.3.2 IMAGE TO TIME SERIES TRANSFORMATION

The process of V2T is the inverse of the T2V procedure. Specifically, V2T converts the reconstruction results from TimeIR back into a time series, serving as the final prediction output. The first step is to identify the maximum value in each column of the image. The pixel corresponding to this maximum value is then set to 1, while all other pixels in that column are assigned a value of 0.

$$\mathbf{F}_{i,j,c} = \begin{cases} 1, & \text{if } \mathbf{F}_{i,j,c} = max(\mathbf{F}_{i,j,c}, axis = 0) \\ 0, & \text{otherwise} \end{cases}$$
(6)

Then, after rescaling the values to the range [-1, 1]. and normaling the values, we can get the predicted time series \hat{y} . The complete T2V and V2T processes are illustrated in Figure 3.

³²⁴ 4 EXPERIMENTS

326

327

328

335

343

Datasets The datasets comprises 7 collections: ETT dataset (including 4 subsets:ETTh1, ETTh2, ETTm1, and ETTm2), Weather, Exchange, and Electricity datasets (Wu et al., 2021; Li et al., 2021). A detailed description of the dataset can be found in the Appendix A.1.

Baseline We will compare VisiTER with 11 latest baselines, including PatchTST (Nie et al., 2023), iTransformer (Liu et al., 2023), Crossformer (Zhang & Yan, 2023), Autoformer (Wu et al., 2021), SparseTSF (Lin et al., 2024), TimesNet (Wu et al., 2023a), DLinear (Zeng et al., 2023), FEDformer (Zhou et al., 2022), Non-Stationary Transformers (Liu et al., 2022c), SCINet (Liu et al., 2022b), and TiDE (Das et al., 2023).

Main results The comprehensive forecasting results are listed in Table 1, with the best results marked in red and the second-best blue, while the visual results are displayed in Figure 5. In the case where the prediction length is greater than 96, our VisiTER model uses the model trained on the prediction length of 96 for the dataset to perform zero-shot prediction. We can see that our model achieve SOTA results on multiple datasets, especially performing particularly well on datasets with fewer variables. Furthermore, we find that the image reconstruction module is able to maintain low traditional MSE and MAE evaluation metrics, while also preserving the reconstruction style.

Table 1: Multivariate forecasting results with prediction lengths $S \in \{96, 192, 336, 720\}$ for others and fixed lookback length T = 96. Results are averaged from all prediction lengths. TimeIR performs **ZERO-SHOT** inference when predicting lengths of $\{192, 336, 720\}$. Avg means further averaged by subsets. Full results are listed in Table 5.

Models	Visi (Ou	FER irs)	iTrans (20	former)23)	Spars (20	eTSF 24)	Patch (20	nTST 23)	Cross (20	former)23)	T (2	DE 023)	Tim (20	esNet 23a)	DLi (20	near 23)	SCI (202	Net 22b)	FEDf (20	ormer 22)	Stati (20	onary 22c)
Metric	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTm1	0.386	0.371	0.407	0.410	0.416	0.408	0.387	<u>0.400</u>	0.513	0.496	0.419	0.419	0.400	0.406	0.403	0.407	0.485	0.481	0.448	0.452	0.481	0.456
ETTm2	0.279	0.323	0.288	0.332	0.288	0.329	0.281	0.326	0.757	0.610	0.358	0.404	0.291	0.333	0.350	0.401	0.571	0.537	0.305	0.349	0.306	0.347
ETTh1	0.431	0.428	0.454	0.447	0.440	0.429	0.469	0.454	0.529	0.522	0.541	0.507	0.458	0.450	0.456	0.452	0.747	0.647	0.440	0.460	0.570	0.537
ETTh2	0.368	0.398	0.383	0.407	0.383	<u>0.402</u>	0.387	<u>0.407</u>	0.942	0.684	0.611	0.550	0.414	0.427	0.559	0.515	0.954	0.723	0.437	0.449	0.526	0.516
ECL	0.194	0.285	0.178	0.270	0.224	0.297	0.205	0.290	0.244	0.334	0.251	0.344	0.192	0.295	0.212	0.300	0.268	0.365	0.214	0.327	0.305	0.349
Exchange	0.334	0.389	0.360	0.403	0.361	0.408	0.367	0.404	0.940	0.707	0.370	0.413	0.416	0.443	0.354	0.414	0.750	0.626	0.519	0.429	0.193	0.296
Weather	0.254	0.277	0.258	0.278	0.276	0.294	0.259	0.281	0.259	0.315	0.271	0.320	0.259	0.287	0.265	0.317	0.292	0.363	0.309	0.360	0.288	0.314

356 357 358

359

360 361

4.1 TRAINING STRATEGY

To avoid the excessively high computational demands of traditional image reconstruction models, we have adopted several strategies during the training process of TimeIR.

Firstly, for a given dataset, we only train the TimeIR for the case of 96-length prediction and 96-length 364 input. When the prediction length is greater than 96, we directly use this trained model for zero-shot prediction, and have achieved very good results. Secondly, since some of the datasets have a large 366 number of variables, we employ a channel sampling operation during training. Specifically, for each 367 training batch, we select a different set of variables to train on. Lastly, we employed a sliding window 368 strategy to align the input, ensuring that when the prediction length exceeds the input length, the size 369 of TimeIR's input remain fixed, so the GPU memory consumption is independent of the prediction 370 length. This allows our model to handle even very long prediction lengths without running into 371 issues of infeasibility. The training strategy for the entire VisiTER framework is described in detail in 372 Appendix A.2.

373 374 4.2 SSIM

375

The SSIM is a metric used to measure the similarity between images (Wang et al., 2004). Unlike traditional pixel-level error metrics, such as MSE, SSIM places greater emphasis on the structural information and perceptual quality of the images. The fundamental idea behind SSIM is to evaluate



Figure 5: Comparison of visual results for time series reconstruction. The experiment focuses on predicting 96 steps with a lookback length of 96, using the ETTh2 dataset, where results from different variables are sampled. Each row represents a sample and each column represents a method. The first half of each image displays the provided time series, while the second half shows the predictions from various models. The last column represents the ground truth. Please zoom in for a closer view.

Table 2: The SSIM results of our model compared to other baselines with prediction lengths $S \in \{96, 192, 336, 720\}$ for others and fixed lookback length T = 96. Results are averaged from all prediction lengths, with the best results highlighted in bold. TimeIR performs **ZERO-SHOT** inference when predicting lengths of $\{192, 336, 720\}$. Higher SSIM values indicate better performance. The complete SSIM results are presented in Table 6.

Model	ETTm1	ETTm2	ETTh1	ETTh2	ECL	Exchange	Weather
Autoformer (Wu et al., 2021)	0.3762	0.4836	0.3867	0.4101	0.6051	0.6168	0.4377
Crossformer (Zhang & Yan, 2023)	0.3800	0.3822	0.4047	0.2620	0.6788	0.3417	0.4199
iTransformer (Liu et al., 2023)	0.4329	0.5025	0.4171	0.4163	0.7015	0.6443	0.5837
TimesNet (Wu et al., 2023a)	0.4292	0.5119	0.4072	0.4013	0.6486	0.6203	0.5720
VisiTER (ours)	0.4553	0.5384	0.4606	0.4395	0.6868	0.6597	0.5941

the quality of two images by comparing their similarity in terms of luminance, contrast, and structure.
For time series images generated by our T2V method, brighter pixel values represent a higher
likelihood, meaning that luminance, contrast, and structure all reflect the authenticity of the time
series. SSIM values range from -1 to 1, where 1 indicates that the two images are identical, while
0 or negative values suggest a low level of similarity. A detailed introduction to SSIM and its
corresponding formulas can be found in Appendix A.4.

We evaluated the performance of our model on various datasets using the SSIM metric. For comparison, we selected four state-of-the-art baselines: Autoformer, Crossformer, TimesNet, and iTransformer. Given that the λ in T2V can impact the SSIM values, we have set λ =100 to facilitate the comparison. The results are presented in Table 2. It can be observed that our VisiTER model outperforms the other models by a significant margin in terms of SSIM across multiple datasets. This demonstrates that the time series predicted by our model exhibits stronger fidelity and structural similarity compared to the baselines.

432 4.3 ABLATIONS 433

435

456

434 4.3.1 EFFECTIVENESS OF DTFE STRUCTURE.

To evaluate the effectiveness of our proposed DTFE architecture, we have conducted comparisons on 436 the ETTh2 dataset against other Transformer-based models, including PatchTST and iTransformer. 437 By directly summing the predicted trends and periodic components from DTFE, we obtain an 438 intermediate result of the VisiTER. Additionally, we have included another variant of the DTFE, 439 which was trained using a different approach and the training strategy are presented in Figure 9. 440 In our current training method, we separately train a periodic feature extractor and a trend feature 441 extractor, and the loss function is the sum of the MSE of the periodic loss and the trend loss. The 442 alternative training method involves directly adding the predicted periodic and trend components, 443 and then computing the MSE loss between the aggregated prediction and the ground truth.

444 The results are presented in the Table 3, which 445 the strategy A denotes calculating the loss sepa-446 rately, while the strategy B refers to calculating 447 the loss after summing the components. We 448 can observe that the DTFE trained using our 449 composite-architecture outperforms the single-450 architecture Transformer models. Furthermore, 451 the DTFE trained using strategy A demonstrates superior performance compared to the tradi-452 tional aggregation-based training strategy B. 453 This indicates that our model not only possesses 454 greater interpretability but also achieves better results. 455

Table 3: Results of the ablation study on the effectiveness of DTFE on ETTh2.

Models	96	192	336	720	Avg
iTransformer	0.297	0.380	0.428	0.427	0.383
PatchTST	0.302	0.388	0.426	0.431	0.387
DTFE(Strategy B)	0.292	0.371	0.409	0.425	0.374
DTFE(Strategy A)	0.286	0.366	0.404	0.420	0.370

4.3.2 Ablation of λ



Figure 6: The first row is the visualization of the same time series under different λ . Others are the results obtained from training with different λ values during the training process, with the last column representing the ground truth. For better visibility, the visualized results only capture the prediction portion in this figure.

482 483

478

479

480

481

In order to compare the influence of different degree of expansion in T2V (λ) on the performance of 484 TimeIR, We conducted experiments with an input length of 96 and a prediction length of 96 on the 485 ETTh2 dataset. In this study, we have fixed the DTFE and ensured that all the inputs are the same,







Figure 8: Results of the ablation study on the zero-shot capability.

while also maintaining consistent training hyperparameters. Furthermore, the SSIM values reported here are measured with $\lambda = 100$ for the sake of uniform comparison. The visualization results are 502 presented in Figure 6, and the experimental results are shown in Figure 7.

It can be observed that as the value of λ increases, both the MSE and SSIM metrics exhibit better 504 performance. However, the visual inspection of a few sample cases reveals that when λ is smaller, 505 the reconstructed models tend to be more extreme, while larger values of λ result in smoother 506 reconstructions. In other words, smaller values of λ generate results that are more stylistically similar 507 to the real data, but this similarity is not necessarily reflected in the numerical metrics. Thus, λ can 508 be considered a hyperparameter that alters the model's ability to fuse styles. A smaller λ leads to 509 reconstructed time series that are more similar in style to the input, but at the cost of ignoring the overall coherence of the time series. Conversely, a larger λ results in less influence from the style of 510 the input sequence. 511

512

495

498 499 500

501

4.3.3 ZERO-SHOT ANALYSIS 513

514 In our experiments, to reduce the computational cost of training, we have only conducted training 515 the TimeIR on the same dataset for sequence length 96 with 96-step prediction. For other cases, 516 we have adopted a zero-shot approach. For prediction lengths greater than the sequence length, our 517 training strategy involves randomly selecting the starting position of the sliding window for training. 518 Additionally, we keep the DTFE model fixed while maintaining the same hyperparameters for the 519 others. The datasets used in this experiment are Weather and ETTh1, and the results are presented in 520 the Figure 8.

521 The results indicate that the zero-shot performance surpasses that of the standard training approach, 522 and this difference becomes more pronounced as the prediction length increases. The underlying 523 reason for this seemingly counterintuitive result is related to the design strategy of our model. Our 524 model is designed to use the input sequence to provide the style features of the time series. However, 525 when the prediction length is significantly longer than the sequence length, the style or characteristics 526 may change. When the sliding window is close to the starting point, the style is more similar, but as the window moves further away, the style can become less similar. This can lead to training instability 527 and poorer performance compared to the zero-shot approach. 528

529 530

5 CONCLUSION

531

532 We propose the Vision-Enhanced Time Series Forecasting by Decomposed Feature Extraction and 533 Composed Reconstruction Framework (VisiTER) that leverages image reconstruction techniques for 534 time series forecasting. The framework consists of two main components: the Decomposed Time 535 Series to Image Generation and the Composed Image to Time Series Generation. It successfully 536 integrates the image modality into time series forecasting. By supplementing the time series data with rich information from the image modality, the prediction results become more reliable and accurate. In future work, we hope to see a growing interest in exploring the use of image modalities within the 538 field of time series forecasting. This integration could uncover new avenues for enhancing predictive models and improving performance.

540 REFERENCES

542 543 544	Cristian Challu, Kin G Olivares, Boris N Oreshkin, Federico Garza Ramirez, Max Mergenthaler Canseco, and Artur Dubrawski. Nhits: Neural hierarchical interpolation for time series forecasting. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , 2023.
545 546 547	Peng Chen, Yingying Zhang, Yunyao Cheng, Yang Shu, Yihang Wang, Qingsong Wen, Bin Yang, and Chenjuan Guo. Pathformer: Multi-scale transformers with adaptive pathways for time series forecasting. <i>ICLR</i> , 2024.
549 550	Abhimanyu Das, Weihao Kong, Andrew Leach, Rajat Sen, and Rose Yu. Long-term forecasting with tide: Time-series dense encoder. <i>arXiv preprint arXiv:2304.08424</i> , 2023.
551 552 553	Guy Demoment. Image reconstruction and restoration: Overview of common estimation structures and problems. <i>IEEE Transactions on Acoustics, Speech, and Signal Processing</i> , 1989.
554 555 556	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. <i>ICLR</i> , 2021.
557 558 559	Jeffrey A Fessler. Optimization methods for magnetic resonance image reconstruction: Key models and optimization algorithms. <i>IEEE signal processing magazine</i> , 2020.
560 561 562	Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. <i>Communications of the ACM</i> , 2020.
563 564 565	Clive William John Granger and Paul Newbold. <i>Forecasting economic time series</i> . Academic press, 2014.
566 567 568	Nima Hatami, Yann Gavet, and Johan Debayle. Classification of time-series images using deep convolutional neural networks. In <i>Tenth international conference on machine vision (ICMV 2017)</i> , 2018.
569 570 571	Wei He, Naoto Yokoya, Longhao Yuan, and Qibin Zhao. Remote sensing image reconstruction using tensor ring completion and total variation. <i>IEEE Transactions on Geoscience and Remote Sensing</i> , 2019.
573 574 575	Pradeep Hewage, Ardhendu Behera, Marcello Trovati, Ella Pereira, Morteza Ghahremani, Francesco Palmieri, and Yonghuai Liu. Temporal convolutional neural (tcn) network for an effective weather forecasting using time-series data from the local weather station. <i>Soft Computing</i> , 2020.
576 577 578 579	Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. <i>Data mining and knowledge discovery</i> , 2019.
580 581	Armin Kappeler, Seunghwan Yoo, Qiqin Dai, and Aggelos K Katsaggelos. Video super-resolution with convolutional neural networks. <i>IEEE transactions on computational imaging</i> , 2016.
582 583 584	Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. <i>ICLR</i> , 2021.
585 586 587	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> , 2018.
588 589 590 591	Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , 2017.
592 593	Jianxin Li, Xiong Hui, and Wancai Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. <i>arXiv: 2012.07436</i> , 2021.

594 595 596	Zekun Li, Shiyang Li, and Xifeng Yan. Time series as images: Vision transformer for irregularly sampled time series. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
597 598 599	Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Im- age restoration using swin transformer. In <i>Proceedings of the IEEE/CVF international conference</i> <i>on computer vision</i> , 2021.
600 601 602	Bryan Lim and Stefan Zohren. Time-series forecasting with deep learning: a survey. <i>Philosophical Transactions of the Royal Society A</i> , 2021.
603 604	Shengsheng Lin, Weiwei Lin, Wentai Wu, Haojun Chen, and Junjie Yang. Sparsetsf: Modeling long-term time series forecasting with 1k parameters. <i>arXiv preprint arXiv:2405.00946</i> , 2024.
605 606 607	Bo Liu, Jiupeng Tang, Haibo Huang, and Xi-Yun Lu. Deep learning methods for super-resolution reconstruction of turbulent flows. <i>Physics of fluids</i> , 2020.
608 609 610	Hongying Liu, Zhubo Ruan, Peng Zhao, Chao Dong, Fanhua Shang, Yuanyuan Liu, Linlin Yang, and Radu Timofte. Video super-resolution based on deep learning: a comprehensive survey. <i>Artificial Intelligence Review</i> , 2022a.
611 612 613	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. <i>NeurIPS</i> , 2022b.
614 615	Yong Liu, Haixu Wu, Jianmin Wang, and Mingsheng Long. Non-stationary transformers: Rethinking the stationarity in time series forecasting. <i>NeurIPS</i> , 2022c.
616 617 618 619	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>arXiv preprint arXiv:2310.06625</i> , 2023.
620 621	Ioannis E Livieris, Emmanuel Pintelas, and Panagiotis Pintelas. A cnn–lstm model for gold price time-series forecasting. <i>Neural computing and applications</i> , 2020.
622 623 624 625	Luis Martín, Luis F Zarzalejo, Jesus Polo, Ana Navarro, Ruth Marchante, and Marco Cony. Prediction of global solar irradiance based on time series analysis: Application to solar thermal power plants energy production planning. <i>Solar Energy</i> , 2010.
626 627	Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. <i>ICLR</i> , 2023.
628 629 630 631	Boris N Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis expansion analysis for interpretable time series forecasting. <i>arXiv preprint arXiv:1905.10437</i> , 2019.
632 633	Sung Cheol Park, Min Kyu Park, and Moon Gi Kang. Super-resolution image reconstruction: a technical overview. <i>IEEE signal processing magazine</i> , 2003.
634 635 636 637	Zheng Qian, Yan Pei, Hamidreza Zareipour, and Niya Chen. A review and discussion of decomposition-based hybrid models for wind energy forecasting applications. <i>Applied energy</i> , 2019.
638 639 640	Mehdi SM Sajjadi, Bernhard Scholkopf, and Michael Hirsch. Enhancenet: Single image super- resolution through automated texture synthesis. In <i>Proceedings of the IEEE international confer-</i> <i>ence on computer vision</i> , 2017.
642 643	David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. <i>International Journal of Forecasting</i> , 2020.
644 645	Alper Tokgöz and Gözde Ünal. A rnn based time series approach for forecasting turkish electricity load. In 2018 26th Signal processing and communications applications conference (SIU), 2018.
647	José F Torres, Dalil Hadjout, Abderrazak Sebaa, Francisco Martínez-Álvarez, and Alicia Troncoso. Deep learning for time series forecasting: a survey. <i>Big Data</i> , 2021.

648 649	A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
650	Huiqiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao. Micn: Multi-scale
651	local and global context modeling for long-term series forecasting. In <i>The eleventh international</i>
652	conference on learning representations, 2025.
653	Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y Zhang, and
655	Jun Zhou. Timemixer: Decomposable multiscale mixing for time series forecasting. <i>arXiv preprint</i>
656	arXiv:2405.14616, 2024.
657 658	Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. <i>IEEE Transactions on Image Processing</i> , 2004.
659 660	Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with Auto-Correlation for long-term series forecasting. <i>NeurIPS</i> , 2021.
661	Hairy Wy Tangga Hy Yang Lin Hang They Lianmin Wang and Mingshang Long. Timesnate
662 663	Temporal 2d-variation modeling for general time series analysis. <i>ICLR</i> , 2023a.
664 665	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> , 2023b.
667	Luoxiao Yang, Yun Wang, Xingi Fan, Israel Cohen, Yue Zhao, and Zijun Zhang. Vitime: A visual
668	intelligence-based foundation model for time series forecasting. arXiv preprint arXiv:2407.07311,
669	2024.
670	Ailing Zeng, Muxi Chen, Lei Zhang, and Oiang Xu. Are transformers effective for time series
671	forecasting? AAAI, 2023.
672	Hei Miss Zhang and Din Dang. A maring in dean learning in medical image reconstruction. Journal
673	of the Operations Research Society of China, 2020.
675	
676	Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. <i>ICLR</i> , 2023.
677	Haovi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
670	Informer: Beyond efficient transformer for long sequence time-series forecasting. In <i>Proceedings</i>
680	of the AAAI conference on artificial intelligence, 2021.
681	Tian Zhou, Ziging Ma, Oingsong Wen, Xue Wang, Liang Sun, and Rong Jin, FEDformer: Frequency
682	enhanced decomposed transformer for long-term series forecasting. <i>ICML</i> , 2022.
683	
684	
685	
686	
687	
688	
689	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

702 IMPLEMENTATION DETAILS А 703

704 A.1 DATASET DESCRIPTIONS 705

706 We conduct experiments on seven real-world datasets to evaluate the performance of the proposed VisiTER, including: 707

- ETTh1,ETTh2: This dataset (Kim et al., 2021) contains electricity transformer data recorded hourly from July 2016 to July 2018. It includes various operational factors that influence electricity consumption and transformer performance.
 - ETTm1,ETTm2: This dataset (Kim et al., 2021) consists of electricity transformer measurements recorded every 15 minutes over the same time span. The higher frequency of data points allows for analysis of more granular trends and patterns in electricity usage.
- Exchange (Wu et al., 2021): This dataset compiles daily exchange rate panel data from eight countries, covering the period from 1990 to 2016. It includes various currency pairs, providing a rich resource for studying financial time series and the effects of economic events on exchange rates.
- 718 • Weather: This dataset (Wu et al., 2021) features 21 meteorological factors, such as tem-719 perature, humidity, and wind speed, collected every 10 minutes from the Weather Station 720 of the Max Planck Biogeochemistry Institute in 2020. It serves as a vital resource for understanding climate patterns and their relationship with other time-dependent variables. 722
 - ECL: This dataset records (Wu et al., 2021) hourly electricity consumption data from 321 clients. It provides insights into consumer behavior and demand patterns, making it useful for load forecasting and energy management studies.

We follow the same data processing and train-validation-test set split protocol used in TimesNet, 726 ensuring a strict chronological order to prevent data leakage. For forecasting settings, we fix the 727 lookback series length at 96 for all datasets, while the prediction length varies among {96, 192, 336, 728 720. Detailed information about the datasets is provided in Table 4. 729

730 731

708

709

710

711

712

713

714

715

716

717

721

723

724

725

A.2 TRAINING PROCESS

732 In the training of our entire framework, we begin by training the DTFE Model. Subsequently, we 733 freeze its parameters and proceed to train the TimeIR model. Specifically, we focus on training the 734 portion of the dataset that corresponds to a prediction length of 96. For any other prediction lengths, 735 we directly employ the model for zero-shot inference. 736

737 A.3 EXPERIMENT DETAILS 738

739 All experiments were conducted on an NVIDIA RTX 4090. We employed the ADAM optimizer and MSE as the loss function. The learning rate for all experiments was set at 0.0001. In Part 1, during 740 the training of DTFE, a batch size of 32 was selected, while in Part 2, the batch size for training 741 TimeIR was set to 5. Both Transformer blocks in DTFE consist of a single layer, whereas the TSTB 742 in TimeIR has two layers. 743

744 745

Table 4: Dataset detailed descriptions. The dataset size is organized in (Train, Validation, Test).

Dataset	Dim	Series Length	Dataset Size	Frequency	Information
ETTm1	7	{96, 192, 336, 720}	(34465, 11521, 11521)	15min	Temperature
ETTm2	7	{96, 192, 336, 720}	(34465, 11521, 11521)	15min	Temperature
ETTh1	7	{96, 192, 336, 720}	(8545, 2881, 2881)	15 min	Temperature
ETTh2	7	{96, 192, 336, 720}	(8545, 2881, 2881)	15 min	Temperature
Exchange	8	{96, 192, 336, 720}	(5120, 665, 1422)	Daily	Economy
Electricity	321	{96, 192, 336, 720}	(18317, 2633, 5261)	Hourly	Electricity
Weather	21	{96, 192, 336, 720}	(36792, 5271, 10540)	10 min	Weather



Figure 9: Different training strategies for DTFE. (a) Utilizing different types of Transformer models
 to separately predict periodicity and trends. (b) A conventional model ensemble method, where
 predictions are directly summed to compute the loss function.

768 A.4 METRIC DETAILS

Regarding metrics, we utilize the mean square error (MSE) and mean absolute error (MAE) for forecasting. The calculations of these metrics are:

$$MSE = \sum_{i=1}^{F} (\mathbf{X}_i - \widehat{\mathbf{X}}_i)^2, \qquad MAE = \sum_{i=1}^{F} |\mathbf{X}_i - \widehat{\mathbf{X}}_i|,$$

where $\mathbf{X}, \widehat{\mathbf{X}} \in \mathbb{R}^{F \times C}$ are the ground truth and prediction results of the future with F time pints and C dimensions. \mathbf{X}_i means the *i*-th future time point. At the same time, we convert the time series into images for Structural Similarity Index (SSIM) testing (Wang et al., 2004). The SSIM relies on three relatively autonomous components: luminance, contrast, and structures. It is widely recognized and better aligned with the requirements of perceptual assessment. SSIM estimates the luminance $\mu_{\mathbf{y}}$ of an image \mathbf{y} as the mean of the intensity, while it estimates the contrast $\sigma_{\mathbf{y}}$ as the standard deviation of the intensity.

$$\mu_{\mathbf{y}} = \frac{1}{N_{\mathbf{y}}} \sum_{p \in \Omega_{\mathbf{y}}} \mathbf{y}_{p},\tag{7}$$

$$\mathbf{y} = \frac{1}{N_{\mathbf{y}} - 1} \sum_{p \in \Omega_{\mathbf{y}}} \left[\mathbf{y}_p - \mu_{\mathbf{y}} \right]^2 \tag{8}$$

To enable the comparison of these entities, a similarity comparison function S is introduced:

 σ

$$S(x, y, c) = \frac{2 \cdot x \cdot y + c}{x^2 + y^2 + c},$$
(9)

The variables x and y are the scalar variables being compared, $c = (k \cdot L)^2$, where $0 < k \ll 1$ is a constant used to avoid instability. Given a real image y and its approximation \hat{y} , the comparison for brightness (C_l) and contrast (C_c) is as follows:

$$C_{l}(\mathbf{y}, \hat{\mathbf{y}}) = S(\mu_{\mathbf{y}}, \mu_{\hat{\mathbf{y}}}, c_{1}) \text{ and } C_{c}(\mathbf{y}, \hat{\mathbf{y}}) = S(\sigma_{\mathbf{y}}, \sigma_{\hat{\mathbf{y}}}, c_{2})$$
(10)

where $c_1, c_2 > 0$. The empirical co-variance

$$\sigma_{\mathbf{y},\hat{\mathbf{y}}} = \frac{1}{N_{\mathbf{y}} - 1} \sum_{p \in \Omega_{\mathbf{y}}} \left(\mathbf{y}_p - \mu_{\mathbf{y}} \right) \cdot \left(\hat{\mathbf{y}}_p - \mu_{\hat{\mathbf{y}}} \right), \tag{11}$$

determines the structure comparison (C_s), expressed as the correlation coefficient between y and \hat{y} :

$$C_s\left(\mathbf{y}, \hat{\mathbf{y}}\right) = \frac{\sigma_{\mathbf{y}, \hat{\mathbf{y}}} + c_3}{\sigma_{\mathbf{y}} \cdot \sigma_{\hat{\mathbf{y}}} + c_3},\tag{12}$$

where $c_3 > 0$. Finally, the SSIM is defined as:

$$SSIM(\mathbf{y}, \hat{\mathbf{y}}) = \left[\mathcal{C}_{l}(\mathbf{y}, \hat{\mathbf{y}})\right]^{\alpha} \cdot \left[\mathcal{C}_{c}(\mathbf{y}, \hat{\mathbf{y}})\right]^{\beta} \cdot \left[\mathcal{C}_{s}(\mathbf{y}, \hat{\mathbf{y}})\right]^{\gamma}$$
(13)

where $\alpha > 0, \beta > 0$ and $\gamma > 0$ are adjustable control parameters for weighting relative importance of all components.



Figure 10: Overall architecture of TimeIR. Initially, shallow feature extraction is performed, followed by deep extraction using multiple TSTB layers, with the Transformer architecture detailed in Figure 11



Figure 11: The overall architecture of the Transformer backbone. The left side of the figure illustrates the basic architecture of the Transformer model, while the right side presents a schematic representation of the self-attention mechanism.



Figure 12: The black line represents the ground truth, the blue line denotes Time Series One, and the red line indicates Time Series Two.

B MODEL DETAILS

866

867

B.1 IMAGE RECONSTRUCTION IN TIME-SERIES FORECASTING

The reason we introduced image reconstruction techniques in time series prediction tasks is that current time series prediction models often lack style continuity. Specifically, the geometric structure of the predicted time series does not align with that of the input time series. This issue arises from the use of the MSE loss function, which only reflects the numerical similarity between the predicted results and the true results, without capturing structural similarity. In other words, two samples may have the same MSE with respect to the ground truth, yet their shapes may differ.

In Figure 12, we present a detailed example: the ground truth is a sine function, while Time Series
One is a sine function that has been shifted both horizontally and vertically, and Time Series Two is a
straight line. By controlling the magnitude of the shift in Time Series One, we can make both time
series have the same MSE as the ground truth. However, it is clear that Time Series Two lacks any
meaningful information, such as periodicity and trends, even though its MSE is identical to that of
Time Series One. In practical applications of time series, we would prefer to use Time Series Two,
which possesses similar periodicity and trends and contains more information.

When ordinary time series prediction models use MSE as the loss function, they focus solely on
numerical similarity while neglecting the geometrical structural similarity in a two-dimensional
space. Therefore, we attempted to introduce image reconstruction models to capture this aspect of
information. Experimental results have shown that our model achieves better reconstruction results,
maintaining a low MSE while providing improved structural integrity.

B.2 Alignment details

To ensure the alignment of \mathbf{P}, \mathbf{T} and \mathbf{X} for concatenation in Part II, we must match their lengths. 889 When the input length is less than the prediction length, we truncate X to the input length. This 890 approach is effective because x serves as a style feature, and altering its length does not significantly 891 impact the stored information. When the input length equals the prediction length, concatenation 892 can be done directly. When the prediction length exceeds the input length, we maintain a sliding 893 window on **P** and **T** with a length equal to the input length. The window moves from the beginning to the end, capturing segments that are then concatenated with X to predict the corresponding output 894 segments. For any excess length after division, an additional window is used to capture the final 895 portion. The reason all windows are concatenated with X is that X acts as a style feature and does 896 not provide any temporal information. Algorithm 1 outlines the complete operation of the model, 897 including a detailed explanation of the input alignment logic. 898

899 900

901

902

886

887

C FULL FORECASTING RESULTS

Table 5 presents a comprehensive comparison of VisiTER with other baselines across seven datasets. It is evident that our model performs exceptionally well on the majority of the datasets.

- 903 904
- 905 906
- 907
- 908
- 909 910
- 911
- 912
- 913
- 914
- 915
- 916
- 917

919 920 921 922 923 924 925 926 Algorithm 1 VisiTER - Overall Architecture. 927 **Require:** Input lookback time series $\mathbf{X} \in \mathbb{R}^{L_1 \times N}$; input Length L_1 ; predicted length L_2 ; variates 928 number N; H is a hyperparameter in T2V; X_i represents the image version of X 929 930 $\triangleright \mathbf{P} \in \mathbb{R}^{L_2 \times N}, \mathbf{T} \in \mathbb{R}^{L_2 \times N}$ 1: $\mathbf{P}, \mathbf{T} = \text{DTFE}(\mathbf{X}).$ 931 $\triangleright \mathbf{P_i} \in \mathbb{R}^{H \times L_2 \times N}$ 2: $\mathbf{P_i} = T2V(\mathbf{P})$. 932 $\triangleright \mathbf{T_i} \in \mathbb{R}^{H \times L_2 \times N}$ 3: $\mathbf{T}_{i} = T2V(\mathbf{P})$. 933 934 $\triangleright \mathbf{X}_{i} \in \mathbb{R}^{H \times L_{1} \times N}$ 4: $\mathbf{X}_{\mathbf{i}} = T2V(\mathbf{X})$. 935 5: if $L_1 \ge L_2$: 936 937 for n in $\{0, ..., N - 1\}$: 6: Processing each variable individually 938 $\triangleright \mathbf{H_n} \in \mathbb{R}^{H \times L_2 \times 3}$ $\mathbf{H_n} = \mathbf{Concat}(X_i[:,:L_2,n], P_i[:,:,n], T_i[:,:,n])$ 7: 939 Image reconstruction using the TimeIR model 8: 940 $\triangleright \, \mathbf{I_n} \in \mathbb{R}^{H \times L_2}$ 941 $\mathbf{I_n} = \texttt{TimeIR}\left(\mathbf{H_n}\right)$ 9: 942 $\triangleright \mathbf{I} \in \mathbb{R}^{H \times L_2 \times N}$ 10: $I = Concat(I_1, \ldots, I_n)$ 943 11: else $L_1 \ge L_2$: 944 945 for l in $\{0, ..., \mathbf{L_2}//\mathbf{L_1} - 2\}$: ▷ Sliding window. 12: 946 for n in $\{0, ..., N - 1\}$: ▷ Processing each variable individually 13: 947 $\triangleright \mathbf{H}_{\mathbf{n}}^{\mathbf{l}} \in \mathbb{R}^{H \times L_1 \times 3}$ 948 14: 949 $\mathbf{H}_{n}^{l} = Concat(\mathbf{X_{i}}, \mathbf{P_{i}}[:, l \times \mathbf{L_{1}}: (l+1) \times \mathbf{L_{1}}, n], \mathbf{T_{i}}[:, l \times \mathbf{L_{1}}: (l+1) \times \mathbf{L_{1}}, n])$ 15: 950 $I_n = \text{TimeIR}(H_n^l)$ $\triangleright \mathbf{I_n} \in \mathbb{R}^{H \times L_1}$ 16: 951 952 $\triangleright \mathbf{I}^{\mathbf{l}} \in \mathbb{R}^{H \times L_1 \times N}$ $I^{l} = Concat(I_{1}, \ldots, I_{n})$ 17: 953 18: for n in $\{0, ..., N - 1\}$: > Processing each variable individually 954 $\triangleright \mathbf{H}_{n}^{l+1} \in \mathbb{R}^{H \times L_{1} \times 3}$ 955 19: 956 $\mathbf{H_n^{l+1}} = Concat(\mathbf{X_{image}}, \mathbf{P_{image}}[:, -\mathbf{L_1}:, n], \mathbf{T_{image}}[:, -\mathbf{L_1}:, n])$ 20: 957 $I_n = \text{TimeIR}(H_n^{l+1})$ $\triangleright \mathbf{I_n} \in \mathbb{R}^{H \times L_1}$ 21: 958 $\triangleright \mathbf{I^{l+1}} \in \mathbb{R}^{H \times L_1 \times N}$ $I^{l+1} = Concat(I_1, \ldots, I_n)$ 959 22: 960 $\triangleright \mathbf{I} \in \mathbb{R}^{H \times L_2 \times N}$ $\mathbf{I} = Concat(\mathbf{I}^0, \dots, \mathbf{I}^l, \mathbf{I}^{l+1}[:, -(\mathbf{L_2} \bmod \mathbf{L_1}):, :])$ 23: 961 $\triangleright \hat{\mathbf{Y}} \in \mathbb{R}^{L_2 \times N}$ 24: $\hat{\mathbf{Y}} = \mathbf{V2T}(\mathbf{I})$ 962 963 25: Return Ŷ \triangleright Return the prediction result $\hat{\mathbf{Y}}$ 964 965 966 967

968

918

- 969
- 970

Table 5: Full results of the long-term forecasting task. We compare extensive competitive models under different prediction lengths following the setting of TimesNet (2023a). The input sequence length is set to 96 for all baselines. *Avg* means the average results from all four prediction lengths.

	-											
Mc	odels	VisiTER (Ours)	iTransformer (2023)	SparseTSF (2024)	PatchTST (2023)	Crossformer (2023)	TiDE (2023)	TimesNet (2023a)	DLinear (2023)	SCINet (2022b)	FEDformer (2022)	Stationary (2022c)
M	etric	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 Avg	0.303 0.326 0.369 0.354 0.400 0.380 0.473 0.423 0.386 0.371	0.3340.3680.3770.3910.4260.4200.4910.4590.4070.410	0.357 0.375 0.394 0.393 0.426 0.414 0.488 0.449	0.329 0.367 0.367 0.385 0.399 0.410 0.454 0.439 0.387 0.400	0.404 0.426 0.450 0.451 0.532 0.515 0.666 0.589 0.513 0.496	0.364 0.387 0.398 0.404 0.428 0.425 0.487 0.461	0.338 0.375 0.374 0.387 0.410 0.411 0.478 0.450	0.345 0.372 0.380 0.389 0.413 0.413 0.474 0.453 0.403 0.407	0.418 0.438 0.439 0.450 0.490 0.485 0.595 0.550 0.485 0.481	0.379 0.419 0.426 0.441 0.445 0.459 0.543 0.490	0.386 0.398 0.459 0.444 0.495 0.464 0.585 0.516 0.481 0.456
ETTm2	96 192 336 720 Avg	0.174 0.255 0.240 0.299 0.302 0.340 0.400 0.398 0.279 0.323	0.180 0.264 0.250 0.309 0.311 0.348 0.412 0.407 0.288 0.332	0.186 0.268 0.248 0.306 0.308 <u>0.343</u> 0.408 0.398 0.288 0.329	$\begin{array}{c} 0.175 & 0.259 \\ 0.241 & 0.302 \\ 0.305 & 0.343 \\ 0.402 & 0.400 \\ 0.281 & 0.326 \end{array}$	0.287 0.366 0.414 0.492 0.597 0.542 1.730 1.042	0.207 0.305 0.290 0.364 0.377 0.422 0.558 0.524 0.358 0.404	0.187 0.267 0.249 0.309 0.321 0.351 0.408 0.403	0.193 0.292 0.284 0.362 0.369 0.427 0.554 0.522 0.350 0.401	0.286 0.377 0.399 0.445 0.637 0.591 0.960 0.735 0.571 0.537	0.203 0.287 0.269 0.328 0.325 0.366 0.421 0.415	0.192 0.274 0.280 0.339 0.334 0.361 0.417 0.413 0.306 0.347
ETTh1	96 192 336 720	0.374 0.383 0.416 0.422 0.459 0.444 0.475 0.461	0.386 0.405 0.441 0.436 0.487 0.458 0.503 0.491	0.386 0.393 0.435 0.422 0.476 0.440 0.460 0.455	0.414 0.419 0.460 0.445 0.501 0.466 0.500 0.488	0.423 0.448 0.471 0.474 0.570 0.546 0.653 0.621	0.479 0.464 0.525 0.492 0.565 0.515 0.594 0.558	0.384 0.402 0.436 0.429 0.491 0.469 0.521 0.500	0.386 0.400 0.437 0.432 0.481 0.459 0.519 0.516	0.654 0.599 0.719 0.631 0.778 0.659 0.836 0.699	$\begin{array}{c} 0.376 \\ 0.420 \\ 0.420 \\ 0.459 \\ 0.506 \\ 0.507 \\$	0.513 0.491 0.534 0.504 0.588 0.535 0.643 0.616
	Avg	0.431 0.428	0.454 0.447	0.440 0.429	0.469 0.454	0.529 0.522	0.541 0.507	0.458 0.450	0.456 0.452	0.747 0.647	7 <u>0.440</u> 0.460	0.570 0.537
ETTh2	96 192 336 720	0.284 0.337 0.364 0.393 0.406 0.423 0.417 <u>0.440</u>	0.297 0.349 0.380 0.400 0.428 0.432 0.427 0.445	$\begin{array}{c} 0.304 & 0.347 \\ 0.385 & 0.396 \\ 0.421 & 0.428 \\ 0.420 & 0.437 \end{array}$	0.302 0.348 0.388 0.400 0.426 0.433 0.431 0.446	0.745 0.584 0.877 0.656 1.043 0.731 1.104 0.763	0.400 0.440 0.528 0.509 0.643 0.571 0.874 0.679	0.340 0.374 0.402 0.414 0.452 0.452 0.462 0.468	0.333 0.387 0.477 0.476 0.594 0.541 0.831 0.657	0.707 0.621 0.860 0.689 1.000 0.744 1.249 0.838	0.358 0.397 0.429 0.439 0.496 0.487 0.463 0.474	0.476 0.458 0.512 0.493 0.552 0.551 0.562 0.560
	Avg	0.368 0.398	0.383 0.407	0.383 0.402	0.387 0.407	0.942 0.684	0.611 0.550	0.414 0.427	0.559 0.515	0.954 0.723	8 0.437 0.449	0.526 0.516
ECL	96 192 336 720	0.170 <u>0.263</u> 0.178 <u>0.271</u> 0.194 <u>0.287</u> 0.233 <u>0.319</u>	0.148 0.240 0.162 0.253 0.178 0.269 0.225 0.317	0.210 0.280 0.206 0.282 0.219 0.296 0.260 0.328	0.181 0.270 0.188 0.274 0.204 0.293 0.246 0.324	0.219 0.314 0.231 0.322 0.246 0.337 0.280 0.363	0.237 0.329 0.236 0.330 0.249 0.344 0.284 0.373	0.168 0.272 0.184 0.289 0.198 0.300 0.220 0.320	0.197 0.282 0.196 0.285 0.209 0.301 0.245 0.333	0.247 0.345 0.257 0.355 0.269 0.369 0.299 0.390	0.193 0.308 0.201 0.315 0.214 0.329 0.246 0.355	0.169 0.273 0.182 0.286 0.200 0.304 0.222 0.321
	Avg	0.194 <u>0.285</u>	0.178 0.270	0.224 0.297	0.205 0.290	0.244 0.334	0.251 0.344	<u>0.192</u> 0.295	0.212 0.300	0.268 0.365	6 0.214 0.327	0.193 0.296
Exchange	96 192 336 720	0.083 0.200 0.167 0.293 0.312 0.404 0.772 0.660	0.086 0.206 0.177 0.299 0.331 0.417 0.847 0.691	0.095 0.218 0.184 0.307 0.324 0.414 0.839 0.691	0.088 <u>0.205</u> 0.176 <u>0.299</u> 0.301 0.397 0.901 0.714	0.256 0.367 0.470 0.509 1.268 0.883 1.767 1.068	0.094 0.218 0.184 0.307 0.349 0.431 0.852 0.698	0.107 0.234 0.226 0.344 0.367 0.448 0.964 0.746	0.088 0.218 0.176 0.315 0.313 0.427 <u>0.839</u> 0.695	0.267 0.396 0.351 0.459 1.324 0.853 1.058 0.797	0.148 0.278 0.271 0.315 0.460 0.427 1.195 0.695	0.111 0.237 0.219 0.335 0.421 0.476 1.092 0.769
	Avg	0.334 0.389	0.360 0.403	0.361 0.408	0.367 <u>0.404</u>	0.940 0.707	0.370 0.413	0.416 0.443	0.354 0.414	0.750 0.626	60.519 0.429	0.461 0.454
Weather	96 192 336 720	0.172 0.214 0.218 0.255 0.274 0.296 0.351 0.345	0.174 0.214 0.221 0.254 0.278 0.296 0.358 0.347	0.197 0.237 0.244 0.273 0.293 0.308 0.368 0.357	0.177 <u>0.218</u> 0.225 0.259 0.278 <u>0.297</u> 0.354 <u>0.348</u>	0.158 0.230 0.206 0.277 0.272 0.335 0.398 0.418	0.202 0.261 0.242 0.298 0.287 0.335 0.351 0.386	0.172 0.220 0.219 0.261 0.280 0.306 0.365 0.359	0.196 0.255 0.237 0.296 0.283 0.335 0.345 0.381	0.221 0.306 0.261 0.340 0.309 0.378 0.377 0.427	0.217 0.296 0.276 0.336 0.339 0.380 0.403 0.428	$\begin{array}{c} 0.173 \ 0.223 \\ 0.245 \ 0.285 \\ 0.321 \ 0.338 \\ 0.414 \ 0.410 \end{array}$
Ì	Avg	0.254 0.277	0.258 0.278	0.276 0.294	0.259 0.281	0.259 0.315	0.271 0.320	0.259 0.287	0.265 0.317	0.292 0.363	80.309 0.360	0.288 0.314

Table 6: The complete SSIM results of our model compared to other baselines with prediction lengths $S \in \{96, 192, 336, 720\}$ for others and fixed lookback length T = 96. TimeIR performs **ZERO-SHOT** inference when predicting lengths of $\{192, 336, 720\}$. Avg means further averaged by subsets. Higher SSIM values indicate better performance.

Models		ETTm1	ETTm2	ETTh1	ETTh2	ECL	Exchange	Weather
	96	0.3829	0.4924	0.4070	0.4186	0.6110	0.6379	0.4381
Autoformer (Wu et al., 2021)	192	0.3498	0.4857	0.3518	0.4091	0.6112	0.6308	0.4191
	336	0.3729	0.4655	0.3936	0.4106	0.6061	0.6141	0.4379
	720	0.3990	0.4909	0.3945	0.4021	0.5655	0.5843	0.4557
	Avg	0.3762	0.4836	0.3867	0.4101	0.6051	0.6168	0.4377
	96	0.3979	0.4594	0.4529	0.3826	0.6988	0.5692	0.4714
	192	0.3907	0.4247	0.4160	0.2423	0.6931	0.5137	0.3913
Crossformer (Zhang & Yan, 2023)	336	0.3718	0.3610	0.3863	0.2163	0.6718	0.1772	0.4409
	720	0.3596	0.2837	0.3639	0.2070	0.6517	0.1067	0.3759
	Avg	0.3800	0.3822	0.4047	0.2620	0.6788	0.3417	0.4199
	96	0.4294	0.5232	0.4485	0.4409	0.7167	0.6971	0.6079
	192	0.4302	0.4991	0.4186	0.4167	0.7079	0.6561	0.5862
iTransformer (Liu et al., 2023)	336	0.4333	0.5029	0.4058	0.4077	0.6984	0.6328	0.5745
	720	0.4386	0.4848	0.3954	0.3998	0.6828	0.5911	0.5661
	Avg	0.4329	0.5025	0.4171	0.4163	0.7015	0.6443	0.5837
	96	0.4267	0.5361	0.4394	0.4352	0.6626	0.6640	0.5884
	192	0.4237	0.5115	0.4206	0.3936	0.6520	0.6329	0.5725
TimesNet (Wu et al., 2023a)	336	0.4259	0.5065	0.3935	0.3953	0.6453	0.6158	0.5637
	720	0.4405	0.4935	0.3754	0.3811	0.6345	0.5686	0.5633
	Avg	0.4292	0.5119	0.4072	0.4013	0.6486	0.6203	0.5720
	96	0.4569	0.5624	0.4796	0.4772	0.6976	0.7156	0.6125
VisiTER (ours)	192	0.4561	0.5410	0.4605	0.4178	0.6952	0.6723	0.5942
	336	0.4513	0.5286	0.4538	0.4310	0.6861	0.6456	0.5867
	720	0.4568	0.5216	0.4485	0.4319	0.6682	0.6053	0.5829
	Avg	0.4553	0.5384	0.4606	0.4395	0.6868	0.6597	0.5941

D MORE VISUALIZATION RESULTS

1080

1081 1082

1083

1084

1085

1086

Figure 13 presents additional visual results. It is clear that our model provides more accurate predictions, particularly when the time series approaches a straight line. In such cases, our model is able to reconstruct it as a straight line, whereas traditional time series prediction models struggle to capture the fluctuations.



Figure 13: More comparison of visual results for time series reconstruction. The experiment focuses on predicting 96 steps with a lookback length of 96, using the ETTh2 and Weather dataset, where results from different variables are sampled. The first half of each image displays the provided time series, while the second half shows the predictions from various models. The last column represents the ground truth. Please zoom in for a closer view.

E EXPECTATIONS FOR FUTURE RESEARCH

This paper introduces the first method that utilizes image reconstruction for time series prediction, leveraging periodic and trend information to reconstruct time series. The core approach involves incorporating the style of the time series through image reconstruction. However, when predicting long time series, the style can change over time as the series lengthens, potentially introducing noise into the supplementary information. Future research could focus on better utilizing the style feature to address this issue.

Additionally, we directly employed traditional Transformer models for both feature prediction and time series image reconstruction. Subsequent work could modify the Transformer architecture to make it more suitable for reconstructing time series images. While our experiments concentrated on long time series prediction, this method is also applicable to other time series tasks, such as short time series prediction and imputation.