# ROSE: REGISTER-ASSISTED GENERAL TIME SE RIES FORECASTING WITH DECOMPOSED FREQUENCY LEARNING

Anonymous authors

Paper under double-blind review

#### ABSTRACT

With the increasing collection of time series data from various domains, there arises a strong demand for general time series forecasting models pre-trained on a large number of time-series datasets to support a variety of downstream prediction tasks. Enabling general time series forecasting faces two challenges: how to obtain unified representations from multi-domian time series data, and how to capture domain-specific features from time series data across various domains for adaptive transfer in downstream tasks. To address these challenges, we propose a Register Assisted General Time Series Forecasting Model with Decomposed Frequency Learning (**ROSE**), a novel pre-trained model for time series forecasting. ROSE employs Decomposed Frequency Learning for the pre-training task, which decomposes coupled semantic information in time series with frequency-based masking and reconstruction to obtain unified representations across domains. We also equip ROSE with a Time Series Register, which learns to generate a register to capture domain-specific representations during pre-training and enhances domainadaptive transfer by selecting related register tokens on downstream tasks. After pretraining on large-scale time series data, ROSE achieves state-of-the-art forecasting performance on 7 real-world benchmarks. Remarkably, it demonstrates competitive or superior few-shot and zero-shot abilities.

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

## 1 INTRODUCTION

033 Time Series Forecasting plays a pivotal role across numerous domains, such as energy, smart 034 transportation, weather, and economics (Qiu et al., 2024). However, training specific models in deep learning for each dataset is costly and requires tailored parameter tuning, whose prediction accuracy 035 may be limited due to data scarcity (Liu et al., 2024). One solution is to pre-train a general model on 036 diverse time series datasets and fine-tune it with a few data for different downstream scenarios or 037 direct predict without fine-tuning. Following this idea, foundation models for time series forecasting have recently raised growing attention, continually scaling up the pre-training datasets and model sizes to improve the generalization performance (Woo et al., 2024; Goswami et al., 2024; Ansari et al., 040 2024). However, an excessively large scale also causes increasing costs in training and inference, 041 which may go against the original intention of a general model, especially in resource-constrained 042 situations. In addition to scaling up, general time series forecasting models can also be designed from 043 the perspectives of pre-training tasks and downstream transfer adaptation. From these two angles, we 044 identify the following challenges.

Obtaining a unified representation from time series data across various domains is challenging. Time series from each domain involve complex temporal patterns, composed of multiple frequency components combined with each other (Zhou et al., 2022), which is frequency superposition. As shown in Figure 1(a), different frequency components contain distinct semantic information. For example, low and high-frequency components represent long-term trends and rapid variations, respectively (Zhang et al., 2022). Furthermore, different datasets exhibit diverse frequency distributions, and the significance of low-frequency and high-frequency components for time series modeling varies across domains(Zhang et al., 2024). As a result, large-scale time series data from different domains introduce even more complex temporal patterns and frequency diversity. Existing pre-training frameworks (Dong et al., 2024; Nie et al., 2022; Lee et al., 2023), such as masked modeling and contrastive learning, were proposed to learn a unified representation from time domain. However, these methods overlook the frequency diversity and complexity exhibited in heterogeneous time series that come from various domains, making it difficult to capture intricate patterns, thus limiting their generalization capabilities.

Adaptive transferring information from multi-domain time series to specific downstream scenarios presents a challenge. Multi-source time series data originate from various domains (Woo 060 et al., 2024), whose data exhibit domain-specific information (Liu et al., 2024). Information from the 061 same or similar domain as the target domain is useful for improving the model's effectiveness in the 062 target task (Chen et al., 2023a). However, as shown in Figure 1(a), existing time series pre-training 063 frameworks (Woo et al., 2024; Liu et al., 2024; Zhou et al., 2024) focus mainly on learning generalized 064 representation during pre-training and overlook domain-specific representation. Thus, they only transfer the same generalized representation to different target domains, called *direct transfer*, which 065 limits the model's effectiveness in specific downstream tasks. Therefore, it is necessary to learn 066 domain-specific information during pre-training and adaptively transfer the specific representations 067 to each target domain, called *adaptive transfer*. Realizing adaptive transfer poses two difficulties: 1) 068 capturing domain-specific information in pre-training. 2) adaptive use of domain-specific information 069 in various downstream tasks.



Figure 1: (a) Pre-training on multi-domain datasets that exhibit combined frequency. Existing general time series forecasting models only extract generalized representations for direct transfer to various downstream target domains. We propose to learn generalized and specific representations during pre-training, and adaptively transfer them to each target domain. (b) The t-SNE visualization of the hidden representations after direct transfer and adaptive transfer: In direct transfer, representations of different domains are mixed, but in adaptive transfer, they show a clear clustering pattern. The detailed experiment setting is in the Appendix A.10.4

088 To address these challenges, we propose a register assisted general time series forecasting model 089 with decomposed frequency learning (ROSE). First, we propose Decomposed Frequency Learning 090 that learns generalized representations to solve the issue with coupled semantic information. We 091 decompose individual time series using the Fourier transform with a novel frequency-based masking 092 method, and then convert it back to the time domain to obtain decoupled time series for reconstruction. It makes complex temporal patterns disentangled, thus benefiting the model to learn generalized 094 representations. **Second**, we introduce Time Series Register (TS-Register) to learn domain-specific information in multi-domain data. By setting up a register, we generate register tokens to learn each domain-specific information during pre-training. In a downstream scenario, the model adaptively 096 selects Top-K vectors from the register that are close to the target domain of interest. During fine-tuning, we adjust the selected register tokens with a novel learnable low-rank matrix, which 098 complements target-specific information to perform more flexible adaptive transfer. As shown in Figure 1(b), adaptive transfer successfully utilizes domain-specific information in multi-domain 100 time series, which contributes to the model's performance in target tasks. The contributions are 101 summarized as follows: 102

102 103 104

- We propose ROSE, a novel light weight general time series forecasting model using multi-domain datasets for pre-training and improving downstream fine-tuning performance and efficiency.
- We propose a novel Decomposed Frequency Learning that employs multi-frequency masking to learn complex general temporal patterns from multi-domain data, empowering the model's generalization capability.

- We propose a novel TS-Register to capture domain-specific information in pre-training and enable adaptive transfer of target-oriented specific information for downstream tasks.
- Our experiments with 7 real-world benchmarks demonstrate that ROSE achieves state-of-the-art performance in full-shot setting and achieves competitive or superior results in few-shot setting, along with impressive transferability in zero-shot setting.
- 113 114 115

109

110

111

112

## 2 RELATED WORK

- 116
- 117 118

## 2.1 TRADITIONAL TIME SERIES FORECASTING

The statistical time series forecasting models like ARIMA (Box and Jenkins, 1968), despite their 119 theoretical support, are limited in modeling nonlinearity. With the rise of deep learning, many 120 RNN-based models (Cirstea et al., 2019; Wen et al., 2017; Salinas et al., 2020) have been proposed, 121 modeling the sequential data with an autoregressive process. CNN-based models (Luo and Wang, 122 2024; Liu et al., 2022a) have also received widespread attention due to their ability to capture local 123 features. MICN (Wang et al., 2022) utilizes TCN to capture both local and global features, while 124 TimesNet (Wu et al., 2022) focuses on modeling 2D temporal variations. However, both RNNs and 125 CNNs struggle to capture long-term dependencies. Transformer-based models (Zhou et al., 2022; 126 Nie et al., 2022; Wu et al., 2021; Liu et al., 2023; Chen et al., 2024), with their attention mechanism, 127 can capture long dependencies and extract global information, leading to widespread applications in 128 long-time series prediction. However, this case-by-case paradigm requires meticulous hyperparameter design for different datasets, and its predictive performance can also be affected by data scarcity. 129

130 131

132

2.2 PRE-TRAINED MODELS FOR TIMES SEIRES FORECASTING

LLMs for time series forecasting: Recent studies have shown that Large Language Models (LLMs) can enhance time series forecasting with limited fine-tuning (Zhou et al., 2024), prompting (Jin et al., 2024; Cao et al., 2024), and modality alignment (Jin et al., 2024). GPT4TS (Zhou et al., 2024) fine-tunes a subset of LLM parameters, and achieves competitive performance by leveraging pre-trained text knowledge. TimeLLM (Jin et al., 2024) transforms time series into text to align with LLM representations. Despite their positive impact, LLMs' high inference costs limit their practicality.

Time series foundation model: Pre-training with multiple sources time series has recently received 140 widespread attention (Rasul et al., 2023; Dooley et al., 2024; Garza and Mergenthaler-Canseco, 141 2023; Kamarthi and Prakash, 2023). MOMENT (Goswami et al., 2024) and MOIRAI (Woo et al., 142 2024) adopt a BERT-style pre-training approach, while Timer (Liu et al., 2024), Chronos (Ansari 143 et al., 2024) and TimsFM (Das et al., 2023a) use a GPT-style pre-training approach, giving rise to 144 improved performance in time series prediction. However, the above methods overlook domain-145 specific information from multi-source data, thus limiting the performance of the models. Different 146 from previous approaches, ROSE pre-trains on large-scale data from various domains and it considers 147 both generalized representations and domain-specific information, which facilitates flexible adaptive 148 transfer in downstream tasks.

149 150 151

## 3 Methodology

152 153 **Problem Definition.** Given a multivariate time series  $\mathbf{X}_t = {\mathbf{x}_{t-L:t}^i}_{i=1}^C$ , where each  $\mathbf{x}_{t-L:t}^i \in \mathbb{R}^L$  is 154 a sequence of observations. *L* denotes the look-back window and *C* denotes the number of channels. 155 The forecasting task is to predict the future values  $\hat{\mathbf{Y}}_t = {\{\hat{\mathbf{x}}_{t:t+F}^i\}_{i=1}^C}$ , where *F* denotes the forecast 156 horizon.  $\mathbf{Y}_t = {\{\mathbf{x}_{t:t+F}^i\}_{i=1}^C}$  is the ground truth of future.

The general time series forecasting model is pre-trained with multi-source datasets  $\mathbf{D}_{\text{pre-train}} = \{(\mathbf{X}_t^j, \mathbf{Y}_t^j)\}_{j=1}^N$ , where N is the number of datasets. For the downstream task, the model is fine-tuned with a training dataset  $\mathbf{D}_{\text{train}} = \{(\mathbf{X}_t^{\text{train}}, \mathbf{Y}_t^{\text{train}})\}$ , and is tested with  $\mathbf{D}_{\text{test}} = \{(\mathbf{X}_t^{\text{test}}, \mathbf{Y}_t^{\text{test}})\}$  to predict  $\hat{\mathbf{Y}}_t^{\text{test}}$ , where  $\mathbf{D}_{\text{pre-train}}, \mathbf{D}_{\text{train}}$  and  $\mathbf{D}_{\text{test}}$  are pairwise disjoint. Alternatively, the model could be directly tested using  $\mathbf{D}_{\text{test}}$  without fine-tuning with  $\mathbf{D}_{\text{train}}$  to predict  $\hat{\mathbf{Y}}_t^{\text{test}}$ .

#### 162 3.1 Architecture 163





179 The time series forecasting paradigm of ROSE contains two steps: *pre-training* and *fine-tuning*. 180 First, ROSE is pre-trained on large-scale datasets from various domains, with two pre-training tasks: 181 reconstruction and prediction. We set up the reconstruction task to help the model understand time 182 series comprehensively and set the prediction task to enhance the model's few-shot and zero-shot 183 abilities. Then, ROSE is fine-tuned with a target dataset in the downstream scenario.

As shown in Figure 2, ROSE employs the encoder-decoder architecture to model time series. The 185 backbone consists of multiple Transformer layers to process sequential information and effectively capture temporal dependencies (Vaswani et al., 2017). The reconstruction decoder and prediction 187 decoder use the same structure as the Transformer encoder. They are used for reconstruction and 188 prediction tasks respectively. ROSE is pre-trained in a channel-independent way, which is widely 189 used in time series forecasting (Nie et al., 2022).

190 Input Representations. To enhance the generalization of ROSE for adaptive transferring from multi-191 domains to different target domains, we model the inputs x with **patch tokens** and **register tokens**. 192 Patch tokens are obtained by partitioning the time series using patching layers (Nie et al., 2022), to 193 preserve local temporal information. Register tokens that capture domain-specific information will 194 be introduced in Section 3.3.

195 196 197

#### 3.2 DECOMPOSED FREQUENCY LEARNING

As shown in Figure 1, time series data are composed of multiple superimposed frequency components, resulting in the overlap of different temporal changes. Furthermore, low-frequency components 199 typically contain information about overall trends and longer-scale variations, and high-frequency 200 components usually contain information about short-term fluctuations and shorter-scale variations, 201 therefore, understanding time series from low and high frequencies separately benefits general time 202 series representation learning. Based on the observations above, we propose a novel frequency-based 203 masked modeling that randomly mask either high-frequency or low-frequency components of a 204 time series multiple times as the key to enable learning of common time series patterns, such as 205 trends and various long and short term fluctuations. Finally, reconstruction task assists the model 206 in comprehending the data from various frequency perspectives, enabling it to learn generalized 207 representations. In contrast, existing frequency masking methods (Zhang et al., 2022), which 208 randomly mask frequencies of a single time series once, show limited forecasting effectiveness due to 209 the lack of common pattern learning from heterogeneous time series that come from various domains.

210 **Multi-frequency masking.** As shown in the green part of Figure 3, given a time series  $\mathbf{x} \in \mathbb{R}^L$ , we 211 utilize the Real Fast Fourier Transform (rFFT) (Brigham and Morrow, 1967) to transform it into the 212 frequency domain, giving rise to  $\mathbf{x}_{\text{freq}} \in \mathbb{C}^{L/2+1}$ .

- 213  $\mathbf{x}_{\mathrm{freq}} = \mathrm{rFFT}(\mathbf{x}).$ (1)214
- To separately model high-frequency and low-frequency information in time series, we sample  $K_{\rm f}$ 215 thresholds  $\tau_1, \tau_2, \tau_3, ..., \tau_{K_f}$  and  $K_f$  random numbers  $\mu_1, \mu_2, \mu_3, ..., \mu_{K_f}$  for multi-frequency masks,

220 221 222

229

230

231

232

233

234

235 236

237



Figure 3: An illustration of decomposed frequency learning. Based on the sampled thresholds, we randomly apply low/high-frequency masking to the time series in the frequency domain and then transform it back to the time domain for reconstruction.

where  $\tau \in \text{Uniform}(0, a)$ , a < L/2 + 1, and  $\mu \in \text{Bernoulli}(p)$ . Each pair of  $\tau_i$  and  $\mu_i$  corresponds to the  $i_{th}$  frequency mask. This generates a mask matrix  $\mathbf{M} \in \{0,1\}^{K_f \times (L/2+1)}$ , where each row corresponds to the  $i_{th}$  frequency mask, each column corresponds to the  $j_{th}$  frequency, and each element  $m_{ij}$  is 0 or 1, meaning that the  $j_{th}$  frequency is masked with the  $i_{th}$  frequency mask or not.

$$m_{ij} = \begin{cases} \mu_i & , if \ j < \tau_i \\ (1 - \mu_i) & , if \ j > \tau_i \end{cases},$$
(2)

238 where  $\tau_i$  and  $\mu_i$  denote the threshold and random number for the  $i_{th}$  frequency domain mask. If 239  $\mu_i = 1$ , it means that frequency components above  $\tau_i$  will be masked, indicating to mask high frequency, as shown by the threshold  $\tau_1$  in Figure 3. Conversely, if  $\mu_i = 0$ , it signifies that frequency 240 components below  $\tau_i$  will be masked, indicating to mask low frequency, exemplified by threshold  $\tau_2$ 241 in Figure 3. 242

243 After obtaining the mask matrix  $\mathbf{M}$ , we replicate  $\mathbf{x}_{\text{freq}} K_{\text{f}}$  times to get the  $\mathbf{X}_{\text{freq}} \in \mathbb{C}^{K_{\text{f}} \times L/2+1}$  and 244 perform element-wise Hadamard product with the mask matrix M to get masked frequency of time 245 series. Then, we use the inverse Real Fast Fourier Transform (irFFT) to convert the results from the 246 frequency domain back to the time domain and get  $K_{\rm f}$  masked sequences  $\mathbf{X}_{\rm mask} = {\{\mathbf{x}_{\rm mask}^i\}_{i=1}^{K_{\rm f}},$ 247 where each  $\mathbf{x}_{\text{mask}}^i \in \mathbb{R}^L$  corresponding to masking with a different threshold  $\tau_i$ . 248

$$\mathbf{X}_{\text{mask}} = \text{irFFT}(\mathbf{X}_{\text{freq}} \odot \mathbf{M}). \tag{3}$$

249 **Representation learning.** As shown in the yellow part of Figure 3, after obtaining the  $K_{\rm f}$  masked 250 sequences  $\mathbf{X}_{\text{mask}}$ , we divide each sequence  $\mathbf{x}_{\text{mask}}^i$  into P non-overlapping patches, and use a linear 251 layer to transforming them into P patch tokens, and thus we get  $\mathcal{X}_{mp} = \{\mathbf{X}_{mp}^i\}_{i=1}^{K_f}$  to capture general 252 information, where each  $\mathbf{X}_{mp}^{i} \in \mathbb{R}^{P \times D}$ , and D is the dimension for each patch token. We replicate 253 the register tokens  $\mathbf{X}_{u}$   $K_{f}$  times to get  $\mathcal{X}_{u} \in \mathbb{R}^{K_{f} \times N_{r} \times D}$ , where  $\mathbf{X}_{u} \in \mathbb{R}^{N_{r} \times D}$  is obtained by 254 inputting the original sequence into the TS-Register, as detailed in Section 3.3. Then, we concatenate 255 the patch tokens  $\mathcal{X}_{mp}$  with the register tokens  $\mathcal{X}_{u}$ , and feed them into the Transformer encoder to 256 obtain the representation of each masked series. These representations are then aggregated to yield a 257 unified representation  $\mathbf{S}_{m} \in \mathbb{R}^{(N_{r}+P) \times D}$ . The aggregator is the averaging operation. 258

$$\mathbf{S}_{m} = \operatorname{Aggregator}(\operatorname{Encoder}(\operatorname{Concatenate}(\mathcal{X}_{mp}, \mathcal{X}_{u}))). \tag{4}$$

260 **Reconstruction task.** After obtaining the representation  $S_m$ , we feed it into the reconstruction 261 decoder, which shares same stucture as the Tranformer encoder, and ultimately reconstruct the original 262 sequence  $\hat{\mathbf{x}} \in \mathbb{R}^L$  through the reconstruction head, which is a linear layer. As frequency domain 263 masking affects the overall time series, we compute the Mean Squared Error (MSE) reconstruction 264 loss for the entire time series.

$$\mathcal{L}_{\text{reconstruction}} = ||\mathbf{x} - \hat{\mathbf{x}}||_2^2.$$
(5)

265 266 267

268

259

By decomposed frequency learning, we can obtain the general representations, and additionally, 269 we propose the TS-Register that learns register tokens as the domain-specific information from the

270 multi-domain datasets for adaptive transfer. It clusters domain-specific information from the multi-271 domain datasets into register tokens and stores such domain-specific information in the register during 272 pre-training. Then, it adaptively selects domain-specific information from the register via a Top-K 273 selection strategy to enhance the performance in the target domain. A novel learnable low-rank matrix 274 is proposed to set to complement the downstream dataset-specific information through fine-tuning.

275 We set up a randomly initialized register  $\mathbf{E} \in \mathbb{R}^{H \times D_r}$  with H cluster center vectors  $\mathbf{e}_i \in \mathbb{R}^{D_r}, i \in \mathbb{R}^{D_r}$ 276  $\{1, 2, \dots, H\}$ . Each of input time series  $\mathbf{x} \in \mathbb{R}^L$  is projected into a data-dependent embedding 277  $\mathbf{x}_{\mathrm{e}} \in \mathbb{R}^{D_{\mathrm{r}}}$  through a linear layer. 278

**Pre-training stage.** As shown in Figure 2(b), we use the register to cluster these data-dependent 279 embeddings, which generate domain-specific information, and store them in pre-training. Specifically, 280 We find a cluster center vector  $\mathbf{e}_{\delta}$  from the register  $\mathbf{E}$  where we use  $\delta$  to denote the cluster that the 281 data-dependent embedding  $\mathbf{x}_{e}$  belongs to. 282

$$\mathcal{L}_{\text{register}} = \|\mathbf{x}_{e} - \mathbf{e}_{\delta}\|_{2}^{2}, \quad \delta = \operatorname*{arg\,min}_{j=1:H} \|\mathbf{x}_{e} - \mathbf{e}_{j}\|_{2}. \tag{6}$$

To update the cluster center vectors in the register E that represents the domain information of the 285 pre-trained datasets, we set the loss function shown in Equation 6 that minimizes the distance between 286 the data-dependent embedding  $\mathbf{x}_{e}$  and the cluster center  $\mathbf{e}_{\delta}$ . To solve the problem that the gradient 287 of the arg min function cannot be backpropagated, we use the stop gradient operation to pass the 288 gradient of  $\mathbf{e}_{\delta}$  directly to  $\mathbf{x}_{e}$ . 289

290 In this way, the vectors in the register E cluster the embeddings of different data and learn the domain-291 specific centers for pre-trained datasets, which can represent domain-specific information. As a vector in the register  $\mathbf{E}, \mathbf{e}_{\delta}$  represents the domain-specific information for input  $\mathbf{x}, \mathbf{e}_{\delta}$  is invariant under 292 small perturbations in  $\mathbf{x}_{e}$  that represents x, which promotes better representation of domain-specific 293 information and robustness of the vectors in the register. This also avoids their over-reliance on 294 detailed information about specific datasets. 295

The cluster center vector  $\mathbf{e}_{\delta}$  is then patched into  $\mathbf{X}_{u} \in \mathbb{R}^{N_{r} \times D}$ , where  $N_{r}$  is the number of the register 296 tokens and D is the dimensionality of Transformer latent space.  $\mathbf{X}_{u}$  is called register tokens, which are used as the prefix of the patch tokens  $\mathbf{X}_{p} \in \mathbb{R}^{P \times D}$  and input for the Transformer encoder to 297 298 299 provide domain-specific information.

300 Fine-tuning stage. As shown in Figure 2(c), after obtaining a register E that contains domain-301 specific information through pre-training, we freeze the register parameters to adaptively use this 302 domain-specific information in the downstream tasks.

303 Since the target domain may not strictly fit one of the upstream domains, we propose a novel 304 embedding learning of the downstream data by employing a Top-K strategy that selects k similar 305 vectors in the register. As shown in Equation 7, the embedding of input time series  $\mathbf{x}_e$  picks the 306 k nearest vectors in the register **E**, and uses their average as  $\bar{\mathbf{e}}_k$  to represent the domain-specific information from the pre-train stage.  $\bar{\mathbf{e}}_k$  is also patched into  $\mathbf{X}_d \in \mathbb{R}^{N_r \times D}$  and is used as the **domain** 307 308 specific register tokens.

Ŀ

312

283

284

$$\bar{\mathbf{e}}_{k} = \frac{1}{k} \sum_{i=1}^{n} \mathbf{e}_{\delta_{i}}, \quad \{\delta_{1}, \cdots, \delta_{k}\} = \operatorname*{argTopk}_{j=1:H} (\frac{1}{\|\mathbf{x}_{e} - \mathbf{e}_{j}\|_{2}}).$$
(7)

Since the downstream data has its own specific information at the dataset level in addition to the 313 domain level, this may not be fully represented by the domain information obtained from the pre-314 trained dataset alone. Therefore, we innovatively set a learnable matrix  $\mathbf{A} \in \mathbb{R}^{N_{r} \times D}$  to adjust  $\mathbf{X}_{d}$  to 315 complement the **specific information of downstream data**. Since the pre-trained model has a very 316 low intrinsic dimension (Aghajanyan et al., 2020), in order to get better fine-tuning results, A is set 317 as a low-rank matrix: 318

$$\mathbf{A} = \mathbf{u} \times \mathbf{v}^{\mathrm{T}},\tag{8}$$

319 where  $\mathbf{u} \in \mathbb{R}^{N_{\mathrm{r}}}$  and  $\mathbf{v} \in \mathbb{R}^{D}$ , and only the vectors  $\mathbf{u}$  and  $\mathbf{v}$  need to be retrained in the fine-tuning 320 step. As illustrated in Equation 9, the register token  $X_r$  of the downstream scenario is obtained by 321 doing the Hadamard product of  $X_d$ , which represents the domain-specific information obtained at 322 the pre-train stage, and A, which represents the downstream dataset-specific information. 323

$$\mathbf{X}_{\mathrm{r}} = \mathbf{X}_{\mathrm{d}} \odot \mathbf{A}. \tag{9}$$

#### 324 3.4 TRAINING 325

326 To improve the model's prediction ability in zero-shot and few-shot setting, we co-train supervised 327 prediction with self-supervised reconstruction that uses multi-frequency masking to learn unified 328 features that are more applicable to the downstream prediction task. We normalize time series by employing the REVIN (Kim et al., 2021) that is commonly used by the state-of-the-art time series models, which first normalizes each input sample and subsequently applies inverse normalization to 330 recover the model's output. 331

332 **Prediction task.** The input time series  $\mathbf{x} \in \mathbb{R}^L$  is sliced into P non-overlapping patches and then mapped to  $\mathbf{X}_p \in \mathbb{R}^{P \times D}$ . Based on common forecasted needs (Qiu et al., 2024), we set up four 333 prediction heads mapping to prediction lengths of {96, 192, 336, 720} to accomplish the prediction 334 task. Patch tokens  $\mathbf{X}_{p}$  are concatenated with the register tokens  $\mathbf{X}_{u}$  and then successively fed into 335 the Transformer encoder to yield the representation  $\mathbf{S} \in \mathbb{R}^{(N_r+P) \times D}$ : 336

$$\mathbf{S} = \text{Encoder}(\text{Concatenate}(\mathcal{X}_{p}, \mathcal{X}_{u})).$$
(10)

We feed the representation S into the prediction decoder and prediction heads to obtain four prediction results  $\hat{\mathbf{Y}}_F$ , where  $F \in \{96, 192, 336, 720\}$ . With the ground truth  $\mathbf{Y}_F$ , the prediction loss  $\mathcal{L}_{\text{prediction}}$ is shown in Equation 11.

$$\mathcal{L}_{\text{prediction}} = \sum_{F \in \{96, 192, 336, 720\}} ||\mathbf{Y}_F - \hat{\mathbf{Y}}_F||_2^2.$$
(11)

Pre-training. The reconstruction task learns generalized features through the Transformer encoder 346 347 and reconstruction decoder. To utilize these features for the prediction task, the parameters of the reconstruction decoder are copied to the prediction decoder during forward propagation. To 348 avoid prediction training affecting the generalization performance of the model, the gradients of the 349 prediction heads are skipped at back-propagation. The overall loss of ROSE in pre-training stage is shown in Equation 12.

$$\mathcal{L}_{\text{pre-train}} = \mathcal{L}_{\text{reconstruction}} + \mathcal{L}_{\text{prediction}} + \mathcal{L}_{\text{register}}.$$
 (12)

**Fine-tuning.** We only perform a prediction task in fine-tuning. Patch tokens  $X_p$  are concatenated with the adjusted register tokens  $X_r$ . For a downstream task with a fixed prediction length, we use the corresponding pre-trained prediction head to fine-tune the model.

356 357 358

359 360

361

355

350

337 338 339

340

341

#### 4 EXPERIMENTS

**Pre-training datasets.** The datasets are crucial for pre-training a general time series forecasting model. In light of this, we gather a considerable amount of publicly available datasets from various 362 domains, including energy, nature, health, transport, web, and economics, etc. The details of these datasets are shown in the Appendix A.1.1. To enhance data utilization, we downsample fine-grained 364 datasets to coarser granularity, resulting in approximately 887 million time points.

365 **Evaluation datasets.** To conduct comprehensive and fair comparisons for different models, we 366 conducted experiments on seven well-known forecasting benchmarks as the target datasets, including 367 Weather, Traffic, Electricity, and ETT (4 subsets), which cover multiple domains. 368

Baselines. We select the state-of-the-art models as our baselines in full-shot and few-shot setting, 369 including four specific models: iTransformer (Liu et al., 2023), PatchTST (Nie et al., 2022), Times-370 Net (Wu et al., 2022), and DLinear (Zeng et al., 2023), and two LLM-based models: GPT4TS (Zhou 371 et al., 2024) and  $S^{2}$ IP-LLM (Pan et al., 2024). In addition, we select five foundation models for 372 comparison in zero-shot setting, including Timer (Liu et al., 2024), MOIRAI (Woo et al., 2024), 373 Chronos (Ansari et al., 2024), TimesFM (Das et al., 2023b), and Moment (Goswami et al., 2024). 374

375 Setup. Consistent with previous works, we adopted Mean Squared Error (MSE) and Mean Absolute Error (MAE) as evaluation metrics. For fair comparison, all methods fix the look-back window L376 = 512 and predict the future values with lengths  $F = \{96, 192, 336, 720\}$ . More implementation 377 details are presented in the Appendix A.1.3.

# 4.1 IN DISTRIBUTION FORECASTING

Setting. In full-shot setting, we fine-tune pre-trained ROSE and the baselines with the full downstream
 data. In few-shot setting, we fine-tune all models with only 10% train data.

382 Full-shot results. As shown in Table 1, we also present the results of the ROSE in 10% few-shot setting. Key observations are summarized as follows. First, as a general forecasting model, ROSE 384 achieves superior performance compared to the six state-of-the-art baselines with full-data training, 385 achieving an average MSE reduction of 15%, which shows that our decomposed frequency learning 386 and register help to learn generalized representations from large-scale datasets and adaptively transfer 387 the multi-domain information to specific downstream scenarios. Second, we observe that ROSE in 388 10% few-shot setting shockingly *improves a large margin as MSE reduction in average exceeding* 12% over the baselines trained with full data. This observation validates the transferability of ROSE 389 pre-trained with large multi-source data. 390

Table 1: The results for ROSE in full-shot setting and 10% few-shot setting, compared with other methods in full-shot setting. The average results of all predicted lengths are listed here.

Model	s R	OSE   ROSE	E (10%)   ITrans	sformer Patc	hTST   Tim	esnet Dli	near GPT	74TS   S <sup>2</sup> II	P-LLM
	<u> </u>	1.001	3 (10 %)   111um	inter inde					DEM
Metric	MSE   MSE	MAE   MSE	MAE   MSE	MAE   MSE	MAE   MSE	MAE   MSE	MAE   MSE	MAE   MSE	MAE
ETTh	0.391	<b>0.414</b>   <u>0.397</u>	0.419 0.439	0.448   0.413	0.434   0.582	0.533   0.416	0.436   0.427	0.426   0.406	0.427
ETTh2	2   0.331	<b>0.374</b>   <u>0.335</u>	0.380 0.374	0.406 0.331	0.381   0.409	0.438   0.508	0.485   0.354	0.394   0.347	0.391
ETTm	1   0.341	<b>0.367</b>   0.349	<u>0.372</u>   0.362	0.391 0.353	0.382   0.490	0.464   0.356	0.378   0.352	0.383   0.343	0.379
ETTm	2   <b>0.246</b>	<b>0.305</b>   <u>0.250</u>	0.308 0.269	0.329   0.256	0.317   0.317	0.358   0.259	0.325   0.266	0.326   0.257	0.319
Weathe	r   <b>0.217</b>	<b>0.251</b>   0.224	0.252 0.233	0.271   0.226	0.264   0.329	0.336   0.239	0.289   0.237	0.270   0.222	0.259
Electrici	ty   0.155	0.248 0.164	0.253 0.164	0.261 0.159	0.253 0.195	0.296 0.166	0.267 0.167	0.263   0.161	0.257
Traffic	0.390	<b>0.264</b>   0.418	<u>0.278</u>   0.397	0.282   <u>0.391</u>	<b>0.264</b>   0.623	0.333   0.433	0.305   0.414	0.294   0.405	0.286

419

421

396 397

399

**Few-shot results.** The results under the 10% few-shot setting are presented in Table 13 in Appendix A.10.2. ROSE outperforms advanced models when training data is scarce in the target domain.

404 Figure 4 shows the performance of 405 pre-trained ROSE and ROSE trained 406 from scratch on ETTh1 and ETTm2 407 with different fine-tuning data percent-408 ages, noting the best baselines in fullshot setting. The pre-trained ROSE 409 shows stable, superior performance 410 even with limited fine-tuning samples. 411 Specifically, the pre-trained ROSE ex-412 ceeds SOTA performance with only 413 1% train data for ETTh1 and 2% for 414 *ETTm2*. Moreover, compared to the 415 ROSE trained from scratch, the pre-416 trained ROSE exhibits a slower de-



Figure 4: The forecasting results of ROSE obtained by training from scratch and fine-tuning from the pre-trained model. The right, upper corner is the best case.

cline in predictive performance with the reduction of fine-tuning data, demonstrating the impressivegeneralization ability of ROSE through pre-training.

## 420 4.2 ZERO-SHOT FORECASTING

422 Setting. In this section, to ensure a fair comparison, we conduct zero-shot predictions for each 423 foundational model on downstream datasets not included in their pre-training data. It is worth noting 424 that, unlike a few foundation models (Woo et al., 2024) that require much longer inputs to achieve 425 better predictive performance, we fix the input length of all baselines to 512 without considering 426 longer input lengths, as many real-world scenarios could offer very limited samples.

Results. As shown in Table 2, ROSE significantly outperforms across the majority of datasets,
 *achieving an average reduction of 15% in Mean Squared Error (MSE)*. In comparison to Timer
 and Moirai, ROSE achieves average MSE reductions of 9% and 6%, respectively, and demonstrates
 a remarkable 43% relative improvement over Moment. Notably, ROSE stands out not only for its
 superior performance but also for its exceptional lightweight and efficient design, which sets it apart
 from other foundational models. Detailed analysis of these aspects will be presented in Section 4.3.

Table 2: The results for ROSE and other foundation models in zero-shot setting. The average results
of all predicted lengths are listed here. We use '-' to indicate that the dataset has been involved in the
model's pre-training, and thus not used for testing.

Models	RO	SE	Tir	ner   MO	IRAI	Chro	onos	Time	esFM	Mor	nent
Metric	MSE	MAE	MSE	MAE   MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
ETTh1	0.401	0.425	<u>0.451</u>	0.463   0.475	0.443	0.560	0.452	0.489	0.444	0.708	0.580
ETTh2	0.346	0.394	<u>0.366</u>	0.408   0.379	<u>0.396</u>	0.392	0.397	0.396	0.405	0.392	0.430
ETTm1	0.525	<u>0.471</u>	0.544	0.476   0.714	0.507	0.636	0.495	0.434	0.419	0.697	0.555
ETTm2	0.299	0.352	0.360	0.386   0.343	0.356	<u>0.313</u>	0.363	0.320	<u>0.353</u>	0.319	0.360
Weather	0.265	<u>0.305</u>	0.292	0.312   <u>0.267</u>	0.300	0.288	0.310	-	-	0.291	0.323
Electricity	0.234	<u>0.320</u>	0.297	0.375   <u>0.241</u>	0.328	0.245	0.312	-	-	0.861	0.766
Traffic	0.588	0.412	0.613	0.407   -	-	<u>0.371</u>	0.370	-	-	1.411	0.804

#### 4.3 MODEL ANALYSIS

446

Efficiency analysis. To exhibit the performance and efficiency advantages of ROSE, we compare its
 parameter count to other foundation models and evaluate their performance and testing time averaged
 on ETTh1 and ETTh2 datasets in zero-setting. Similarly, for each specific model, we evaluate its
 parameter count as well as its performance in full-shot setting and training-to-testing time averaged
 on the same datasets. Specific implementation details and results can be found in Appendix A.2.

452 As shown in Figure 5, ROSE is a lightweight 453 general model with 7.4M parameters and short inference time, which are only about one-tenth 454 of the second fastest/smallest foundation model 455 (Timer). Importantly, ROSE uses the least num-456 ber of parameters among foundation models, 457 with its parameter count approaching that of spe-458 cific models, while exhibiting superior zero-shot 459 performance. This is attributed to our proposed 460 decomposed frequency learning that enhances 461 the model's comprehension of time series. Con-462 currently, the TS-Register achieves the adaptive 463 transfer thus efficiently adapting to downstream tasks without the need of scaling up to achieve 464 strong generalizability. Compared to foundation 465 models with large scale, ROSE may better meet 466



Figure 5: Model performance, number of parameter and efficiency comparison.

the need for general models in real scenarios that require high computational and parameter efficiency
 as well as high prediction accuracy with scarce downstream data.

Visualization of TS-Register. To validate the TS-Register's capability to transfer domain-specific
 information adaptively from pre-training datasets to target datasets, we visualize the cosine similarity
 of register vector selections from datasets across different domains. As shown in Figure 6(a), the
 cosine similarity is higher for datasets within the same domain and lower between different domains.
 We also visualize the register vector selections from different datasets in Figures 6(b) and (c), where
 datasets from the same domain show similar visualizations. This confirms the TS-Register's capability
 of adaptive transfer from multi-source to target datasets across various domains.

476 Scalability and Sensitivity. The scalability analysis of ROSE's model size and pre-training data size 477 are presented in Appendix A.4. The sensitivity analyses for the upper bound a of the thresholds, the 478 number of masked series  $K_f$ , the number of register tokens  $N_r$ , the size of register H and number of 479 selections k in Top-K strategy are presented in Appendix A.5.

480

482

481 4.4 ABLATION STUDIES

483 Model architecture. To validate effectiveness of our model design, we perform ablation studies on
 484 TS-Register, prediction tasks, and reconstruction task in 10% few-shot setting. Table 3 shows the
 485 impact of each module. The TS-Register leverages multi-domain information during pre-training,
 aiding adaptive transfer to downstream datasets, as further discussed in Section 4.3. The prediction



Figure 6: Visualization of TS-Register. The calculation of cosine similarity is in the Appendix A.3.

tasks enhance performance in data-scarce situations. Without it, performance significantly drops on
ETTh1 and ETTh2 with limited samples. Without the reconstruction task, our model shows negative
transfer effects on ETTm1 and ETTm2, likely due to the prediction task making the model more
susceptible to pre-training data biases.

Masking method. To further validate the effectiveness of decomposed frequency learning, we replace the multi-frequency masking with different masking methods, including two mainstream time-domain methods: patch masking (Nie et al., 2022) and multi-patch masking (Dong et al., 2024), as well as random frequency masking (Chen et al., 2023b). The results in Table 4 show that random frequency masking and patch masking led to negative transfer on ETTm1 and ETTm2, likely due to significant disruption of the original time series, causing overfitting. In contrast, multi-patch masking and multi-frequency masking resulted in positive transfer across all datasets by preventing excessive disruption. Multi-frequency masking achieved better results, demonstrating its ability to help the model understand temporal patterns from a multi-frequency perspective. We also compare with some other pre-training tasks in Table 14 in Appendix A.10.3. 

Table 3: Ablations on key components of model architecture, including TS-register, prediction taskand reconstruction task. The average results of all predicted lengths are listed here.

		ETTm1	ETTm2	ETTh1	ETTh2	
Design		MSE   MAE	MSE   MAE	MSE MAE	MSE   MAE	
ROSE		0.349   0.372	0.250 0.308	0.397 0.419	0.335   0.380	
	TS-Register	0.354   0.378	0.256 0.312	0.418 0.427	0.355   0.390	
w/o	Prediction Task	0.360   0.384	0.257   0.314	0.422 0.438	0.372   0.410	
	Reconstruction Task	0.387   0.403	0.269 0.327	0.412 0.428	0.361   0.399	
	From scratch	0.371   0.391	0.261   0.318	0.470 0.480	0.400   0.425	

Table 4: Ablations on decomposed frequency learning, where we replace Multi-freq masking with other masking methods. The average results of all predicted lengths are listed here.

Design			Tm1	ET	Tm2	ETTh1	ETTh2
			MAE	MSE	MAE	MSE   MAE	MSE   MAE
ROSE			0.372	0.250	0.308	0.397   0.419	0.335 0.380
	Random Freq Masking	0.381	0.397	0.261	0.324	0.410 0.427	0.374 0.405
Replace	Multi-Patch Masking	0.356	0.379	0.259	0.316	0.404   0.426	0.349 0.389
Multi-Frequency Masking	Patch Masking	0.378	0.400	0.261	0.319	0.408 0.432	0.375 0.407
From scratch			0.391	0.261	0.318	0.470   0.480	0.400 0.425

## 5 CONCLUSION AND FUTURE WORK

In this work, we propose ROSE, a novel general model, addressing the challenges of leveraging multi-domain datasets for enhancing downstream prediction task performance. ROSE utilizes decomposed frequency learning and TS-Register to capture generalized and domain-specific representations, en-abling improved fine-tuning results, especially in data-scarce scenarios. Our experiments demonstrate ROSE's superior performance over baselines with both full-data and few-data fine-tuning, as well as its impressive zero-shot capabilities. Future efforts will concentrate on expanding pre-training datasets and extending ROSE's applicability across diverse time series analysis tasks, e.g., classification. We provide our code at https://anonymous.4open.science/r/ROSE-A235.

#### 540 REPRODUCIBILITY 6 541

Our work meets reproducibility requirements. Specifically, you can download our evaluation datasets in a standardized format from the public link (Wang et al., 2024a): https://drive.google. com/drive/folders/13Cq1KYOlzM5C7K8qK8NfC-F3EYxkM3D2, and we provide our code at https://anonymous.4open.science/r/ROSE-A235.

## References

542

543

544

546 547

548

552

553

554

565

566

567

572

578

579

- Xiangfei Qiu, Jilin Hu, Lekui Zhou, Xingjian Wu, Junyang Du, Buang Zhang, Chenjuan Guo, Aoying 549 Zhou, Christian S Jensen, Zhenli Sheng, et al. Tfb: Towards comprehensive and fair benchmarking 550 of time series forecasting methods. arXiv preprint arXiv:2403.20150, 2024. 551
  - Yong Liu, Haoran Zhang, Chenyu Li, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. Timer: Transformers for time series analysis at scale. arXiv preprint arXiv:2402.02368, 2024.
- Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo. 555 Unified training of universal time series forecasting transformers. arXiv preprint arXiv:2402.02592, 556 2024.
- 558 Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski. 559 Moment: A family of open time-series foundation models. arXiv preprint arXiv:2402.03885, 560 2024.
- 561 Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, 562 Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al. 563 Chronos: Learning the language of time series. arXiv preprint arXiv:2403.07815, 2024.
  - Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In International conference on machine learning, pages 27268-27286. PMLR, 2022.
- 568 Xiang Zhang, Ziyuan Zhao, Theodoros Tsiligkaridis, and Marinka Zitnik. Self-supervised contrastive 569 pre-training for time series via time-frequency consistency. Advances in Neural Information 570 Processing Systems, 35:3988–4003, 2022. 571
- Xingyu Zhang, Siyu Zhao, Zeen Song, Huijie Guo, Jianqi Zhang, Changwen Zheng, and Wenwen Qiang. Not all frequencies are created equal: Towards a dynamic fusion of frequencies in time-573 series forecasting. In ACM Multimedia 2024, 2024. 574
- 575 Jiaxiang Dong, Haixu Wu, Haoran Zhang, Li Zhang, Jianmin Wang, and Mingsheng Long. Simmtm: 576 A simple pre-training framework for masked time-series modeling. Advances in Neural Information 577 Processing Systems, 36, 2024.
  - Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. arXiv preprint arXiv:2211.14730, 2022.
- 581 Seunghan Lee, Taeyoung Park, and Kibok Lee. Learning to embed time series patches independently. 582 arXiv preprint arXiv:2312.16427, 2023.
- 583 Liyue Chen, Linian Wang, Jinyu Xu, Shuai Chen, Weiqiang Wang, Wenbiao Zhao, Qiyu Li, and Leye 584 Wang. Knowledge-inspired subdomain adaptation for cross-domain knowledge transfer. In Pro-585 ceedings of the 32nd ACM International Conference on Information and Knowledge Management, 586 pages 234-244, 2023a. 587
- Tian Zhou, Peisong Niu, Liang Sun, Rong Jin, et al. One fits all: Power general time series analysis 588 by pretrained lm. Advances in neural information processing systems, 36, 2024. 589
- George EP Box and Gwilym M Jenkins. Some recent advances in forecasting and control. Journal of 591 the Royal Statistical Society. Series C (Applied Statistics), 17(2):91–109, 1968. 592
- Razvan-Gabriel Cirstea, Bin Yang, and Chenjuan Guo. Graph attention recurrent neural networks for correlated time series forecasting. In MileTS19@KDD, 2019.

594 595	Ruofeng Wen, Kari Torkkola, Balakrishnan Narayanaswamy, and Dhruv Madeka. A multi-horizon quantile recurrent forecaster. <i>arXiv preprint arXiv:1711.11053</i> , 2017.
596 597	David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. <i>International journal of forecasting</i> , 36(3):
599 599	1181–1191, 2020.
600	Densher Luc and Yas Wars. Madamtan A madam and smallting structure for some al time
601	bongnao Luo and Xue wang. Modernich: A modern pure convolution structure for general time series analysis. In The Twelfth International Conference on Learning Representations 2024
602	series analysis. In the tweigh international Conjerence on Learning Representations, 2024.
603	Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. Scinet:
604 605	Time series modeling and forecasting with sample convolution and interaction. Advances in Neural Information Processing Systems, 35:5816–5828, 2022a.
606	Huigiang Wang, Jian Peng, Feihu Huang, Jince Wang, Junhui Chen, and Yifei Xiao, Micn: Multi-scale
607 608	local and global context modeling for long-term series forecasting. In <i>The Eleventh International Conference on Learning Representations</i> , 2022.
609	Haixu Wu Tengge Hu Yong Liu Hang Zhou Jianmin Wang and Mingsheng Long Timesnet:
01U 611	Temporal 2d-variation modeling for general time series analysis. In <i>The eleventh international</i>
612	conference on learning representations, 2022.
613	Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long, Autoformer: Decomposition transformers
614	with auto-correlation for long-term series forecasting. Advances in neural information processing
615	systems, 34:22419–22430, 2021.
616	
617	itransformer: Inverted transformers are effective for time series forecasting arYiv preprint
618	arXiv:2310.06625, 2023.
619	
620	Peng Chen, Yingying Zhang, Yunyao Cheng, Yang Shu, Yihang Wang, Qingsong Wen, Bin Yang,
622	forecasting. 2024.
623	
624 625 626	uan Liang, Yuan-Fang Li, Shirui Pan, et al. Time-Ilm: Time series forecasting by reprogramming large language models. In <i>The Twelfth International Conference on Learning Representations</i> ,
627	2024.
628	Defu Cao, Furong Jia, Sercan O Arik, Tomas Pfister, Yixiang Zheng, Wen Ye, and Yan Liu. Tempo:
629	Prompt-based generative pre-trained transformer for time series forecasting. In <i>The Twelfth</i> International Conference on Learning Representations 2024
630	merhanohai Conjerence on Learning Representations, 2024.
632	Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Arian Khorasani, George Adamopoulos,
633	Rishika Bhagwatkar, Marin Bilos, Hena Ghonia, Nadhir Vincent Hassen, Anderson Schnei- dar et al. Lea llema: Tawarda foundation models for time series forecesting. arXiv propriet
634	arXiv:2310.08278, 2023
635	<i>urxiv.2510.00270, 2025.</i>
636	Samuel Dooley, Gurnoor Singh Khurana, Chirag Mohapatra, Siddartha V Naidu, and Colin White.
637	Forecastpfn: Synthetically-trained zero-shot forecasting. Advances in Neural Information Process-
638	ing Systems, 36, 2024.
639	Azul Garza and Max Mergenthaler-Canseco. Timegpt-1. arXiv preprint arXiv:2310.03589, 2023.
640	Harshavardhan Kamarthi and B Aditya Prakash Large pre-trained time series models for cross-
642	domain time series analysis tasks. arXiv preprint arXiv:2311.11413, 2023.
643	
644	Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for time-series forecasting. arXiv preprint arXiv:2310.10688, 2023a
645	ane series foreeasting. <i>arriv preprint arriv</i> ,2510.10000, 2025a.
646 647	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, 2017.

648 649	E Oran Brigham and RE Morrow. The fast fourier transform. IEEE spectrum, 4(12):63-70, 1967.
650 651	Armen Aghajanyan, Luke Zettlemoyer, and Sonal Gupta. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. <i>arXiv preprint arXiv:2012.13255</i> , 2020.
652	
653	Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-
654	International Conference on Learning Representations 2021
655	merhanonal conjerence on Learning Representations, 2021.
656	Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series
657	forecasting? In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 37, pages
658	11121–11128, 2023.
659	Zijie Pan, Yushan Jiang, Sahil Garg, Anderson Schneider, Yuriy Nevmyvaka, and Dongjin Song. $S^2$
660	ip-llm: Semantic space informed prompt learning with llm for time series forecasting. In Forty-first
662	International Conference on Machine Learning, 2024.
663	Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
664	time-series forecasting. arXiv preprint arXiv:2310.10688, 2023b.
665	
666	Muxi Chen, Zhijian Xu, Ailing Zeng, and Qiang Xu. Fraug: Frequency domain augmentation for
667	time series forecasting. arXiv preprint arXiv:2302.09292, 2023b.
668	Yuxuan Wang, Haixu Wu, Jiaxiang Dong, Yong Liu, Mingsheng Long, and Jianmin Wang. Deep time
669	series models: A comprehensive survey and benchmark. arXiv preprint arXiv:2407.13278, 2024a.
670	Pakabitha Codahawa Christoph Paramair Cooffray I Wahh, Bah I Hundman, and Pahla Montara
671	Manso Monash time series forecasting archive arXiv preprint arXiv:2105.06643, 2021
672	
673	Anthony Bagnall, Hoang Anh Dau, Jason Lines, Michael Flynn, James Large, Aaron Bostrom, Paul
674	Southam, and Eamonn Keogh. The uea multivariate time series classification archive, 2018. arXiv
675	preprint arXiv:1811.00073, 2018.
677	Hoang Anh Dau, Anthony Bagnall, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh
678	Gharghabi, Chotirat Ann Ratanamahatana, and Eamonn Keogh. The ucr time series archive.
679	IEEE/CAA Journal of Automatica Sinica, 6(6):1293–1305, 2019.
680	Shuyi Zhang, Bin Guo, Anlan Dong, Jing He, Ziping Xu, and Song Xi Chen, Cautionary tales on
681	air-quality improvement in beijing. Proceedings of the Royal Society A: Mathematical, Physical
682	and Engineering Sciences, 473(2205):20170457, 2017.
683	Viba Wang, Vy Han, Haichusi Wang, and Viang Zhang. Contrast avarything: A hierorchical
684	contrastive framework for medical time-series Advances in Neural Information Processing
685	Systems, 36, 2024b.
686	
687	Minhao Liu, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia Lai, Lingna Ma, and Qiang Xu. Scinet:
688	Information Processing Systems 35:5816–5828, 2022b
689	<i>Information 1 rocessing Systems</i> , 55.5610–5626, 26226.
601	Michael W McCracken and Serena Ng. Fred-md: A monthly database for macroeconomic research.
692	Journal of Business & Economic Statistics, 34(4):574–589, 2016.
693	Souhaib Ben Taieb, Gianluca Bontempi, Amir F Ativa, and Antti Soriamaa. A review and comparison
694	of strategies for multi-step ahead time series forecasting based on the nn5 forecasting competition.
695	Expert systems with applications, 39(8):7067–7083, 2012.
696	Adam Paszka, Sam Gross, Francisco Massa, Adam Larar, Jamas Pradhury, Gragory Chanan, Travor
697	Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style
698	high-performance deep learning library. Advances in neural information processing systems. 32.
699	2019.
700	Diadamile D Kingma and Kimmy Da Adams A mothed for starburght with the King at
701	Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.

702 703 704	Spyros Makridakis. M4 dataset, 2018. https://github.com/M4Competition/ M4-methods/tree/master/Dataset,.
705	Xiyuan Zhang, Ranak Roy Chowdhury, Jingbo Shang, Rajesh Gupta, and Dezhi Hong. Towards diverse and coherent augmentation for time-series forecasting. In <i>ICASSP</i> 2023-2023 <i>IEEE</i>
700	International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5, IEEE,
707 708	2023.
709	Zhihan Yue. Yujing Wang, Juanyong Duan, Tianmeng Yang, Congrui Huang, Yunhai Tong, and
710	Bixiong Xu. Ts2vec: Towards universal representation of time series. In <i>Proceedings of the AAAI</i>
711	Conference on Artificial Intelligence, volume 36, pages 8980–8987, 2022.
712	
713	
714	
715	
716	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	

## 756 A APPENDIX

# A.1 IMPLEMENTATION DETAILS 759

## 760 A.1.1 PRE-TRAINING DATASETS

761 We use multi-source datasets in pre-training which contain subsets of Monash (Godahewa et al., 762 2021), UEA (Bagnall et al., 2018) and UCR (Dau et al., 2019) time series datasets, as well as 763 some other time series classical datasets (Zhang et al., 2017; Wang et al., 2024b; Liu et al., 2022b; 764 McCracken and Ng, 2016; Taieb et al., 2012). The final list of all pre-training datasets is shown in 765 Table 6. There is no overlap between the pre-training datasets and the target datasets. It is worth 766 noting that the dataset weather in the pre-training dataset is a univariate dataset, which is different 767 to the multivariate dataset weather in the target task. The pre-trained datasets can be categorized 768 into 6 different domains according to their sources: Energy, Nature, Health, Transport, and Web. The sampling frequencies of the datasets show a remarkable diversity, ranging from millisecond 769 samples to monthly samples, which reflects the diverse application scenarios and complexity of the 770 real world. For all pre-training datasets, we split them into univariate sequences and train them in a 771 channel-independent manner. 772

773 774

784

792

A.1.2 EVALUATION DATASETS

775 We use the following 7 multivariate time-series datasets for downstream fine-tuning and forecast-776 ing: ETT datasets<sup>1</sup> contain 7 variates collected from two different electric transformers from July 777 2016 to July 2018. It consists of four subsets, of which ETTh1/ETTh2 are recorded hourly and 778 ETTm1/ETTm2 are recorded every 15 minutes. Traffic<sup>2</sup> contains road occupancy rates measured 779 by 862 sensors on freeways in the San Francisco Bay Area from 2015 to 2016, recorded hourly. Weather<sup>3</sup> collects 21 meteorological indicators, such as temperature and barometric pressure, for Germany in 2020, recorded every 10 minutes. Electricity<sup>4</sup> contains the electricity consumption of 781 321 customers from July 2016 to July 2019, recorded hourly. We split each evaluation dataset into 782 train-validation-test sets and detailed statistics of evaluation datasets are shown in Table 5. 783

Dataset	ETTm1	ETTm2	ETTh1	ETTh2	Traffic	Weather	Electricity
Variables	7	7	7	7	862	21	321
Timestamps	69680	69680	17420	17420	17544	52696	26304
Split Ratio	6:2:2	6:2:2	6:2:2	6:2:2	7:1:2	7:1:2	7:1:2

Table 5: The statistics of evaluation datasets.

### A.1.3 SETTING

We implemented ROSE in PyTorch (Paszke et al., 2019) and all the experiments were conducted on 8 NVIDIA A800 80GB GPU. We used ADAM (Kingma and Ba, 2014) with an initial learning rate of  $5 \times 10^{-4}$  and implemented learning rate decay using the StepLR method to implement learning rate decaying pre-training. By default, ROSE contains 3 encoder layers and 3 decoder layers with head number of 16 and the dimension of latent space D = 256. The patch size for patching is set to 64.

**Pre-training.** We use  $N_r = 3$  as the number of register tokens and P = 8 as the path tokens. We set the input length to 512 for the supervised prediction task with target lengths of 96, 192, 336, and 720. We also set the input length to 512 and mask number  $K_f = 4$ . The batch size is set to 8192 in pre-training.

Fine-tuning. We fix the lookback window to 512, and perform predictions with target lengths of 96, 192, 336, and 720, respectively. The number of register tokens  $N_r$  and patch tokens P is the same as in pre-training, and the parameter k = 3 in TopK is set when selection vectors are performed in the register.

<sup>&</sup>lt;sup>1</sup>https://github.com/zhouhaoyi/ETDataset

<sup>808 &</sup>lt;sup>2</sup>https://pems.dot.ca.gov/

<sup>&</sup>lt;sup>3</sup>https://www.bgc-jena.mpg.de/wetter/

<sup>&</sup>lt;sup>4</sup>https://archive.ics.uci.edu/ml/datasets/ElectricityLoadDiagrams20112014

811			1	C	
812	Domain	Dataset	Frequency	Time Pionts	Source
813		Aus. Electricity Demand	Half Hourly	1155264	Monash(Godahewa et al., 2021)
814		Wind	4 Seconds	7397147	Monash(Godahewa et al., 2021)
815	Energy	Wind Farms	Minutely	172178060	Monash(Godahewa et al., 2021)
816		Solar	10 Minutes	7200720	Monash(Godahewa et al., 2021)
817		Solar Power	4 Seconds	7397222	Monash(Godahewa et al., 2021)
818		London Smart Meters	Half Hourly	166527216	Monash(Godahewa et al., 2021)
819		Phoneme	-	2160640	UCRDau et al. (2019)
820		EigenWorms	-	27947136	UEA(Bagnall et al., 2018)
821		PRSA	Hourly	4628448	(Zhang et al., 2017)
822		Temperature Rain	Daily	23252200	Monash(Godahewa et al., 2021)
823		StarLightCurves	-	9457664	UCR(Dau et al., 2019)
824	Natura	Worms	0.033 Seconds	232200	UCR(Dau et al., 2019)
825	Nature	Saugeen River Flow	Daily	23741	Monash(Godahewa et al., 2021)
826		Sunspot	Daily	73924	Monash(Godahewa et al., 2021)
827		Weather	Daily	43032000	Monash(Godahewa et al., 2021)
828		KDD Cup 2018	Daily	2942364	MonashGodahewa et al. (2021)
220		US Births	Daily	7305	Monash(Godahewa et al., 2021)
029		MotorImagery	0.001 Seconds	72576000	UEA(Bagnall et al., 2018)
030		SelfRegulationSCP1	0.004 Seconds	3015936	UEA(Bagnall et al., 2018)
001		SelfRegulationSCP2	0.004 Seconds	3064320	UEA(Bagnall et al., 2018)
832	Health	AtrialFibrillation	0.008 Seconds	38400	UEA(Bagnall et al., 2018)
833		PigArtPressure	-	624000	UCR(Dau et al., 2019)
834		PIGCVP	-	624000	UCR(Dau et al., 2019)
835		TDbrain	0.002 Seconds	79232703	(Wang et al., 2024b)
836		Pems03	5 Minute	9382464	(Liu et al., 2022b)
837		Pems04	5 Minute	5216544	(Liu et al., 2022b)
838		Pems07	5 Minute	24921792	(Liu et al., 2022b)
839	Transport	Pems08	5 Minute	3035520	(Liu et al., 2022b)
840		Pems-bay	5 Minute	16937700	(Liu et al., 2022b)
841		Pedestrian_Counts	Hourly	3132346	Monash(Godahewa et al., 2021)
842	Web	Web Traffic	Daily	116485589	Monash(Godahewa et al., 2021)
843		FRED_MD	Monthly	77896	(McCracken and Ng, 2016)
844	Economic	Bitcoin	Daily	75364	Monash(Godahewa et al., 2021)
845		NN5	Daily	87801	(Taieb et al., 2012)
846					

## Table 6: List of pretraining datasets.

### A.1.4 BASELINES

We select the state-of-the-art models as our baselines in full-shot and few-shot setting, including four specific models: iTransformer (Liu et al., 2023), PatchTST (Nie et al., 2022), TimesNet (Wu et al., 2022), and DLinear (Zeng et al., 2023), and two LLM-based models: GPT4TS (Zhou et al., 2024) and  $S^2$ IP-LLM (Pan et al., 2024). In addition, we selected five foundation models for comparison in zero-shot setting, including Timer (Liu et al., 2024), MOIRAI (Woo et al., 2024), Chronos (Ansari et al., 2024), TimesFM (Das et al., 2023b) and Moment (Goswami et al., 2024). The specific code base for these models is listed in Table 7:

A.2 EFFICIENCY ANALYSIS

As an important aspect of foundation models, inference efficiency is crucial. Therefore, we evaluate
the testing time of ROSE and five foundation models in the ETTh1 and ETTh2 dataset in zero-shot
setting. Similarly, we evaluate the time of the entire process of training, validation, and testing for
four specific models in the same datasets in full-shot setting. The above experiments all set the batch
size to 32. The specific results are shown in Table 8 and Figure 5. We observe that ROSE maintains
its advantage in zero-shot performance while also being significantly faster compared to the baselines,

866	Model Types	Models	Code Repositories
867		iTransformer	https://github.com/thuml/iTransformer
868		PatchTST	https://github.com/yuqinie98/PatchTST
869	Small Model	TimesNet	https://github.com/thuml/TimesNet
870		Dlinear	https://github.com/cure-lab/LTSF-Linear
871		Timer	https://github.com/thuml/Large-Time-Series-Model
872		MOIRAI	https://github.com/redoules/moirai
873		Chronos	https://github.com/amazon-science/chronos-forecasting
874	Foundation Model	TimesFM	https://github.com/google-research/timesfm/
875		Moment	https://anonymous.4open.science/r/BETT-773F/README.md
876		GPT4TS	https://github.com/DAMO-DI-ML/NeurIPS2023-One-Fits-All
877	LLM-based Model	S2IP-LLM	https://github.com/panzijie825/S2IP-LLM
878			

### Table 7: Code repositories for baselines.

even being approximately ten times faster than the second-fastest foundation model, Timer. This raises our reflection on whether time-series foundation models require extremely large parameter sizes and whether existing time-series foundation models have validated their architectures' scaling laws on time-series data.

Table 8: Efficiency analysis.

Model	Parameters	Pre-train datasize	Averaged time
ROSE	7.4M	0.89B	0.652s
MOIRAI	311M	27B	7.920s
Timer	67.4M	1B	5.989s
Chronos	46M	84B	176s
TimesFM	200M	100B	10.5s
Moment	385M	1.13B	13s
Itransformer	3.8M	-	34.18s
PatchTST	3.2M	-	35.47s
TimesNet	1.8M	-	146s
Dlinear	0.018M	-	24.06s

### A.3 CALCULATION OF COSINE SIMILARITY

In Figure 6, we visualize the cosine similarity of register vector selections. For each sample in a dataset, during the inference process, k vectors are selected from the register based on the Top-K strategy. We iterate through all samples in the dataset and count the number of times each vector is selected, which allows us to obtain a record vector of length equivalent to the size of the register for each dataset. The  $i_{th}$  position in the record vector represents the number of times the  $i_{th}$  vector in the register has been selected by samples in the dataset. For a pair of datasets, we can obtain a unique record vector for each dataset, and then we are able to calculate the cosine similarity of two vectors.

906 907

864

865

879 880

883

885

A.4 SCALABILITY

Scalability is crucial for a general model, enabling significant performance improvements by expanding pre-training data and model sizes. To investigate the scalability of ROSE, we increased both the model size and dataset size and evaluated its predictive performance on four ETT datasets.

Model size. Constrained by computational resources, we use 40% pre-training datasets. The results are shown in Figure 7(a) and (b). When maintaining the model dimension, we increased the model layers, increasing model parameters from 2M to 4.5M. This led to 10.37% and 9.34% improvements in the few-shot scenario with 5% and 10% downstream data, respectively.

Data size. When keeping the model size, we increase the size of the pre-training datasets from 178M to 887M. The results are shown in Figure 7(c) and (d). The performance of our model steadily improves with the increase in dataset size and achieves improvements of 7.4% and 4.8% respectively.



Figure 7: (a)/(b): Larger ROSE demonstrates better performance on downstream forecasting. (c)/(d): ROSE pre-trained on larger datasets demonstrates better performance on downstream forecasting.

929 A.5 SENSITIVITY

927

928

930

931

932

933

934

958

959

960

961

962

963

964

965

We perform the sensitivity analyses for the upper bound a of the thresholds, the number of masked series  $K_f$ , the number of register tokens  $N_r$ , the size of register H and the number of selections k in Top-K strategy. All the sensitivity experiments present the average results on the four ETT datasets: ETTh1, ETTh2, ETTm1 and ETTm2 under 10% few-shot setting.

Number of masked series. As described in Section 3.2, we propose decomposed frequency learning, 935 which employs multiple thresholds to randomly mask high and low frequencies in the frequency do-936 main, thereby decomposing the original time series into multiple frequency components. This allows 937 the model to understand the time series from multiple frequency perspectives. In this experiment, we 938 study the influence of the number of masked series  $K_{\rm f}$  on downstream performance. We train ROSE 939 with 1, 2, 3, 4, 5, or 6 mask series. We report the results of this analysis in Figure 8(a). We find that as 940 the number of masked sequences increases, the downstream performance gradually improves. This is 941 because the model can better understand the time series from the decomposed frequency components, 942 which enhances the model's generalization ability. However, more masked series do not bring better 943 downstream performance. This could be due to an excessive number of masked sequences leading to 944 information redundancy. In all our experiments, we keep 4 mask series.

945 Number of register tokens. The TS-register module presented in Section 3.3 supports the con-946 figuration of an arbitrary number of register tokens. In Figure 8(b), we visualize the relationship 947 between the performance on the ETT datasets under a 10% few-shot setting and the number of 948 register tokens. It is observed that when the number of register tokens ranges from 1 to 6, the 949 model's performance remains relatively stable, with an optimal outcome achieved when the number 950 is set to 3. This phenomenon may be because when the number of register tokens is too small, they contain insufficient domain-specific information, which limits their effectiveness in enhancing the 951 model's performance. Conversely, an excess of register tokens may introduce redundant information, 952 hindering the accurate representation of domain-specific information. Additionally, we compared 953 the results without the adjustment of a low-rank matrix on the register tokens and found that the 954 incorporation of a low-rank matrix adjustment led to improvements across all quantities of register 955 tokens. This finding underscores the significance of utilizing a low-rank matrix to supplement the 956 register tokens with downstream data-specific information. 957



Figure 8: (a): Analysis of the number of masked series. (b): Analysis of the number of register tokens.

969 Thresholds upper bound. Figure 9(a) illustrates the relationship between threshold upper bound and
 970 model performance. We have observed that the upper bound of the threshold has a minimal impact on
 971 the model's performance. Generally, the information density is higher in low-frequency components
 compared to high-frequency ones. Therefore, the upper bound of the threshold should be biased

18

towards the low-frequency range to balance the information content between low-frequency and high-frequency components. However, this bias should not be excessive. Our experiments indicate that an upper bound of L/10 performs worse than L/5 as an overly left-skewed threshold results in insufficient information in the low-frequency range, making the reconstruction task either too difficult or too simple. Based on our findings, we recommend using L/5 as the upper bound for the threshold.

977 Register size. Figure 9(b) illustrates the relationship between register size and model performance. The register size determines the upper limit of domain-specific information that the register can store.
979 We can observe that there is a significant improvement in the model effect when the register size is increased from 32 to 128. When the register size exceeds 128, the improvement of the model effect with the increase of register size is no longer obvious. Therefore, we believe that 128 is an appropriate register size for the current pre-training datasets.

**Number of selections in Top-K strategy.** Figure 9(c) illustrates the relationship between the number of selections k in Top-K strategy and model performance when we use the register to realize adaptive transfer of domain-specific information in downstream tasks. It can be seen that the model effect performance peaks at 3 tokens at k = 3, which has some advantages over selecting once (k = 1), indicating that the TopK strategy can compensate for the problem of incomplete matching of upstream and downstream domains to some extent. However, too large k will also introduce redundant information and limit the accuracy of domain-specific information transfer.



Figure 9: (a): Analysis of the threshold upper bound. (b): Analysis of the size of register. (c): Analysis of the number of selections in Top-K strategy.

### A.6 SHORT-TERM FORECASTING

We also try to apply ROSE to short-term forecasting on the M4 (Makridakis, 2018) dataset, which contains the yearly, quarterly and monthly collected univariate marketing data. We follow Times-Net's (Wu et al., 2022) setting and metrics (SMAPE, MASE and OWA) for testing. As shown in Table 9, ROSE also exhibits competitive performance on the M4 dataset compared to the baselines.

		ROSE		iTr	ansformer		Р	atchTST		1	limesnet			Dlinear		•	GPT4TS		5	2IP-LLM	
Metric	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA	SMAPE	MASE	OWA
Yearly	13.302	3.014	0.833	13.238	2.952	0.823	16.766	4.331	1.018	13.387	2.996	0.786	16.965	4.283	1.058	13.531	3.015	0.793	13.413	3.024	0.792
Quarterly	9.998	1.165	0.885	10.001	1.278	0.949	12.132	1.513	0.966	10.100	1.182	0.890	12.145	1.520	1.106	10.100	1.194	0.898	10.352	1.228	0.922
Monthly	12.650	0.915	0.866	13.399	1.031	0.949	13.428	0.997	0.948	12.670	0.933	0.933	13.514	1.037	0.956	12.894	0.956	0.897	12.995	0.970	0.910
Others	4.668	3.126	1.020	6.558	4.511	1.401	6.667	4.834	1.417	4.891	3.302	1.035	6.709	4.953	1.487	4.940	3.228	1.029	4.805	3.247	1.071

#### A.7 MODEL GENERALITY

We evaluate the effectiveness of our proposed multi-frequency masking on Transformer-based models and CNN-based models, whose results are shown in Table 10. It is notable that multifrequency masking consistently improves these forecasting models. Specifically, it achieves average improvements of 6.3%, 3.7%, 1.5% in Autoformer (Wu et al., 2021), TimesNet (Wu et al., 2022), and PatchTST (Nie et al., 2022), respectively. This indicates that multi-frequency Masking can be widely utilized across various time series forecasting models to learn generalized time series representations and improve prediction accuracy. 

Table 10: Performance of multi-frequency masking.

Datasets	ETTm1	ETTm2	ETTh1	ETTh2
Metric	MSE   MAE	MSE   MAE	MSE   MAE	MSE   MAE
Autoformer	0.600   0.521	0.328   0.365	0.493   0.487	0.452   0.458
+Multi-frequency Masking	0.549 0.488	0.306   0.349	0.474   0.478	0.406 0.425
TimesNet	0.400   0.406	0.291   0.333	0.458   0.450	0.414   0.427
+Multi-frequency Masking	0.386 0.398	0.282 0.324	0.446   0.438	0.386 0.403
PatchTST	0.353   0.382	0.256   0.317	0.413   0.434	<b>0.331</b>   0.381
+Multi-frequency Masking	0.347 0.372	0.252 0.308	0.405   0.424	0.337 0.379

#### A.8 **RESULTS DEVIATION**

We have conducted ROSE three times with different random seeds and have recorded the standard deviations for both the full-shot setting and the 10% few-shot setting, as illustrated in Table 11. As the baselines didn't report deviations in the original paper, we only reported the deviations of the PatchTST in the full-shot setting as a comparison. It can be observed that ROSE exhibits stable performance. 

Table 11: Results deviation.

1055					Results devi	auon.		
1056	Models	RC	DSE	ROSE	(10%)	Patch	nTST	confidence interval
1057	Metric	MSE	MAE	MSE	MAE	MSE	MAE	-
1058	ETTm1	0.342±0.003	0.367±0.002	$0.349 \pm 0.003$	0.372±0.002	$0.349 {\pm}~0.004$	0.383±0.003	99%
1059	ETTm2	$0.246 \pm 0.002$	$0.303 \pm 0.004$	$0.249 \pm 0.002$	$0.308 \pm 0.002$	$0.255 {\pm} 0.002$	$0.314 \pm 0.003$	99%
1060	ETTh1	$0.392 {\pm} 0.004$	$0.413 \pm 0.004$	$0.397 {\pm} 0.003$	$0.419 \pm 0.003$	$0.411 \pm 0.003$	$0.432 \pm 0.005$	99%
1000	ETTh2	$0.330 {\pm} 0.003$	$0.374 \pm 0.002$	$0.335 {\pm} 0.004$	$0.380 \pm 0.003$	$0.348 {\pm} 0.004$	$0.390 \pm 0.004$	99%
1061	Traffic	$0.391 \pm 0.008$	$0.266 \pm 0.005$	$0.418 \pm 0.011$	$0.278 \pm 0.006$	$0.404 \pm 0.009$	$0.283 \pm 0.002$	99%
1062	Weather	$0.217 \pm 0.008$	$0.250 \pm 0.007$	$0.224 \pm 0.007$	$0.252 \pm 0.009$	$0.223 \pm 0.011$	$0.263 \pm 0.014$	99%
1000	Electricity	$0.156 \pm 0.007$	$0.249 \pm 0.009$	$0.164 \pm 0.004$	$0.253 \pm 0.004$	$0.163 \pm 0.009$	$0.261 \pm 0.013$	99%
1063								
1064								
1065								
1005								
1066								
1067								
1068								
1069								
1070								
1071								

1080 A.9 VISUALIZATION

1113

1114

1082 A.9.1 VISUALIZATION ANALYSIS

To showcase the benefits of cross-domain pre-training, we performed visualizations in both the zero-shot setting and full-shot setting.

Zero-shot: We pre-train the baselines iTransformer and PatchTST on the energy domain dataset
ETTm1 and test their zero-shot performance on two different domains (weather, traffic). ROSE,
without fine-tuning, is evaluated to the same two test-sets. As shown in the Figure 10, we find that
the baselines generally perform worse during domain shift due to their poor generalization. However,
ROSE excels in scenarios across all domains, which demonstrates the benefits of cross-domain
pre-training for improving generalization.

Full-shot: We train the baselines on the train-set of downstream dataset ETTh2 and fine-tune ROSE
on the same train-set. As shown in the Figure 11, We find that the baselines is limited by data diversity,
leading to poor performance on patterns which rarely appear. However, ROSE excels in these cases,
as the cross-domain pre-training allows ROSE to learn diverse temporal patterns, and helps ROSE to
predict the patterns which rarely appear in the downstream train-set well.



Figure 10: Visualization comparison of ROSE with cross-domain pre-training and other SOTA baselines in the zero-shot setting for three domain datasets.



Figure 11: Visualization comparison of ROSE with cross-domain pre-training and other SOTA baselines in the full-shot setting for rare and common patterns.

# 1134 A.9.2 VISUALIZATION SHOWCASE

To provide a distinct comparison among different models, we present visualizations of the forecasting results on the ETTh2 dataset and the weather dataset in different settings, as shown in Figures 12 to Figures 15, given by the following models: DLinear (Zeng et al., 2023), TimesNet (Wu et al., 2022), iTransformer (Liu et al., 2023), and PatchTST (Nie et al., 2022). Among the methods, ROSE demonstrates the most accurate prediction ability.



Figure 12: Visualization of input-512 and predict-336 forecasting results on the ETTh2 dataset in full-shot setting.



Figure 13: Visualization of input-512 and predict-336 forecasting results on the ETTh2 dataset in 10% few-shot setting.



Figure 14: Visualization of input-512 and predict-336 forecasting results on the weather dataset in full-shot setting.



Figure 15: Visualization of input-512 and predict-336 forecasting results on the weather dataset in 10% few-shot setting.

# 1242 A.10 FULL RESULTS

# 1244 A.10.1 Full-shot results 1245

Model	s	RC	DSE	ITrans	former	Patch	nTST	Tim	esnet	Dli	near	GP1	T4TS	S
Metric	:	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	M
	96	0.354	0.385	0.386	0.405	0.370	0.400	0.470	0.470	0.367	<u>0.396</u>	0.376	0.397	<u>0.3</u>
	192	0.389	0.407	0.424	0.440	0.413	0.429	0.568	0.523	<u>0.400</u>	<u>0.417</u>	0.416	0.418	0.4
ETTh1	336	0.406	0.422	0.449	0.460	0.422	0.440	0.595	0.547	0.428	0.439	0.442	0.433	<u>0.4</u>
	720	0.413	0.443	0.495	0.487	0.447	0.468	0.694	0.591	0.468	0.491	0.477	<u>0.456</u>	<u>0.</u>
	avg	0.391	0.414	0.439	0.448	0.413	0.434	0.582	0.533	0.416	0.436	0.427	0.426	<u>0.</u>
	96	0.265	0.320	0.297	0.348	0.274	0.337	0.351	0.399	0.302	0.368	0.285	0.342	0.
	192	<u>0.328</u>	0.369	0.371	0.403	0.341	0.382	0.394	0.429	0.404	0.433	0.354	0.389	0.
ETTh2	336	0.353	0.391	0.404	0.428	0.329	0.384	0.415	0.443	0.511	0.498	0.373	0.407	0.
	720	0.376	0.417	0.424	0.444	0.379	0.422	0.477	0.481	0.815	0.640	0.406	0.441	0.
	avg	0.331	0.374	0.374	0.406	0.331	<u>0.381</u>	0.409	0.438	0.508	0.485	0.354	0.394	<u>0.</u>
	96	0.275	0.328	0.300	0.353	0.293	0.346	0.405	0.421	0.303	0.346	0.292	0.262	<u>0</u> .
	192	0.324	0.358	0.345	0.382	0.333	0.370	0.508	0.473	0.335	0.365	0.332	0.301	0.
FTTm1	336	0.354	0.377	0.374	0.398	0.369	0.392	0.523	0.479	0.365	0.384	0.366	0.341	<u>0</u> .
211111	720	0.411	0.407	0.429	0.430	0.416	0.420	0.523	0.484	0.418	0.415	0.417	0.401	0.
	avg	0.341	0.367	0.362	0.391	0.353	0.382	0.490	0.464	0.356	0.378	0.352	0.383	0.
	96	0.157	0.243	0.175	0.266	0.166	0.256	0.233	0.305	0.164	0.255	0.173	0.262	0
	192	0.213	0.283	0.242	0.312	0.223	0.296	0.265	0.328	0.224	0.304	0.229	0.301	<u>0</u> .
FTTm2	336	0.266	0.319	0.282	0.340	0.274	0.329	0.379	0.392	0.277	0.339	0.286	0.341	0
211112	720	0.347	0.373	0.378	0.398	0.362	0.385	0.390	0.407	0.371	0.401	0.378	0.401	0.
	avg	0.246	0.305	0.269	0.329	0.256	0.317	0.317	0.358	0.259	0.325	0.266	0.326	0.
	96	0.145	0.182	0.159	0.208	0.149	0.198	0.193	0.244	0.170	0.230	0.162	0.212	0.
	192	0.183	0.226	0.200	0.248	0.194	0.241	0.320	0.329	0.212	0.267	0.204	0.248	0.
Weather	336	0.232	0.267	0.253	0.289	0.245	0.282	0.363	0.366	0.257	0.305	0.254	0.286	0.
weather	720	0.309	0.327	0.321	0.338	0.314	0.334	0.440	0.404	0.318	0.356	0.326	0.337	0.
	avg	0.217	0.251	0.233	0.271	0.226	0.264	0.329	0.336	0.239	0.289	0.237	0.270	<u>0</u> .
	96	0.125	0.220	0.138	0.237	0.129	0.222	0.182	0.287	0.141	0.241	0.139	0.238	0.
	192	0.142	0.235	0.157	0.256	0.147	0.240	0.193	0.293	0.154	0.254	0.153	0.251	0
Electricity	336	0.162	0.252	0.167	0.264	0.163	0.259	0.196	0.298	0.168	0.271	0.169	0.266	0
Bieeurienty	720	0.191	0.284	0.194	0.286	0.197	0.290	0.209	0.307	0.203	0.303	0.206	0.297	0
	avg	0.155	0.248	0.164	0.261	0.159	0.253	0.195	0.296	0.166	0.267	0.167	0.263	0.
	96	0.354	0.252	0.363	0.265	0.360	0.249	0.611	0.323	0.411	0.294	0.388	0.282	0.
	192	0.377	0.257	0.385	0.273	0.379	0.256	0.609	0.327	0.421	0.298	0.407	0.290	0.
Traffic	336	0.396	0.262	0.396	0.277	0.392	0.264	0.616	0.335	0.431	0.304	0.412	0.294	0.
maille	720	0.434	0.283	0.445	0.312	0.432	0.286	0.656	0.349	0.468	0.325	0.450	0.312	0.
	avg	0.390	0.264	0.397	0.282	0.391	0.264	0.623	0.333	0.433	0.305	0.414	0.294	0.

# 1296 A.10.2 Few-shot results 1297

1298 1299					Table	13: F	ull res	ults in	10% 1	few-sh	not set	ting				
1300	Models		R	OSE	ITrans	former	Pate	hTST	Tim	esnet	Dli	near	GPI	TATS	S <sup>2</sup> IP	JIM
1301	Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
1302		96	0.367	0.395	0.442	0.464	0.458	0.463	0.579	0.522	1.355	0.816	0.458	0.456	0.481	0.474
1303		192	0.399	0.416	0.476	0.475	0.481	0.490	0.641	0.553	1.210	0.825	0.570	0.516	0.518	0.491
130/	ETTh1	336	0.405	0.423	0.486	0.482	0.465	0.475	0.721	0.582	1.487	0.914	0.608	0.535	0.664	0.570
1205	21111	720	0.416	0.443	0.509	0.506	<u>0.478</u>	<u>0.492</u>	0.630	0.574	1.369	0.826	0.725	0.591	0.711	0.584
1206		avg	0.397	0.419	0.478	0.482	<u>0.470</u>	<u>0.480</u>	0.643	0.558	1.355	0.845	0.590	0.525	0.593	0.529
1207		96	0.273	0.332	0.333	0.385	0.350	0.389	0.378	0.413	1.628	0.724	<u>0.331</u>	<u>0.374</u>	0.354	0.400
1307		192	0.334	0.376	0.402	0.428	0.416	0.426	0.463	0.460	1.388	0.713	0.402	<u>0.411</u>	<u>0.400</u>	0.423
1308	ETTh2	336	0.358	0.397	0.438	0.452	<u>0.401</u>	0.429	0.507	0.495	1.595	0.772	0.406	0.433	0.442	0.450
1309		720	0.376	0.417	0.466	0.477	0.436	0.457	0.516	0.501	1.664	0.857	0.449	0.464	0.480	0.486
1310		avg	0.335	0.380	0.410	0.436	0.401	0.425	0.466	0.467	1.569	0.766	0.397	0.421	0.419	0.439
1311		96	0.287	0.336	0.353	0.392	0.317	0.363	0.481	0.446	0.454	0.475	0.390	0.404	0.388	0.401
1312		192	0.351	0.362	0.385	0.410	0.351	0.308	0.621	0.491	0.575	0.548	0.429	0.423	0.422	0.421
1313	ETTm1	720	0.302	0.379	0.422	0.432	0.370	0.398	0.521	0.479	0.773	0.031	0.409	0.439	0.430	0.430
1314		720   avg	0.410	0.412	0.494	0.472	0.455	0.450	0.571	0.508	0.945	0.710	0.309	0.498	0.554	0.435
1315		96	0.159	0.247	0.183	0.420	0.170	0.259	0.545	0.401	0.000	0.575	0.188	0.441	0.192	0.433
1316		192	0.217	0.287	0.247	0.320	0.226	0.297	0.212	0.353	0.923	0.658	0.251	0.309	0.246	0.313
1317	ETT2	336	0.269	0.322	0.300	0.353	0.284	0.333	0.328	0.364	1.407	0.822	0.307	0.346	0.301	0.340
1318	E11m2	720	0.357	0.377	0.385	0.408	0.363	0.382	0.456	0.440	1.626	0.905	0.426	0.417	0.400	0.403
1319		avg	0.250	0.308	0.279	0.340	0.261	0.318	0.323	0.362	1.112	0.715	0.293	0.335	0.284	0.332
1320		96	0.145	0.184	0.189	0.229	0.166	0.217	0.199	0.248	0.230	0.318	0.163	0.215	0.159	0.210
1321		192	0.190	0.227	0.239	0.269	0.211	0.257	0.249	0.285	0.357	0.425	0.210	0.254	0.200	0.251
1222	Weather	336	0.245	0.269	0.294	0.308	0.261	0.296	0.297	0.316	0.464	0.493	0.256	0.292	0.257	0.293
1000		720	0.317	0.328	0.366	0.356	0.328	0.342	0.367	0.361	0.515	0.532	0.321	0.339	<u>0.317</u>	0.335
1323		avg	0.224	0.252	0.272	0.291	0.242	0.278	0.278	0.303	0.391	0.442	0.238	0.275	<u>0.233</u>	0.272
1324		96	0.135	0.226	0.184	0.276	0.161	0.256	0.279	0.359	0.227	0.334	<u>0.139</u>	0.237	0.143	0.243
1325		192	0.150	0.240	0.192	0.284	0.163	0.257	0.282	0.363	0.265	0.366	<u>0.156</u>	0.252	0.159	0.258
1326	Electricity	336	0.100	0.258	0.216	0.308	0.173	0.266	0.289	0.367	0.339	0.41/	0.175	0.270		0.269
1327		/20	0.205	0.290	0.265	0.347	0.221	0.313	0.333	0.399	0.482	0.478	0.233		0.230	0.315
1328		06	0.104	0.233	0.214	0.304	0.180	0.275	0.290	0.372	0.528	0.399	0.170	0.209	0.175	0.271
1329		192	0.390	0.270	0.473	0.319	0.439	0.313	0.705	0.393	0.710	0.480	0.426	0.201	0.412	0.295
1330	Traffa	336	0.417	0.277	0.491	0.329	0.448	0.318	0.863	0.456	0.723	0.481	0.434	0.303	0.427	0.316
1331	manic	720	0.452	0.294	0.536	0.361	0.478	0.320	0.928	0.485	0.673	0.436	0.487	0.337	0.469	0.325
1332		avg	0.418	0.278	0.490	0.331	0.447	0.312	0.801	0.430	0.680	0.446	0.440	0.310	0.427	0.307
1333																<u>.</u>
1334																
1335																
1336																
1337																
1338																
1339																
1340																
13/1																
12/0																
1042																
1343																
1344																
1345																
1346																
1347																
1348																
1349																

# 1350 A.10.3 ABLATION STUDY RESULTS

**Novelty of decomposed frequency learning.** Frequency masking is not a new concept, but past approaches randomly mask frequencies of a single time series once (Chen et al., 2023b; Zhang et al., 2023), which show limited forecasting effectiveness due to the lack of common pattern learning from heterogeneous time series that come from various domains. While the multi-frequency masking we proposed randomly mask either high-frequency or low-frequency components of a time series multiple times as the key to enable learning of common time series patterns, such as trends and various long and short term fluctuations. Moreover, different from utilizing frequency masking as a way of data augmentation to enhance the diversity of input data (Chen et al., 2023b; Zhang et al., 2023), we combine multi-frequency masking with reconstruction task as a novel pre-training framework, that learns a universal and unified feature representation by comprehending the data from various frequency perspectives, thereby enabling it to learn generalized representations.

Difference between frequency-domain masking and time-domain noise addition. Multi-frequency masking and reconstruction are not equivalent to the pre-training methods of adding noise and denoising (Noise). Due to the sparsity of time series, the process of adding noise and denoising may potentially disrupt the information of original time series (Dong et al., 2024). In contrast, multi-frequency masking not only preserves the series from such disruption but also helps the model understand temporal patterns from a multi-frequency perspective, thereby helping the model to learn general features better.

Other pre-training tasks. Based on the two points above, we conduct experiments to compare two other pre-training tasks: 1) using frequency-domain augmentation only for data expansion without reconstruction task (*Aug*); 2) replacing multi-frequency masking and reconstruction task with adding time-domian noise and denoise task (*Noise*). As shown in Table 14, we find that ROSE is significantly more effective than *Aug* and *Noise*, which demonstrates the effectiveness of multi-frequency masking and reconstruction task in learning generalized features.

			ET	Tm1	ET	Tm2	ET	Th1	1
D	esign	Pred_len	MSE	MAE	MSE	MAE	MSE	MAE	M
		96	0.287	0.336	0.159	0.247	0.367	0.395	0.2
		192	0.331	0.362	0.217	0.287	0.399	0.416	0.3
R	OSE	336	0.362	0.379	0.269	0.322	0.405	0.423	0.3
		720	0.416	0.412	0.357	0.377	0.416	0.443	0.3
		avg	0.349	0.372	0.250	0.308	0.397	0.419	0.3
		96	0.330	0.370	0.170	0.262	0.391	0.399	0.3
		192	0.352	0.392	0.232	0.298	0.406	0.430	0.3
	Random Frequency Masking	336	0.390	0.389	0.276	0.342	0.411	0.432	0.4
		720	0.452	0.438	0.366	0.392	0.432	0.447	0.4
		avg	0.381	0.397	0.261	0.324	0.410	0.427	0.3
		96	0.302	0.348	0.168	0.257	0.377	0.408	0.2
		192	0.336	0.367	0.228	0.297	0.404	0.423	0.3
Replace	Multi-Patch Masking	336	0.364	0.385	0.277	0.328	0.405	0.420	0.3
man-riequency masking	 	/20	0.423	0.410	0.304	0.381	0.431	0.433	0.3
		avg	0.550	0.379	0.239	0.310	0.404	0.420	
		102	0.318	0.300	0.100	0.239	0.388	0.412	0.
		336	0.355	0.500	0.220	0.290	0.402	0.422	0.2
	Patch Masking	720	0.500	0.100	0.279	0.388	0.111	0.155	0.4
		avg	0.378	0.400	0.261	0.319	0.408	0.432	0.3
		96	0.304	0.357	0.178	0.266	0.376	0.405	0.2
		192	0.343	0.379	0.254	0.318	0.409	0.429	0.3
	Aug	336	0.373	0.400	0.299	0.354	0.435	0.453	0.3
		720	0.444	0.432	0.387	0.408	0.452	0.471	0.4
		avg	0.366	0.392	0.279	0.336	0.418	0.439	0.3
Other pre-training tasks		96	0.303	0.355	0.172	0.261	0.370	0.405	0.2
pre duining uoio		192	0.342	0.376	0.221	0.292	0.403	0.427	0.3
	Noise	336	0.368	0.393	0.272	0.325	0.420	0.439	0.3
		720	0.423	0.422	0.367	0.386	0.442	0.462	0.4
		avg	0.359	0.387	0.258	0.316	0.409	0.433	0.3
		96	0.297	0.345	0.164	0.252	0.379	0.399	0.2
		192	0.334	0.367	0.221	0.290	0.419	0.420	0.3
	TS-Register	336	0.360	0.384	0.275	0.325	0.438	0.442	0.3
		720	0.424	0.416	0.364	0.379	0.435	0.448	0.4
		avg	0.354	0.378	0.256	0.312	0.418	0.427	0.3
		96	0.301	0.348	0.166	0.255	0.380	0.407	0.2
		192	0.343	0.374	0.221	0.291	0.410	0.426	0.3
w/o	Prediction Task	336	0.3/4	0.393	0.275	0.327	0.440	0.443	0.4
		/20	0.424	0.420	0.366	0.384	0.458	0.470	
			0.300	0.384	0.23/	0.314	0.422	0.438	
	1	102	0.329	0.3/1	0.1/5	0.205	0.574	0.399	1 0.2
		336	0.303	0.391	0.233	0.304	0.407	0.422	
	Reconstruction Task		0.461	0.442	0.379	0.396	0.430	0.453	0.
	avg	0.387	0.403	0.269	0.327	0.412	0.428	0.3	
	l	96	0.301	0.357	0.171	0.260	0.419	0.439	
		192	0.358	0.385	0.223	0.294	0.438	0.457	
From	Scratch	336	0.390	0.396	0.282	0.336	0.484	0.484	0.4
FIOI	i Gerateli	720	0.436	0.427	0.366	0.380	0.540	0.538	0.4
			1	1	1	1	1	1	

#### A.10.4 ZERO-SHOT RESULTS

1461         North Nor	1460				Та	Table 15: Full results in zero-shot setting.									
	1461										0				
1463         Metic         MSE         MAE         MSE         MAE         MAE<	1462	Models	3	ROSI	E_512	Tiı	ner	MO	IRAI	Chr	onos	Time	esFM	Mor	nent
1464         96         0.382         0.408         0.414         0.439         0.405         0.407         0.402         0.405         0.432         0.405         0.432         0.405         0.432         0.405         0.432         0.405         0.432         0.405         0.432         0.405         0.435         0.435         0.435         0.530         0.443         0.405         0.432         0.432         0.438         0.705         0.539           1466         1470         0.420         0.442         0.455         0.451         0.500         0.432         0.440         0.705         0.539           1469         0.402         0.442         0.455         0.451         0.443         0.500         0.432         0.449         0.449         0.339         0.334         0.320         0.434         0.336         0.334         0.320         0.431         0.336         0.334         0.322         0.441         0.423         0.411         0.432         0.441         0.435         0.441         0.435         0.441         0.435         0.441         0.443         0.441         0.443         0.441         0.443         0.441         0.445         0.445         0.445         0.445         0.445	1463	Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
1465         147         1480         0.420         0.440         0.455         0.458         0.428         0.561         0.443         0.492         0.438         0.716         0.579           1466         720         0.420         0.425         0.450         0.509         0.451         0.509         0.510         0.512         0.472         0.705         0.577           1469         720         0.401         0.425         0.451         0.445         0.430         0.630         0.438         0.401         0.377         0.416           1470         0.430         0.355         0.303         0.338         0.300         0.338         0.301         0.339         0.330         0.438         0.410         0.397         0.384         0.420         0.440         0.421         0.450         0.440         0.421         0.450         0.430 <td< td=""><td>1464</td><td></td><td>96</td><td>0.382</td><td>0.408</td><td>0.414</td><td>0.439</td><td>0.405</td><td>0.397</td><td>0.494</td><td>0.409</td><td>0.432</td><td>0.405</td><td>0.706</td><td>0.561</td></td<>	1464		96	0.382	0.408	0.414	0.439	0.405	0.397	0.494	0.409	0.432	0.405	0.706	0.561
1466         ETTH1         336         0.404         0.425         0.463         0.509         0.454         0.580         0.460         0.519         0.488         0.705         0.583           1466         120         0.420         0.442         0.445         0.463         0.474         0.605         0.448         0.605         0.448         0.605         0.495         0.512         0.477         0.705         0.597           1469         0         0.80         0.326         0.305         0.335         0.305         0.338         0.305         0.338         0.305         0.338         0.301         0.334         0.301         0.335         0.335         0.335         0.335         0.335         0.337         0.410         0.422         0.430         0.386         0.422           1477         1472         3.36         0.332         0.397         0.340         0.345         0.434         0.346         0.392         0.397         0.341         0.430         0.386         0.422           1477         1476         9         0.512         0.460         0.450         0.441         0.430         0.356         0.331         0.555         0.513         0.471         0.577	1465		192	0.400	0.420	0.440	0.455	0.458	0.428	0.561	0.443	0.492	0.438	0.716	0.579
1467         720         0.420         0.447         0.496         0.492         0.494         0.605         0.495         0.512         0.447         0.705         0.597           1468         96         0.208         0.362         0.305         0.335         0.336         0.500         0.432         0.443         0.560         0.452         0.449         0.444         0.708         0.580           1470         1470         1470         0.336         0.336         0.336         0.338         0.306         0.338         0.311         0.334         0.432         0.416         0.334         0.422         0.439         0.410         0.337         0.416         0.439         0.417         0.436         0.334         0.422         0.439         0.431         0.336         0.334         0.422         0.439         0.413         0.430         0.336         0.334         0.421         0.439         0.431         0.445         0.439         0.431         0.431         0.431         0.445         0.425         0.434           1477         1476         133         0.470         0.570         0.440         0.479         0.555         0.513         0.410         0.609         0.507         0.	1466	ETTh1	336	0.404	0.426	0.455	0.463	0.509	0.454	0.580	0.460	0.519	0.458	0.705	0.583
1468         avg         0.401         0.422         0.451         0.443         0.452         0.452         0.489         0.444         0.708         0.589           1469         192         0.336         0.335         0.302         0.338         0.306         0.338         0.301         0.338         0.311         0.336         0.338         0.336         0.336         0.336         0.336         0.336         0.338         0.306         0.338         0.306         0.338         0.306         0.338         0.336         0.337         0.416         0.337         0.416         0.337         0.416         0.337         0.416         0.337         0.416         0.337         0.416         0.337         0.416         0.337         0.430         0.336         0.425         0.433         0.366         0.430         0.336         0.430         0.336         0.430         0.366         0.430         0.366         0.431         0.460         0.530           1477         1476         120         0.552         0.410         0.570         0.530         0.536         0.536         0.531         0.535         0.513         0.401         0.425         0.437         0.536         0.536         0.536	1467		720	0.420	0.447	0.496	0.496	0.529	0.494	0.605	0.495	0.512	0.477	0.705	0.597
1469         96         0.208         0.362         0.305         0.335         0.308         0.306         0.338         0.311         0.345         0.373         0.416           1470         192         0.336         0.335         0.305         0.306         0.384         0.390         0.410         0.394         0.401         0.396         0.436	1468		avg	0.401	0.425	0.451	0.463	0.475	0.443	0.560	0.452	0.489	0.444	0.708	0.580
1470         192         0.336         0.385         0.365         0.406         0.369         0.384         0.396         0.391         0.401         0.397         0.384         0.422           1471         336         0.339         0.372         0.413         0.377         0.410         0.423         0.417         0.436         0.430         0.386         0.422         0.437           1473         1474         1476         0.396         0.392         0.397         0.396         0.492         0.430         0.450         0.422         0.401         0.423         0.413         0.450         0.423         0.437         0.450         0.432         0.450         0.422         0.400         0.423         0.414         0.423         0.407         0.414         0.442         0.430         0.423         0.413         0.400         0.502         0.433         0.470         0.514         0.443         0.460         0.507         0.551         0.413         0.413         0.401         0.609         0.507         0.551         0.413         0.410         0.469         0.507         0.555           1477         170         0.552         0.470         0.517         0.501         0.501	1469		96	0.298	0.362	0.305	0.355	0.303	0.338	0.306	0.338	0.311	0.345	0.373	0.416
1471         ETTH2         336         0.335         0.399         0.378         0.413         0.397         0.410         0.423         0.417         0.436         0.430         0.386         0.425           1472         720         0.395         0.432         0.414         0.457         0.447         0.450         0.432         0.437         0.450         0.425         0.437         0.450         0.422         0.444           1473         wg         0.346         0.340         0.422         0.660         0.476         0.514         0.437         0.413         0.400         0.422         0.460         0.413         0.415         0.413         0.400         0.422         0.460         0.414         0.415         0.413         0.401         0.600         0.550           1476         136         0.522         0.470         0.570         0.490         0.501         0.501         0.445         0.419         0.601         0.551           1477         1476         wg         0.525         0.471         0.544         0.476         0.714         0.576         0.327         0.293         0.330         0.555         0.513         0.430         0.435         0.431         0.330	1470		192	0.336	0.385	0.365	0.406	0.369	0.384	0.396	0.394	0.401	0.397	0.384	0.422
1472       1720       0.395       0.432       0.414       0.457       0.447       0.450       0.442       0.439       0.437       0.450       0.432       0.434         1473       avg       0.346       0.394       0.366       0.408       0.379       0.396       0.392       0.397       0.396       0.437       0.366       0.374       0.650       0.514       0.443       0.366       0.374       0.650       0.514       0.443       0.461       0.401       0.650       0.557       0.548       0.707       0.6400       0.500       0.600       0.570       0.444       0.476       0.514       0.443       0.461       0.476       0.516       0.600       0.507       0.444       0.470       0.570       0.440       0.570       0.443       0.470       0.570       0.444       0.470       0.561       0.560       0.573       0.513       0.470       0.710       0.557         1477       1476       avg       0.522       0.471       0.544       0.470       0.530       0.535       0.513       0.407       0.703       0.555       0.513       0.407       0.425       0.434       0.425       0.434       0.425       0.434       0.425       0.434	1471	ETTh2	336	0.353	0.399	0.378	0.413	0.397	0.410	0.423	0.417	0.436	0.430	0.386	0.426
intra         intr<         intr<         intr<         intr<         intr<         intr         intr<         intr< <thi< td=""><td>1472</td><td></td><td>720</td><td>0.395</td><td>0.432</td><td>0.414</td><td>0.457</td><td>0.447</td><td>0.450</td><td>0.442</td><td>0.439</td><td>0.437</td><td>0.450</td><td>0.425</td><td>0.454</td></thi<>	1472		720	0.395	0.432	0.414	0.457	0.447	0.450	0.442	0.439	0.437	0.450	0.425	0.454
1173         96         0.512         0.460         0.440         0.422         0.660         0.476         0.514         0.443         0.366         0.374         0.679         0.559           1475         1476         192         0.512         0.462         0.505         0.458         0.707         0.500         0.608         0.475         0.413         0.401         0.690         0.550           1476         1720         0.552         0.470         0.570         0.490         0.730         0.515         0.607         0.445         0.429         0.701         0.557           1477         1478         96         0.224         0.309         0.285         0.216         0.282         0.202         0.293         0.189         0.257         0.230         0.388           1480         1420         0.266         0.333         0.265         0.312         0.308         0.301         0.388         0.313         0.355         0.464         0.448         0.423         0.424           1481         1482         1482         1482         0.330         0.328         0.331         0.335         0.331         0.332         0.331         0.332         0.444         0.423	1473		avg	0.346	0.394	0.366	0.408	0.379	<u>0.396</u>	0.392	0.397	0.396	0.405	0.392	0.430
1475         192         0.512         0.462         0.505         0.458         0.707         0.500         0.608         0.475         0.413         0.401         0.690         0.550           1476         336         0.523         0.470         0.570         0.490         0.730         0.515         0.690         0.507         0.445         0.429         0.701         0.557           1477         1478         0.490         0.524         0.783         0.556         0.513         0.445         0.401         0.690         0.555           1479         0.522         0.471         0.544         0.476         0.714         0.536         0.495         0.434         0.419         0.697         0.555           1480         96         0.224         0.300         0.285         0.216         0.233         0.265         0.331         0.361         0.368         0.373         0.555         0.561         0.339         0.369         0.339         0.361         0.336         0.320         0.335         0.381         0.339         0.361         0.336         0.320         0.335         0.331         0.361         0.332         0.346         0.331         0.361         0.332         0.345	1470		96	0.512	0.460	0.440	0.422	0.660	0.476	0.514	0.443	0.366	0.374	0.679	0.544
ETTm1         336         0.523         0.470         0.570         0.490         0.730         0.515         0.690         0.507         0.445         0.429         0.701         0.557           1476         720         0.552         0.490         0.659         0.534         0.758         0.536         0.733         0.555         0.513         0.470         0.719         0.569           1476         avg         0.525         0.471         0.544         0.476         0.714         0.507         0.636         0.495         0.434         0.419         0.697         0.555           1479         avg         0.525         0.471         0.544         0.476         0.714         0.507         0.636         0.495         0.449         0.495         0.445         0.429         0.308         0.308           1480         mag         0.226         0.333         0.265         0.327         0.294         0.330         0.265         0.333         0.360         0.338         0.319         0.360           1481         mag         0.206         0.335         0.341         0.356         0.313         0.360         0.332         0.340         0.322         0.313         0.332	1/175		192	0.512	0.462	0.505	0.458	0.707	0.500	0.608	0.475	0.413	0.401	0.690	0.550
1477       720       0.552       0.490       0.659       0.534       0.758       0.536       0.733       0.555       0.513       0.470       0.719       0.569         1477       avg       0.525       0.471       0.544       0.476       0.714       0.507       0.636       0.495       0.434       0.419       0.697       0.555         1479       avg       0.525       0.471       0.545       0.224       0.309       0.203       0.285       0.216       0.282       0.202       0.293       0.189       0.257       0.230       0.308         1480       192       0.266       0.333       0.265       0.327       0.294       0.330       0.286       0.348       0.277       0.325       0.281       0.338         1481       1482       avg       0.299       0.352       0.407       0.405       0.410       0.494       0.439       0.409       0.425       0.444       0.423       0.424         1483       avg       0.260       0.190       0.236       0.188       0.250       0.209       0.244       -       -       0.216       0.216       0.217         1484       192       0.230       0.235 <td< td=""><td>1475</td><td>ETTm1</td><td>336</td><td>0.523</td><td>0.470</td><td>0.570</td><td>0.490</td><td>0.730</td><td>0.515</td><td>0.690</td><td>0.507</td><td>0.445</td><td>0.429</td><td>0.701</td><td>0.557</td></td<>	1475	ETTm1	336	0.523	0.470	0.570	0.490	0.730	0.515	0.690	0.507	0.445	0.429	0.701	0.557
iiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiiii	1470		720	0.552	0.490	0.659	0.534	0.758	0.536	0.733	0.555	0.513	0.470	0.719	0.569
1478       96       0.224       0.309       0.203       0.285       0.216       0.282       0.202       0.293       0.189       0.257       0.230       0.308         1430       192       0.266       0.333       0.265       0.327       0.294       0.330       0.286       0.348       0.277       0.325       0.285       0.338         1480       336       0.310       0.358       0.319       0.361       0.368       0.373       0.355       0.386       0.350       0.381       0.399       0.369         1481       482       avg       0.299       0.352       0.400       0.499       0.429       0.425       0.446       0.448       0.423       0.424         482       avg       0.299       0.352       0.360       0.348       0.356       0.313       0.363       0.320       0.353       0.319       0.360         1483       484       424       avg       0.299       0.288       0.261       0.293       0.223       0.209       0.244       -       -       0.216       0.271         1484       486       423       0.225       0.323       0.311       0.322       0.284       0.322       0.331 <td>1477</td> <td></td> <td>avg</td> <td>0.525</td> <td>0.471</td> <td>0.544</td> <td>0.476</td> <td>0.714</td> <td>0.507</td> <td>0.636</td> <td>0.495</td> <td>0.434</td> <td>0.419</td> <td>0.697</td> <td>0.555</td>	1477		avg	0.525	0.471	0.544	0.476	0.714	0.507	0.636	0.495	0.434	0.419	0.697	0.555
1479       192       0.266       0.333       0.265       0.327       0.294       0.330       0.286       0.348       0.277       0.325       0.285       0.338         1480       336       0.310       0.358       0.319       0.361       0.368       0.373       0.355       0.386       0.350       0.381       0.339       0.369         1481       1482       avg       0.299       0.352       0.407       0.405       0.410       0.494       0.439       0.409       0.425       0.464       0.448       0.423       0.424         1483       avg       0.299       0.352       0.360       0.386       0.343       0.356       0.313       0.363       0.320       0.353       0.319       0.360         1483       1484       192       0.239       0.286       0.188       0.250       0.209       0.244       -       -       0.216       0.271         1486       1486       336       0.279       0.315       0.332       0.323       0.301       0.332       -       -       0.264       0.305         1486       1486       avg       0.265       0.305       0.292       0.312       0.265       0.302	1478		96	0.224	0.309	0.203	0.285	0.216	0.282	0.202	0.293	0.189	0.257	0.230	0.308
1480       ETTm2       336       0.310       0.358       0.319       0.361       0.368       0.373       0.355       0.386       0.350       0.381       0.39       0.369         1481       1482       720       0.395       0.407       0.405       0.410       0.494       0.439       0.409       0.425       0.464       0.448       0.423       0.464         1482       avg       0.299       0.352       0.360       0.386       0.313       0.363       0.320       0.353       0.319       0.360         1483       avg       0.299       0.352       0.360       0.188       0.250       0.209       0.244       -       -       0.216       0.271         1484       192       0.239       0.288       0.261       0.293       0.287       0.284       0.288       -       -       0.264       0.306         1485       336       0.279       0.315       0.332       0.340       0.288       0.371       0.288       0.374       -       -       0.369       0.330       0.380         1486       avg       0.265       0.305       0.292       0.312       0.212       0.300       0.288       0.310	1479		192	0.266	0.333	0.265	0.327	0.294	0.330	0.286	0.348	0.277	0.325	0.285	0.338
1481       720       0.395       0.407       0.405       0.410       0.494       0.409       0.425       0.464       0.448       0.423       0.424         1482       avg       0.299       0.352       0.360       0.386       0.343       0.356       0.313       0.363       0.320       0.353       0.319       0.360         1483       96       0.200       0.260       0.190       0.236       0.188       0.250       0.209       0.244       -       -       0.216       0.271         1484       192       0.239       0.288       0.237       0.284       0.254       0.288       -       -       0.0264       0.306         1485       336       0.279       0.315       0.332       0.340       0.282       0.323       0.301       0.332       -       -       0.0264       0.306         1486       avg       0.265       0.305       0.292       0.312       0.267       0.300       0.288       0.310       -       -       0.0369       0.332         1488       avg       0.265       0.305       0.292       0.312       0.212       0.301       0.194       0.266       -       -       0.844 </td <td>1480</td> <td>ETTm2</td> <td>336</td> <td>0.310</td> <td>0.358</td> <td>0.319</td> <td>0.361</td> <td>0.368</td> <td>0.373</td> <td>0.355</td> <td>0.386</td> <td>0.350</td> <td>0.381</td> <td>0.339</td> <td>0.369</td>	1480	ETTm2	336	0.310	0.358	0.319	0.361	0.368	0.373	0.355	0.386	0.350	0.381	0.339	0.369
1482       avg       0.299       0.352       0.360       0.386       0.343       0.356       0.313       0.363       0.320       0.353       0.319       0.360         1483       96       0.200       0.260       0.190       0.236       0.188       0.250       0.209       0.244       -       -       0.216       0.271         1484       192       0.239       0.288       0.261       0.293       0.237       0.284       0.254       0.288       -       -       0.264       0.306         1486       336       0.279       0.315       0.332       0.340       0.282       0.323       0.301       0.332       -       -       0.369       0.386         1487       336       0.265       0.305       0.292       0.312       0.267       0.300       0.288       0.310       -       -       0.369       0.389         1488       avg       0.265       0.305       0.292       0.312       0.212       0.301       0.194       0.266       -       -       0.844       0.761         1489       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.329	1481		720	0.395	0.407	0.405	0.410	0.494	0.439	0.409	0.425	0.464	0.448	0.423	0.424
1483       96       0.200       0.260       0.190       0.236       0.188       0.209       0.244       -       -       0.216       0.271         1484       192       0.239       0.288       0.261       0.293       0.237       0.284       0.254       0.288       -       -       0.264       0.306         1485       336       0.279       0.315       0.332       0.340       0.282       0.323       0.301       0.332       -       -       0.313       0.336         1486	1482		avg	0.299	0.352	0.360	0.386	0.343	0.356	<u>0.313</u>	0.363	0.320	0.353	0.319	0.360
1484       192       0.239       0.288       0.261       0.293       0.237       0.284       0.254       0.288       -       -       0.264       0.306         1485       1486       336       0.279       0.315       0.332       0.340       0.282       0.323       0.301       0.332       -       -       0.313       0.336         1486       1487       avg       0.265       0.305       0.292       0.312       0.267       0.300       0.288       0.310       -       -       0.369       0.380         1488       avg       0.265       0.305       0.292       0.312       0.267       0.300       0.288       0.310       -       -       0.369       0.323         1489       96       0.209       0.307       0.210       0.312       0.212       0.301       0.194       0.266       -       -       0.844       0.761         1490       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.289       -       -       0.862       0.766         1490       192       0.219       0.321       0.245       0.331       0.244       0.321 <td< td=""><td>1483</td><td></td><td>96</td><td>0.200</td><td>0.260</td><td>0.190</td><td>0.236</td><td><u>0.188</u></td><td>0.250</td><td>0.209</td><td>0.244</td><td>-  </td><td>  -</td><td>0.216</td><td>0.271</td></td<>	1483		96	0.200	0.260	0.190	0.236	<u>0.188</u>	0.250	0.209	0.244	-	-	0.216	0.271
1485       Weather       336       0.279       0.315       0.332       0.340       0.282       0.323       0.301       0.332       -       -       0.313       0.336         1486       720       0.340       0.357       0.385       0.381       0.359       0.345       0.388       0.374       -       -       0.369       0.380         1487       avg       0.265       0.305       0.292       0.312       0.267       0.300       0.288       0.310       -       -       0.291       0.323         1488       avg       0.265       0.307       0.210       0.312       0.212       0.301       0.194       0.266       -       -       0.844       0.761         1489       96       0.209       0.307       0.210       0.312       0.212       0.301       0.194       0.266       -       -       0.844       0.761         1490       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.321       -       -       0.862       0.766         1491       366       0.236       0.479       0.282       0.358       0.312       0.312       - <td< td=""><td>1484</td><td></td><td>192</td><td><u>0.239</u></td><td><u>0.288</u></td><td>0.261</td><td>0.293</td><td>0.237</td><td>0.284</td><td>0.254</td><td>0.288</td><td>-</td><td>  -</td><td>0.264</td><td>0.306</td></td<>	1484		192	<u>0.239</u>	<u>0.288</u>	0.261	0.293	0.237	0.284	0.254	0.288	-	-	0.264	0.306
1486       720       0.340       0.357       0.385       0.381       0.359       0.345       0.388       0.374       -       -       0.369       0.380         1487       avg       0.265       0.305       0.292       0.312       0.267       0.300       0.288       0.310       -       -       0.291       0.323         1488       96       0.209       0.307       0.210       0.312       0.212       0.301       0.194       0.266       -       -       0.844       0.761         1489       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.289       -       -       0.844       0.761         1490       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.289       -       -       0.862       0.766         1491       336       0.236       0.330       0.244       0.321       -       -       0.888       0.774         1492       0.273       0.328       0.456       0.479       0.282       0.358       0.312       -       -       0.861       0.766         1493       avg <td>1485</td> <td>Weather</td> <td>336</td> <td>0.279</td> <td>0.315</td> <td>0.332</td> <td>0.340</td> <td>0.282</td> <td><u>0.323</u></td> <td>0.301</td> <td>0.332</td> <td>-  </td> <td>  -</td> <td>0.313</td> <td>0.336</td>	1485	Weather	336	0.279	0.315	0.332	0.340	0.282	<u>0.323</u>	0.301	0.332	-	-	0.313	0.336
1487       avg       0.265       0.305       0.292       0.312       0.267       0.300       0.288       0.310       -       -       0.291       0.323         1488       96       0.209       0.307       0.210       0.312       0.212       0.301       0.194       0.266       -       -       0.844       0.761         1489       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.289       -       -       0.844       0.762         1490       336       0.236       0.330       0.284       0.372       0.245       0.333       0.244       0.321       -       -       0.862       0.762         1491       336       0.236       0.330       0.284       0.372       0.245       0.333       0.244       0.321       -       -       0.862       0.766         1492       0.273       0.328       0.456       0.479       0.282       0.358       0.312       -       -       0.861       0.766         1493       avg       0.234       0.320       0.297       0.375       0.241       0.328       0.245       0.312       -       -       1.8	1486		720	0.340	0.357	0.385	0.381	<u>0.359</u>	0.345	0.388	0.374	-	-	0.369	0.380
1488       96       0.209       0.307       0.210       0.312       0.212       0.301       0.194       0.266       -       -       0.844       0.761         1489       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.289       -       -       0.850       0.762         1490       192       0.236       0.330       0.284       0.372       0.245       0.333       0.244       0.321       -       -       0.862       0.766         1491       1492       120       0.234       0.320       0.297       0.375       0.241       0.328       0.321       -       -       0.862       0.766         1492       1492       120       0.234       0.320       0.297       0.375       0.241       0.328       0.324       0.312       -       -       0.863       0.774         1493       1494       192       0.575       0.406       0.561       0.385       -       -       0.579       0.412       -       -       1.403       0.802         1493       192       0.575       0.406       0.561       0.385       -       -       0.579	1487		avg	0.265	0.305	0.292	0.312	0.267	0.300	0.288	0.310	-	-	0.291	0.323
1489       192       0.219       0.315       0.239       0.337       0.225       0.320       0.218       0.289       -       -       0.850       0.762         1490       336       0.236       0.330       0.284       0.372       0.245       0.333       0.244       0.321       -       -       0.862       0.766         1491       720       0.273       0.328       0.456       0.479       0.282       0.358       0.324       0.371       -       -       0.862       0.766         1492       avg       0.234       0.320       0.297       0.375       0.241       0.328       0.324       0.371       -       -       0.868       0.774         1492       avg       0.234       0.320       0.297       0.375       0.241       0.328       0.324       0.312       -       -       0.861       0.766         1493       96       0.572       0.407       0.526       0.368       -       -       0.562       0.378       -       -       1.390       0.800         1494       192       0.575       0.406       0.561       0.385       -       -       0.594       0.420       -	1488		96	0.209	0.307	0.210	0.312	0.212	<u>0.301</u>	0.194	0.266	-	-	0.844	0.761
1490       Electricity       336       0.236       0.330       0.284       0.372       0.245       0.333       0.244       0.321       -       -       0.862       0.766         1491       720       0.273       0.328       0.456       0.479       0.282       0.358       0.324       0.371       -       -       0.862       0.766         1492       avg       0.234       0.320       0.297       0.375       0.241       0.328       0.324       0.312       -       -       0.868       0.774         1493       avg       0.234       0.320       0.297       0.375       0.241       0.328       0.245       0.312       -       -       0.861       0.766         1493       96       0.572       0.407       0.526       0.368       -       -       0.562       0.378       -       -       1.390       0.800         1494       192       0.575       0.406       0.561       0.385       -       -       0.579       0.412       -       -       1.403       0.802         1495       336       0.588       0.411       0.614       0.412       -       -       0.594       0.420	1489		192	<u>0.219</u>	0.315	0.239	0.337	0.225	0.320	0.218	0.289	-	-	0.850	0.762
1491       720       0.273       0.328       0.456       0.479       0.282       0.358       0.324       0.371       -       -       0.888       0.774         1492       avg       0.234       0.320       0.297       0.375       0.241       0.328       0.245       0.312       -       -       0.888       0.774         1493       96       0.572       0.407       0.526       0.368       -       -       0.562       0.378       -       -       1.390       0.800         1493       192       0.575       0.406       0.561       0.385       -       -       0.579       0.412       -       -       1.403       0.802         1494       336       0.588       0.411       0.614       0.412       -       -       0.594       0.420       -       1.403       0.802         1496       336       0.588       0.411       0.614       0.412       -       -       0.594       0.420       -       1.415       0.808         1496       .928       0.588       0.411       0.614       0.412       -       -       0.594       0.420       -       1.437       0.808       .492	1490	Electricity	336	0.236	0.330	0.284	0.372	0.245	0.333	0.244	0.321	-	-	0.862	0.766
1492       avg       0.234       0.320       0.297       0.375       0.241       0.328       0.245       0.312       -       -       0.861       0.766         1493       96       0.572       0.407       0.526       0.368       -       -       0.562       0.378       -       -       1.390       0.800         1494       192       0.575       0.406       0.561       0.385       -       -       0.579       0.412       -       -       1.403       0.802         1495       336       0.588       0.411       0.614       0.412       -       -       0.579       0.412       -       -       1.403       0.802         1496       336       0.588       0.411       0.614       0.412       -       -       0.594       0.420       -       1.415       0.804         1496       720       0.618       0.422       0.749       0.464       -       -       0.723       0.472       -       -       1.437       0.808         1497       avg       0.588       0.412       0.613       0.407       -       -       0.615       0.421       -       -       1.411       0.804 </td <td>1491</td> <td></td> <td>720</td> <td>0.273</td> <td>0.328</td> <td>0.456</td> <td>0.479</td> <td>0.282</td> <td><u>0.358</u></td> <td>0.324</td> <td>0.371</td> <td>-</td> <td>-  </td> <td>0.888</td> <td>0.774</td>	1491		720	0.273	0.328	0.456	0.479	0.282	<u>0.358</u>	0.324	0.371	-	-	0.888	0.774
1493       96       0.572       0.407       0.526       0.368       -       -       0.562       0.378       -       -       1.390       0.800         1494       1495       192       0.575       0.406       0.561       0.385       -       -       0.579       0.412       -       -       1.403       0.802         1496       336       0.588       0.411       0.614       0.412       -       -       0.594       0.420       -       -       1.415       0.804         1496       720       0.618       0.422       0.749       0.464       -       -       0.723       0.472       -       -       1.437       0.808         1497       avg       0.588       0.412       0.613       0.407       -       -       0.615       0.421       -       -       1.411       0.804	1492		avg	0.234	0.320	0.297	0.375	0.241	0.328	0.245	0.312	-	-	0.861	0.766
1494       192       0.575       0.406       0.561       0.385       -       -       0.579       0.412       -       -       1.403       0.802         1495       336       0.588       0.411       0.614       0.412       -       -       0.594       0.420       -       -       1.415       0.804         1496       720       0.618       0.422       0.749       0.464       -       -       0.723       0.472       -       -       1.437       0.808         1497       avg       0.588       0.412       0.613       0.407       -       -       0.615       0.421       -       -       1.411       0.804	1493		96	0.572	0.407	0.526	0.368	-	-	<u>0.562</u>	<u>0.378</u>	-	-	1.390	0.800
1495       Traffic $336$ $0.588$ $0.411$ $0.614$ $0.412$ $  0.594$ $0.420$ $  1.415$ $0.804$ 1496       720 $0.618$ $0.422$ $0.749$ $0.464$ $  0.723$ $0.472$ $ 1.437$ $0.808$ 1497       avg $0.588$ $0.412$ $0.613$ $0.407$ $  0.615$ $0.421$ $ 1.411$ $0.804$	1494		192	<u>0.575</u>	0.406	0.561	0.385	-	-	0.579	0.412	-	-	1.403	0.802
1496         720         0.618         0.422         0.749         0.464         -         -         0.723         0.472         -         -         1.437         0.808           1497         avg         0.588         0.412         0.613         0.407         -         -         0.615         0.421         -         -         1.411         0.804	1495	Traffic	336	0.588	0.411	0.614	0.412	-	-	<u>0.594</u>	0.420	-	-	1.415	0.804
1497 avg 0.588 0.412 0.613 0.407 - 0.615 0.421 - 1.1411 0.804	1496		720	0.618	0.422	0.749	0.464	-	-	0.723	0.472	-	-	1.437	0.808
	1497		avg	0.588	<u>0.412</u>	0.613	0.407	-	-	<u>0.615</u>	0.421	-	-	1.411	0.804

## 1512 A.11 THE DETAILED SETTING OF T-SNE VISUALIZATION

In Figure1(b), we select three datasets (Pems08, PSRA, Electricity) from the transport, nature, and
energy domains respectively, and compared the differences in hidden representations between direct
transfer and adaptive transfer. Specifically, direct transfer refers to the case where domain specific
information is not considered, while adaptive transfer considers domain specific information that is
learned by register tokens. We visualize the output of the encoder's hidden representations using
t-SNE.

1520

1521 A.12 ADDITIONAL RESULTS

#### 1523 A.12.1 THE RESULTS WITH LOOK-BACK WINDOW L = 961524

To demonstrate ROSE excels not only with fixed inputs, we evaluate it using input lengths significantly
shorter than the 512 employed during pre-training. As shown in Table 16, ROSE still achieves the
state-of-the-art performance, indicating the effective transfer of pre-trained knowledge with a shorter
look-back window.

Table 16: The results for ROSE and other baselines in full-shot setting with look-back window of 96.
 The average results of all predicted lengths are listed here.

-							
Models	ROSE	ITransformer	PatchTST	Timesnet	Dlinear	GPT4TS	S <sup>2</sup> IP-LLM
Metric	MSE MA	E   MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE
ETTm1	0.389 0.38	<b>89</b>   0.407   0.410	0.387 0.400	0.400   0.406	0.403 0.407	0.389 0.397	0.390   0.399
ETTm2	0.272 0.32	21   0.288   0.332	<u>0.281</u> <u>0.326</u>	0.291   0.333	0.350 0.401	0.285   0.331	0.278 0.327
ETTh1	0.432 0.42	26   0.454   0.447	0.469 0.454	0.458   0.450	0.456 0.452	0.447   0.436	<u>0.444</u> <u>0.431</u>
ETTh2	0.376 0.39	<b>3</b>   0.383   0.407	0.387 0.407	0.414   0.427	0.559 0.515	0.381   0.408	<u>0.378</u> <u>0.402</u>
Weather	0.257 0.27	<b>76</b>   <u>0.258</u>   <u>0.278</u>	0.259 0.281	0.259   0.287	0.265 0.317	0.264   0.284	0.266   0.284
Electricity	0.176 0.26	<b>58</b>   <u>0.178</u>   <u>0.270</u>	0.205 0.290	0.192   0.296	0.354 0.414	0.205   0.290	0.195   0.285
Traffic	0.440 0.27	<b>76</b>   <b>0.428</b>   <u>0.282</u>	0.481 0.304	0.620 0.336	0.625 0.383	0.488 0.317	0.467 0.305

### A.12.2 THE SHORT-TERM FORECASTING IN ZERO-SHOT SETTING.

We evaluate ROSE for short-term prediction in zero-shot setting following Moment (Goswami et al., 2024). For extremely short input sequences in the M4 dataset, we adapt them by applying padding.
As shown in Table 17, ROSE demonstrates competitive performance compared to other baselines.

Table 17: The results on short-term forecasting in zero-shot setting, measured using SMAPE.

Models	ROSE	Moment	GPT4TS   '	TimesNet   N-BEATS
M4 Yearly	14.08	14.84	14.80	14.40   <u>14.18</u>
M4 Quarterly	<u>11.79</u>	12.02	11.77	13.21   12.25
M4 Monthly	15.33	15.80	15.36	15.67   15.24

1555 1556

1544

1548 1549

1550 1551

1552 1553 1554

1557

A.12.3 THE COMPARISON WITH THE TIME SERIES REPRESENTATION METHOD

To demonstrate the effectiveness of the pre-trained model in comparison to existing time series representation learning methods, we select recent time series representation learning methods, including SimMTM (Dong et al., 2024), TS2Vec (Yue et al., 2022) and TF-C (Zhang et al., 2022). Table 18 illustrates the performance of ROSE in full-shot and 10% few-shot settings, in comparison with that of time series representation learning methods in full-shot settings. It can be observed that ROSE outperforms the representation learning methods in full-shot setting, and achieves a competitive performance even in 10% few-shot setting, which substantiates the effectiveness of ROSE as a pre-trained model.

Models   RO	SE   ROSE	E (10%)   Sim	MTM	Time2Vec	TF-C
Metric   MSE	MAE   MSE	MAE   MSE	MAE	MSE   MAE	MSE   MAE
ETTm1   0.341	<b>0.367</b> 0.349	0.372 0.341	0.377	0.691   0.547	0.732   0.652
ETTm2   0.246	<b>0.305</b> <u>0.250</u>	0.308 0.258	0.315	0.316   0.351	1.721   0.922
ETTh1   0.391	<b>0.414</b>   <u>0.397</u>	<u>0.419</u>   0.401	0.423	0.426   0.436	0.614   0.601
ETTh2   0.331	<b>0.374</b>   <u>0.335</u>	<u>0.380</u>   0.342	0.384	0.423   0.459	0.387   0.374
Weather   0.217	<b>0.251</b>   <u>0.224</u>	<u>0.252</u> <u>0.224</u>	0.262	0.231   0.264	0.286   0.349
Electricity   0.155	<b>0.248</b>   0.164	0.253 0.161	0.254	0.203   0.283	0.355   0.389
Traffic   0.390	<b>0.264</b> 0.418	0.278   0.393	0.268	0.450   0.330	0.702   0.443

Table 18: The results for ROSE in full-shot setting and 10% few-shot setting, compared with other time series representation learning methods in full-shot setting. The average results of all predicted lengths are listed here.

#### ADDITIONAL MODEL ANALYSIS A.13

#### A.13.1 THE ANALYSIS OF MULTIPLE PREDICTION HEADS

**Pre-train**: We compared the full-shot performance of using four prediction heads versus a single prediction head in calculating the loss across four different prediction lengths during pre-training. As shown in the Table 19, using four prediction heads consistently yields better performance across all lengths. This phenomenon may be because training four prediction heads in pre-training allows the model to focus on forecasts of different lengths, enhancing accuracy in specific ranges. Conversely, a single prediction head must handle both short-term and long-term forecasts, requiring the model to balance accuracy across various lengths, which could limit its performance on multiple prediction lengths.

Table 19: The results in full-shot setting of pre-training with four prediction heads or a single prediction head in calculating the loss across four different prediction lengths.

Models		ET	Th1	ET	Th2	ET	Гm1	ET	Гm2
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
	96	0.354	0.385	0.265	0.320	0.275	0.328	0.157	0.243
	192	0.389	0.407	0.328	0.369	0.324	0.358	0.213	0.283
pre-train w four heads	336	0.406	0.422	0.353	0.391	0.354	0.377	0.266	0.319
I	720	0.413	0.443	0.376	0.417	0.411	0.407	0.347	0.373
	avg	0.391	0.414	0.331	0.374	0.341	0.367	0.246	0.305
	96	0.357	0.388	0.269	0.325	0.277	0.330	0.158	0.245
	192	0.394	0.411	0.333	0.372	0.324	0.356	0.215	0.287
pre-train w one head	336	0.415	0.422	0.359	0.395	0.360	0.380	0.272	0.325
1	720	0.420	0.444	0.390	0.425	0.421	0.410	0.355	0.385
	avg	0.397	0.416	0.338	0.379	0.346	0.369	0.250	0.310

Inference: During inference, when the prediction length is covered by multiple heads, we select a prediction head whose output length is the closest to the prediction length. For example, if the prediction length is 48, we select only the prediction head whose output length is 96, even though the other three heads could also perform the prediction length of 48 by cropping. 

To demonstrate the effectiveness of this practice, we evaluate the prediction performance for different prediction lengths that are covered by multiple heads, whose results are shown in Table 20 and Table 21. As an example, the first column indicates that using prediction heads of 96, 192, 336, and 720 respectively, to predict with the length of 48. It is evident that the strategy of choosing the closest prediction head allows the model to adapt to various prediction lengths and achieve SOTA performance.

1622						
1623		48	96	192	336	720
1624	prediction heads\ prediction lengths	MSE   MAE	MSE   MAE	MSE   MAE	MSE   MAE	MSE   MAE
1625	head of 96	0.325 0.364	0.354 0.385			
1626	head of 192	0.327   0.365	0.358 0.385	0.389 0.407	-   -	
1627	head of 336	0.327   0.365	0.361 0.386	0.399 0.410	0.406   0.422	-   -
1628	head of 720	0.329 0.365	0.359 0.389	0.401 0.411	0.429 0.425	0.413 0.443

# Table 20: The prediction performance of ROSE in ETTh1 for different prediction lengths that are covered by multiple heads.

Table 21: The prediction performance of ROSE in ETTm2 for different prediction lengths that are covered by multiple heads.

	48	96	192	336	720
prediction heads\ prediction lengths	MSE   MAE				
head of 96	0.120   0.213	0.157   0.243	-   -	-   -	-   -
head of 192	0.120   0.213	0.159   0.245	0.213   0.281	-   -	-   -
head of 336	0.122   0.214	0.159   0.244	0.216   0.283	0.266   0.319	-   -
head of 720	0.124   0.216	0.159   0.245	0.216   0.284	0.269 0.321	0.347   0.373

# 1640 A.13.2 MULTI-LOSSES BALANCING

1642 To validate ROSE's robustness of the multi-losses, we introduce a hyper-parameter  $\lambda \in (0, 1)$ , and 1643 define  $\mathcal{L}_{pretrain} = \lambda \mathcal{L}_{reconstruction} + (1 - \lambda) \mathcal{L}_{prediction} + \mathcal{L}_{register}$ . Since the register loss only 1644 constraints the parameter updates of register, its gradient does not influence to the backbone of the 1645 model. Therefore, the register loss does not cause an imbalance in the training of the model. We 1646 vary  $\lambda$ 's value and report the results in full-shot setting. As shown in Table 22, ROSE is not sensitive 1647 to changes of  $\lambda$ , thus balancing the loss of the model is not challenging. Therefore, our final loss 

Table 22: The results in full-shot setting of ROSE pre-trained with different  $\lambda$ . The average results of all predicted lengths are listed here.

$\lambda$	0.2	0.4	0.6	0.8	Standard Deviation
Metric	MSE   MAE	MSE   MAE	MSE   MAE	MSE   MAE	MSE MAE
ETTh1	0.3973   0.4199	0.3978   0.4193	0.3978 0.4205	0.3996 0.423	0.0008 0.0014
ETTh2	0.3339   0.379	0.3347   0.3802	0.3369 0.3822	0.3351 0.383	0.0011 0.0016
ETTm1	0.3500   0.3733	0.3512   0.3747	0.3492   0.3717	0.3479 0.3719	0.0012 0.0012
ETTm2	0.2538   0.3111	0.2534   0.3095	0.2512 0.3092	0.2505 0.3092	0.0014 0.0007

# 1674 A.13.3 THE ANALYSIS OF DECOMPOSED FREQUENCY LEARNING

To demonstrate the effectiveness of decomposed frequency learning in capturing unified representations, we pre-trained the model using multi-frequency masking, patch masking and multi-patch masking. We visualize the reconstruction performance of the three methods in out-of-distribution (OOD) scenarios. As shown in Table 16, the model pre-trained with multiple frequency masking exhibits greater robustness to complex temporal patterns, confirming that decomposed frequency learning can assist in learning unified representations.



Figure 16: The visualization of the reconstruction performance of models pre-trained with three methods (multi-frequency masking, patch masking and multi-patch masking) respectively in out-of-distribution (OOD) scenarios.