

---

# Abstain Mask Retain Core: Time Series Prediction by Adaptive Masking Loss with Representation Consistency

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Time series forecasting plays a pivotal role in critical domains such as energy management  
2 and financial markets. Although deep learning-based approaches (e.g.,  
3 MLP, RNN, Transformer) have achieved remarkable progress, the prevailing "long-  
4 sequence information gain hypothesis" exhibits inherent limitations. Through  
5 systematic experimentation, this study reveals a counterintuitive phenomenon:  
6 appropriately truncating historical data can paradoxically enhance prediction ac-  
7 curacy, indicating that existing models learn substantial redundant features (e.g.,  
8 noise or irrelevant fluctuations) during training, thereby compromising effective  
9 signal extraction. Building upon information bottleneck theory, we propose an  
10 innovative solution termed Adaptive Masking Loss with Representation Consis-  
11 tency (AMRC), which features two core components: 1) Dynamic masking loss,  
12 which adaptively identified highly discriminative temporal segments to guide gra-  
13 dient descent during model training; 2) Representation consistency constraint,  
14 which stabilized the mapping relationships among inputs, labels, and predictions.  
15 Experimental results demonstrate that AMRC effectively suppresses redundant  
16 feature learning while significantly improving model performance. This work  
17 not only challenges conventional assumptions in temporal modeling but also pro-  
18 vides novel theoretical insights and methodological breakthroughs for developing  
19 efficient and robust forecasting models. We have made our code available at  
20 <https://anonymous.4open.science/r/AMRC/>.

21 

## 1 Introduction

22 Time series forecasting, as a pivotal technology in critical domains such as energy management  
23 and financial markets, directly influences decision-making quality and economic efficiency [11,  
24 18, 12, 19, 22]. Recent breakthroughs in deep learning have driven revolutionary advancements in  
25 time series prediction. Contemporary frameworks including Multilayer Perceptron (MLP)-based  
26 architectures [17, 30, 7, 27, 4, 28], Recurrent Neural Networks (RNNs) with their variants [13, 21, 9],  
27 and attention mechanism-based models exemplified by the Transformer [20, 33, 32, 16, 34, 2, 6], have  
28 achieved remarkable breakthroughs in modeling complex temporal patterns through the construction  
29 of elaborate hierarchical temporal dependencies.

30 Current mainstream forecasting models predominantly adhere to the "long-sequence information gain  
31 hypothesis," which posits that extending historical data length enhances the availability of temporal  
32 dependencies [31, 15]. However, through systematic experimental analysis, this study challenges  
33 this conventional assumption. As shown in Table 1, we observed a counterintuitive phenomenon  
34 across multiple benchmark datasets and diverse model architectures: appropriately truncating early  
35 segments of input sequences can significantly improve prediction accuracy. This finding reveals a

36 critical issue in modern predictive models: during training, models inadvertently capture a substantial  
37 number of redundant features. These features not only fail to enhance performance but also interfere  
38 with the learning process, thereby limiting the models' potential to achieve optimal results.

39 Through systematic analysis, we have identified two typical manifestations of redundant features and  
40 their underlying mechanisms. First, input truncation optimization experiments (as shown in Figure  
41 2b and Table 1) demonstrate that selectively masking partial historical data can significantly improve  
42 model prediction performance. This phenomenon reveals the current model's inefficient utilization  
43 of long historical windows. Second, representation similarity analysis (as illustrated in Figure 2a)  
44 shows that both the model's prediction results and intermediate embeddings exhibit an abnormally  
45 concentrated distribution, which significantly deviates from the natural dispersion characteristics  
46 of the input and label. Collectively, these observations indicate that existing models exhibit low  
47 efficiency when processing long historical windows, often encoding substantial noise or irrelevant  
48 variables rather than truly predictive signals.

49 Building upon information bottleneck theory [24, 25, 23, 10], this study proposes an innovative  
50 method called Adaptive Masking Loss with Representation Consistency (AMRC). The core method-  
51 ology comprises: 1) An adaptive masking mechanism that dynamically identifies key segments with  
52 high discriminative power in sequential data and leverages these informative segments to guide the  
53 gradient optimization process (as illustrated in Fig 3) ; 2) A representation consistency constraint that  
54 establishes stable mapping relationships among the input feature space, label space, and predicted  
55 outputs, thereby effectively enhancing the model's generalization capability. Experimental results (as  
56 shown in Table 2) demonstrate that the AMRC method significantly reduces the complexity of the  
57 training solution space by suppressing the model's reliance on redundant features, fully exploits the  
58 performance potential of the model architecture, and consequently improves prediction accuracy.

59 The primary contributions of this study include:

- 60 • **Theoretical Insight:** Through rigorous experimental validation, We demonstrate that existing time  
61 series forecasting models are prone to learning redundant features, which in turn constrain their  
62 performance. Building on the theory of information bottlenecks, we construct a novel theoretical  
63 framework for time series modeling and propose an innovative optimization pathway, offering a  
64 new theoretical perspective for advancing the field of time series forecasting.
- 65 • **Methodological Innovation:** We propose an optimization framework Adaptive Masking Loss with  
66 Representation Consistency. By dynamically selecting discriminative temporal segments to guide  
67 gradient descent (as illustrated in Figure 1) while enforcing input-label-prediction consistency,  
68 our method effectively suppresses redundant feature learning. Extensive experiments demonstrate  
69 consistent performance gains across diverse benchmarks and architectures.

70 Our work advances the understanding of temporal pattern learning mechanisms while offering a  
71 practical pathway to enhance the efficiency and reliability of time series forecasting systems.

## 72 2 Analysis of Redundant Feature Learning

73 Given a multivariate time series  $\mathbf{X} \in \mathbb{R}^{T \times D}$ , where  $T$  is the number of timesteps and  $D$  is the  
74 number of variables, the objective of time series forecasting is to learn a mapping function  $f_\theta$  that  
75 transforms historical observations  $\mathbf{X}_{t-L:t} \in \mathbb{R}^{L \times D}$  (where  $L$  denotes the input length ) into future  
76 values  $\mathbf{X}_{t+1:t+H} \in \mathbb{R}^{H \times D}$  (where  $H$  represents the forecasting horizon).

77 Conventional time series forecasting models follow the long-sequence information gain hypothesis[3,  
78 33, 5, 29], which holds that increasing the input length  $L$  improves forecasting accuracy. However,  
79 our experiments (Table 1) on multiple standard benchmarks reveal a counterintuitive result: truncating  
80 the input—such as masking the first  $k$  timesteps—often improves forecasting performance, which is  
81 measured by Mean Squared Error (MSE). We found that models tend to learn redundant features,  
82 which degrade model performance even after convergence. This finding is supported by two key  
83 observations:

### 84 2.1 Input Truncation Optimization

85 Based on the baseline model configuration (input length  $L = 48$ , forecasting horizon  $H = 48$ ), we  
86 design an input truncation comparative experiment by applying a masking operator  $\mathcal{M}_k(\cdot)$  to the

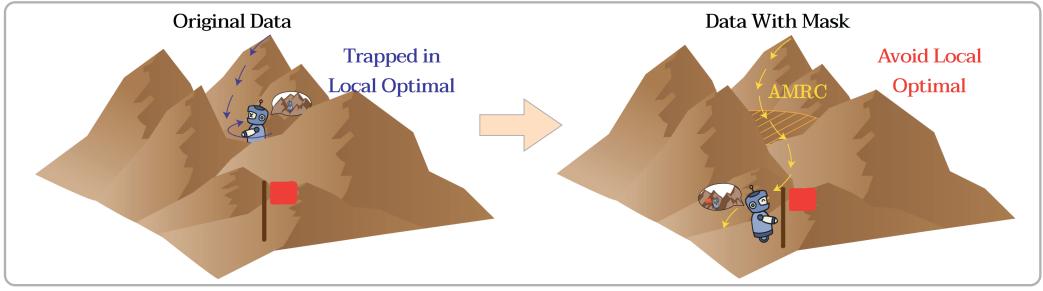


Figure 1: Illustration of the effect of AMRC method. Without regularization, the model tends to overfit redundant input features, leading to suboptimal convergence. By suppressing redundant input features, AMRC restructures the optimization landscape, promoting more efficient representation learning and facilitating better convergence.

87 input sequence. When we have an input sequence of length  $L$  at time step  $t$ , denoted as  $\mathbf{X}_t^{(L)}$ , the  
88 masking operator  $\mathcal{M}_k(\cdot)$  is mathematically defined as:

$$\mathcal{M}_k(\mathbf{X}_t^{(L)}) = \begin{cases} 0 & \text{if } i \leq k \\ \mathbf{X}_t^{(L)} & \text{otherwise} \end{cases} \quad (1)$$

89 Here,  $k \in \{1, \dots, L\}$  denotes the masking step size.

90 To probe redundant features, we employ an Optimal Masking strategy: Given an input sequence of  
91 length  $L$ , we generate  $L$  masked variants  $\{\mathcal{M}_k(\mathbf{X}_t^{(L)})\}_{k=1}^L$  (zero-padded to preserve dimensionality).  
92 For instance,  $k = 5$  yields  $L' = 43$  (first 5 positions zeroed). The optimal mask length  $k^*$  is selected  
93 as the configuration minimizing MSE, thereby defining the theoretical upper bound for redundancy  
94 elimination:

$$k^* = \arg \min_{k \in \{1, 2, \dots, L\}} \mathbb{E} \left[ \left\| f_\theta(\mathcal{M}_k(\mathbf{X}_t^{(L)})) - \mathbf{Y}_t^{(H)} \right\|^2 \right] \quad (2)$$

Table 1: Performance Gains via Optimal Masking Across Time Series Models. Ratio quantifies the percentage of training samples demonstrating prediction error reduction through Optimal Masking, calculated as *number of masked series/number of total series*  $\times 100\%$

Model	Metric	ETTh1			ETTh2			Solar-Energy			Weather		
		MSE	MSE*	Ratio	MSE	MSE*	Ratio	MSE	MSE*	Ratio	MSE	MSE*	Ratio
SOFTS	Train Set	0.278	<b>0.254</b>	56.54%	0.318	<b>0.259</b>	61.65%	0.182	<b>0.155</b>	11.80%	0.421	<b>0.400</b>	45.10%
	Test Set	0.408	<b>0.365</b>	64.24%	0.326	<b>0.303</b>	28.73%	0.293	<b>0.184</b>	41.58%	0.205	<b>0.185</b>	54.93%
iTransformer	Train Set	0.298	<b>0.270</b>	57.87%	0.315	<b>0.261</b>	64.19%	0.410	<b>0.281</b>	61.97%	0.436	<b>0.389</b>	62.98%
	Test Set	0.413	<b>0.289</b>	60.07%	0.329	<b>0.299</b>	32.16%	0.395	<b>0.271</b>	68.43%	0.209	<b>0.170</b>	80.26%
PatchTST	Train Set	0.343	<b>0.303</b>	65.57%	0.329	<b>0.269</b>	69.35%	0.366	<b>0.277</b>	35.89%	0.227	<b>0.180</b>	45.55%
	Test Set	0.424	<b>0.402</b>	65.51%	0.327	<b>0.298</b>	42.46%	0.374	<b>0.344</b>	51.66%	0.215	<b>0.180</b>	42.43%
TSMixer	Train Set	0.372	<b>0.342</b>	55.79%	0.544	<b>0.431</b>	73.96%	0.233	<b>0.195</b>	26.30%	0.363	<b>0.348</b>	37.57%
	Test Set	0.402	<b>0.372</b>	59.19%	0.324	<b>0.289</b>	42.13%	0.288	<b>0.250</b>	40.12%	0.222	<b>0.195</b>	70.88%
TimeMixer	Train Set	0.290	<b>0.262</b>	57.96%	0.309	<b>0.251</b>	59.36%	0.142	<b>0.112</b>	13.58%	0.403	<b>0.353</b>	63.93%
	Test Set	0.393	<b>0.366</b>	58.04%	0.318	<b>0.285</b>	44.52%	0.288	<b>0.253</b>	36.25%	0.197	<b>0.172</b>	66.13%

95 As demonstrated in Table 1, the experimental results confirm that masked models consistently  
96 achieve lower MSE, with more than 50% of samples exhibiting improved predictive performance  
97 ( $\text{Ratio} > 50\%$ ). Notably, the phenomenon of redundancy learning shows strong architecture-agnostic  
98 characteristics. On the Weather dataset, both iTransformer (a Transformer-based model) and TSMixer  
99 (an MLP-based model) demonstrate similar relative improvements: iTransformer achieves an MSE  
100 reduction from **0.209** to **0.170** ( $-18.7\%$ ), while TSMixer improves from **0.222** to **0.195** ( $-12.2\%$ ).  
101 These results indicate that the effectiveness of our masking strategy is not dependent on specific  
102 model architectures.

103 **2.2 Representation Similarity Paradox**

104 To further investigate the redundant feature learning phenomenon, we apply t-SNE to project the  
 105 SOFTS model's high-dimensional representations of the input, embedding, prediction, and label onto  
 106 a 2D plane (Fig. 2a), after normalizing all features to the  $[0, 1]$  range.

107 As illustrated in Fig.2a, Normalized input ( $\mathbf{Z}_{\text{in}} \in \mathbb{R}^L$ ) and output ( $\mathbf{Z}_{\text{out}} \in \mathbb{R}^H$ ) embeddings show  
 108 a clear contrast: inputs remain dispersed, while embeddings and preds cluster tightly despite large  
 109 differences in their corresponding labels. This suggests that the model encodes redundant, task-  
 110 irrelevant features that misrepresent semantic relationships and distort the input-output mapping.

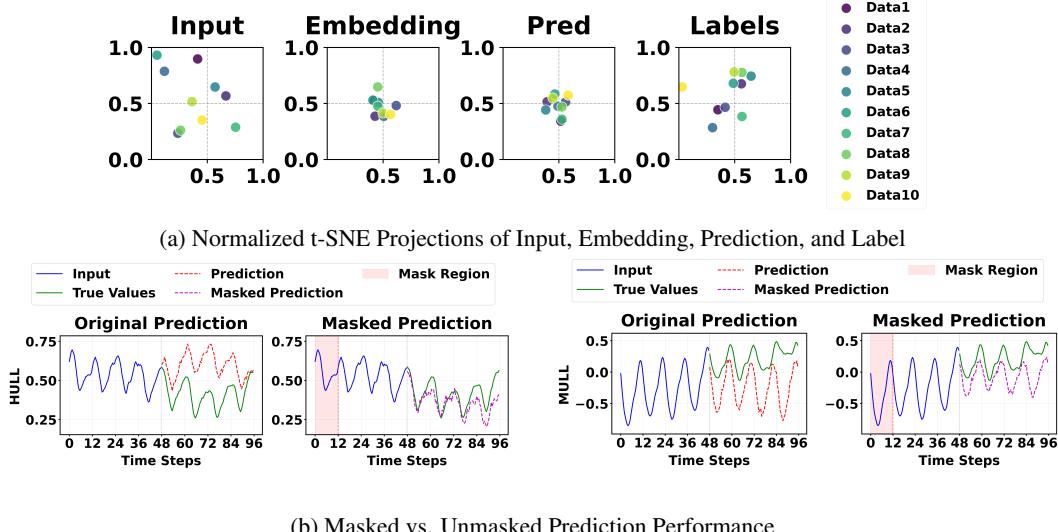


Figure 2: Embedding Distributions and Masking Effects of Our Method.

111 **2.3 Information Bottleneck Constraints on Redundancy**

112 According to the Information Bottleneck (IB) Theory [23], a neural network functions like a bottleneck  
 113 that compresses input information during feature extraction. It discards irrelevant or noisy details and  
 114 retains only the components most relevant to the overall task. For a time series forecasting model, let  
 115 the input be denoted by  $X$ , the latent representation by  $Z$ , and the prediction target by  $Y$ . The model  
 116 aims to learn a representation  $Z$  that maximally preserves information relevant to  $Y$ . This objective  
 117 can be formally expressed as maximizing the mutual information between  $Z$  and  $Y$ :

$$I(Z, Y; \theta) = \int dx dy p(z, y | \theta) \log \frac{p(z, y | \theta)}{p(z | \theta)p(y | \theta)}. \quad (3)$$

118 Due to inherent limitations in the data and model capacity, the amount of information that can be  
 119 extracted and transmitted during training is bounded. Consequently, the representation capacity is  
 120 subject to an upper information constraint  $I_c$ . Based on this, the objective of the time series prediction  
 121 model can be equivalently formulated as the following constrained optimization problem:

$$\max_{\theta} I(Z, Y; \theta) \quad \text{s.t.} \quad I(X, Z; \theta) \leq I_c. \quad (4)$$

122 This constrained optimization problem can be transformed into an unconstrained form using the  
 123 method of Lagrange multipliers, leading to the maximization of the following objective[1]:

$$R_{\text{IB}}(\theta) = I(Z; Y; \theta) - \beta I(Z; X; \theta). \quad (5)$$

124 There are two implementation paths under this objective: one is to maximize the mutual information  
 125  $I(Z; Y)$  between  $Z$  and  $Y$ ; the other is to minimize the mutual information  $I(Z; X)$  between  $Z$  and  $X$ .  
 126 Most current sequential prediction models focus on improving  $I(Z; Y)$  through iterative training, but  
 127 have not explicitly optimized performance by penalizing redundant features via minimizing  $I(Z; X)$ .  
 128 Therefore, we propose an adaptive loss function that aims to minimize the mutual information  
 129 between  $X$  and  $Z$ , offering a novel optimization path for improving the performance of sequential  
 130 prediction models.

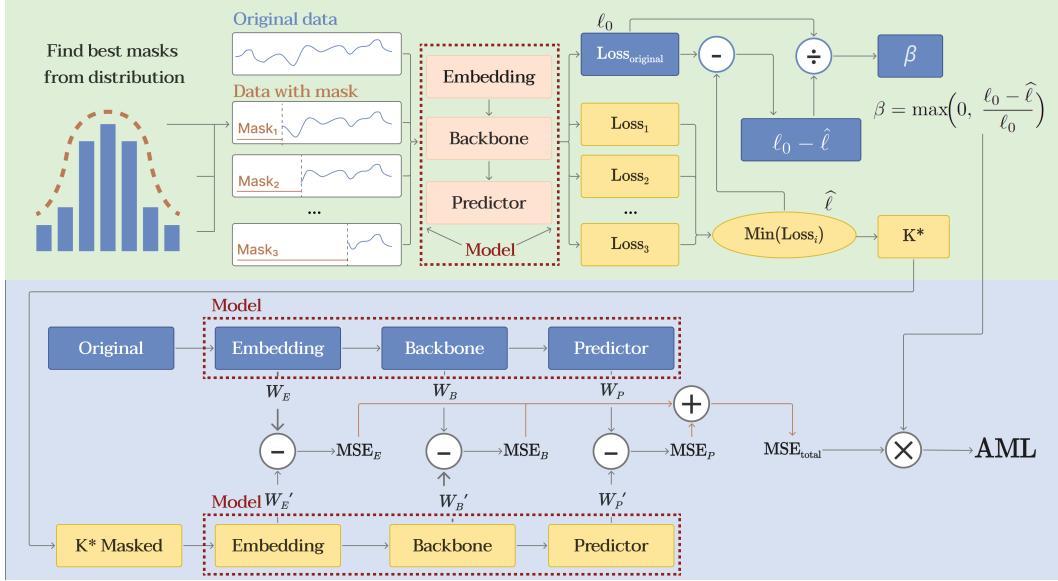


Figure 3: Overview of the Adaptive Masking Loss (AML) framework. The upper half illustrates how the optimal mask length  $K^*$  is selected by evaluating prediction losses over sampled masks. A weighting coefficient  $\beta$  is computed based on the gain over the unmasked loss. The lower half shows the AML loss, calculated as the sum of representation differences between the original input and the  $K^*$  masked input across embedding, backbone, and predictor layers.

### 3 Proposed Method

#### 3.1 Adaptive Masking Loss (AML)

As discussed in Section 2.1, applying ideal masking to input data reduces the information  $I(X)$  while improving prediction accuracy. This indicates that the representation  $Z_{k^*}$ , generated by encoder  $p_\theta$  from masked features  $X_{t,k^*}$ , contains less redundancy and better approximates the minimal sufficient statistics (i.e., with smaller  $I(X, Z_{k^*}; \theta)$ ). Based on this insight, we propose the **Adaptive Masking Loss (AML)** to explicitly reduce mutual information  $I(X, Z; \theta)$  by guiding the encoder’s output representation  $Z$  toward  $Z_{k^*}$ , thereby suppressing redundant feature learning and unleashing model potential. The overall framework of AML is illustrated in Figure 3.

##### 3.1.1 Implementation

The exhaustive search for optimal mask  $k^*$  by enumerating all possible mask lengths  $k \in \{1, \dots, L\}$  results in prohibitive  $O(L)$  time complexity for long sequences. We therefore adopt an efficient stochastic approximation strategy:

- 1. Random Mask Generation:** Independently sample  $m$  mask indices  $\{k_s\}_{s=1}^m$  from uniform distribution  $d(k) = \text{Uniform}\{1, \dots, L\}$ , each generating a masked variant:

$$\tilde{X}_{t,s}^{(L)} = \mathcal{M}_{k_s}(X_t^{(L)}) \quad (6)$$

- 2. Loss Evaluation:** Compute prediction losses for both masked and original data:

$$\ell_s = \mathcal{L}(f_\theta(\tilde{X}_{t,s}^{(L)}), Y_t^{(H)}) \quad (7)$$

$$\ell = \mathcal{L}(f_\theta(X_t^{(L)}), Y_t^{(H)}) \quad (8)$$

- 3. Optimal Representation Selection:** If  $\exists \ell_s < \ell$ , the corresponding representation  $\tilde{Z}_s = p_\theta(\tilde{X}_{t,s}^{(L)})$  satisfies  $I(X_t^{(L)}, \tilde{Z}_s) < I(X_t^{(L)}, Z)$ , where  $Z = p_\theta(X_t^{(L)})$  is the original representation. The optimal mask variant is selected by:

$$s^* = \arg \max_s (\ell - \ell_s) \quad (9)$$

150 **3.1.2 Loss Formulation**

151 To promote compact and informative representations, AML minimizes the distance between the  
 152 original representation  $Z$  and the optimal masked variant  $\tilde{Z}_{s^*}$ :

$$\mathcal{L}_{\text{AML}} = \beta \cdot \frac{1}{D_1 \times D_2} \|Z - \tilde{Z}_{s^*}\|^2 \quad (10)$$

153 where the adaptive weight  $\beta = \max(0, (\ell - \ell_{s^*})/\ell)$  dynamically scales the optimization intensity,  
 154 ensuring stronger influence from mask variants with greater loss reduction.

155 **3.2 Embedding Similarity Penalty (ESP)**

156 Time series forecasting models often encounter two issues: semantic inconsistency, where seman-  
 157 tically similar inputs lead to substantially different predictions, and representation collapse, where  
 158 dissimilar inputs result in nearly identical outputs. Both problems reduce the robustness and gener-  
 159 alization ability of the model. To address these issues, we introduce a regularization strategy that  
 160 compares, for each pair of samples within a mini-batch, the geometry of the embedding space with  
 161 that of the output space.

162 **Pairwise distances.** For a batch  $\mathcal{B} = \{(X_i, Y_i)\}_{i=1}^n$  we denote by  $Z_i = f_{\text{enc}}(X_i) \in \mathbb{R}^{L \times D}$  the  
 163 encoder output and keep the ground-truth  $Y_i \in \mathbb{R}^{P \times D}$ . The (normalised) squared Frobenius distances  
 164 are

$$\Delta_{ij}^E = \frac{1}{L \times D} \|Z_i - Z_j\|_F^2, \quad \Delta_{ij}^O = \frac{1}{P \times D} \|Y_i - Y_j\|_F^2, \quad 1 \leq i, j \leq n. \quad (11)$$

165 **Consistency penalty.** Ideally  $\Delta_{ij}^E$  and  $\Delta_{ij}^O$  should match: semantically similar inputs ( $\Delta_{ij}^E \approx 0$ )  
 166 ought to produce similar outputs ( $\Delta_{ij}^O \approx 0$ ), and vice versa. Deviation is quantified element-wise  
 167 through

$$P_{ij} = \text{ReLU}(\Delta_{ij}^E - \Delta_{ij}^O) + \text{ReLU}(\Delta_{ij}^O - \Delta_{ij}^E) = |\Delta_{ij}^E - \Delta_{ij}^O|_+, \quad (12)$$

168 where  $\text{ReLU}(x) = \max(0, x)$  and  $|\cdot|_+$  denotes the non-negative part. The **Embedding-Similarity**  
 169 **Penalty** then reads

$$\mathcal{L}_{\text{ESP}} = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P_{ij}. \quad (13)$$

170 Equation (13) back-propagates smooth, unbiased gradients that jointly reshape the encoder and the  
 171 predictor so that input and output manifolds remain geometrically aligned. The detailed implemen-  
 172 tation of the Embedding Similarity Penalty (ESP) is provided as pseudocode in Appendix D.

173 **3.3 Overall Training Objective**

174 Section 3.1 introduced the Adaptive Masking Loss  $\mathcal{L}_{\text{AML}}$  that discourages the learning of redundant  
 175 temporal prefixes, while Section 3.2 proposed the Embedding-Similarity Penalty  $\mathcal{L}_{\text{ESP}}$  to enforce  
 176 semantic-behavioural consistency. Combined with the standard prediction loss  $\mathcal{L}_{\text{pred}}$  (e.g., MSE  
 177 between the forecast  $\hat{Y}$  and the target  $Y$ ), our final objective is

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{pred}} + \lambda_{\text{AML}} \mathcal{L}_{\text{AML}} + \lambda_{\text{ESP}} \mathcal{L}_{\text{ESP}}, \quad (14)$$

178 where  $\lambda_{\text{AML}}, \lambda_{\text{ESP}} > 0$  control the strength of each auxiliary term. Minimizing (14) jointly (i)  
 179 identifies the informative prefix for every sequence, (ii) preserves the intrinsic topology of the data,  
 180 and (iii) improves predictive accuracy and interpretability without adding inference-time overhead.

181 **4 Experiment**

182 **4.1 Experiment Setup**

183 **Datasets.** We evaluate our proposed method using seven widely recognized benchmark datasets for  
 184 multivariate time series forecasting: **ETTh1**, **ETTh2**, **ETTm1**, **ETTm2**, **Solar-Energy**, **Electricity**,  
 185 and **Weather**. These datasets encompass a variety of application scenarios with different temporal  
 186 resolutions, seasonality patterns, and dynamic structures. Detailed descriptions of each dataset,  
 187 including their specific characteristics and collection periods, are provided in the appendix F.

188 **Task formulation.** In our experimental setup, the forecasting task is formulated as a sequence-to-  
 189 sequence regression problem, applicable to multivariate time series. Each model is trained to predict  
 190 a future sequence  $\mathbf{Y}_t^{(H)} \in \mathbb{R}^{H \times D}$  from a fixed-length historical input sequence  $\mathbf{X}_t^{(48)} \in \mathbb{R}^{48 \times D}$ ,  
 191 where  $H$  denotes the prediction length and  $D$  is the number of variables. We adopt multiple prediction  
 192 horizons  $H \in \{48, 72, 96, 120, 144, 168, 192\}$ .

193 **Baselines.** Our method is compared against five diverse baseline models: **SOFTS** [8], **iTransformer**  
 194 [14], **PatchTST** [16], **TSMixer** [7], and **TimeMixer** [26]. These baselines are implemented us-  
 195 ing their official codebases and recommended hyperparameters to ensure a fair comparison under  
 196 consistent experimental conditions.

197 **Implementation details.** All models are implemented in PyTorch and trained on a single NVIDIA  
 198 A100 80GB GPU. To ensure a fair comparison and allow both baseline models and those augmented  
 199 with our proposed modules to fully exploit their capacity, we train each model for up to 100 epochs  
 200 using the Adam optimizer with an initial learning rate of  $1 \times 10^{-4}$ , a cosine annealing scheduler, and  
 201 a batch size of 32. Early stopping is applied based on validation loss with a patience of 20 epochs.  
 202 The best-performing checkpoint on the validation set is selected for final evaluation on the test set.

203 **Hyperparameter selection.** For the AML, the input sequence prefix length is configured as  $L = 48$ ,  
 204 with the mask sampling cardinality parameterized as  $m = 12$ . We fix both  $\lambda_{\text{AML}}$  and  $\lambda_{\text{ESP}}$  to 1 for  
 205 all experiments. These settings follow standard benchmark configurations commonly used in time  
 206 series forecasting.

## 207 4.2 Forecasting Results

208 We present the forecasting performance of our method—Adaptive Masking Loss with Representation  
 209 Consistency (AMRC)—in comparison with five representative baseline models across seven widely  
 210 used time series benchmark datasets. Table 2 reports the Mean Squared Error (MSE) and Mean  
 211 Absolute Error (MAE) for each model, both with and without the incorporation of AMRC.

Table 2: Performance Comparison of Time Series Forecasting Models With and Without AMRC. In  
 the experimental results, we highlighted in bold the parts where the AMRC model improved by more  
 than 0.05 in MSE and MAE metrics compared to the baseline model. The detailed hyperparameter  
 configurations for each model can be found in Appendix C. Full results are listed in Appendix E.1

Model	ETTh1		ETTh2		ETTm1		ETTm2		Solar-Energy		Electricity		Weather		
	Metric	MSE	MAE												
SOFTS	Original	0.408	0.414	0.326	0.359	0.484	0.434	0.210	0.285	0.293	0.314	0.169	0.255	0.205	0.234
	AMRC	<b>0.389</b>	<b>0.393</b>	<b>0.311</b>	0.362	<b>0.475</b>	<b>0.423</b>	<b>0.198</b>	<b>0.265</b>	0.290	0.309	<b>0.162</b>	<b>0.244</b>	<b>0.196</b>	<b>0.220</b>
iTransformer	Original	0.413	0.415	0.329	0.362	0.517	0.448	0.213	0.290	0.395	0.352	0.176	0.260	0.209	0.237
	AMRC	<b>0.402</b>	<b>0.399</b>	0.324	<b>0.356</b>	<b>0.502</b>	<b>0.447</b>	0.211	<b>0.280</b>	0.392	<b>0.342</b>	<b>0.163</b>	<b>0.239</b>	<b>0.201</b>	<b>0.221</b>
TimeMixer	Original	0.393	0.408	0.318	0.355	0.466	0.429	0.209	0.285	0.288	0.317	0.194	0.279	0.197	0.237
	AMRC	0.388	<b>0.401</b>	0.316	<b>0.339</b>	<b>0.447</b>	<b>0.405</b>	0.204	<b>0.269</b>	0.284	0.317	<b>0.188</b>	0.277	<b>0.186</b>	<b>0.228</b>
PatchTST	Original	0.424	0.424	0.327	0.358	0.461	0.422	0.211	0.287	0.374	0.382	0.211	0.283	0.215	0.280
	AMRC	<b>0.411</b>	<b>0.415</b>	<b>0.319</b>	0.356	0.456	<b>0.413</b>	<b>0.196</b>	<b>0.271</b>	<b>0.361</b>	0.376	0.207	0.285	0.210	<b>0.264</b>
TSMixer	Original	0.402	0.412	0.324	0.357	0.440	0.413	0.201	0.279	0.288	0.314	0.172	0.258	0.222	0.288
	AMRC	<b>0.386</b>	<b>0.397</b>	0.319	<b>0.340</b>	<b>0.432</b>	0.412	0.196	<b>0.257</b>	<b>0.280</b>	0.313	0.169	<b>0.247</b>	<b>0.212</b>	<b>0.281</b>

212 **Consistent Performance Gains.** Across all models and datasets, our method consistently yields  
 213 performance improvements. For example, the MSE of the SOFTS model decreases from 0.408  
 214 to 0.389 on the ETTh1 dataset. Similar trends are observed in iTransformer, where the MSE on  
 215 Electricity drops from 0.176 to 0.163. The enhancements demonstrate that AMRC effectively  
 216 mitigates redundant or noisy temporal segments, thereby improving prediction stability and accuracy.

217 **Architecture-Agnostic Effectiveness.** AMRC delivers significant performance gains not only on  
 218 Transformer-based architectures such as iTransformer and PatchTST, but also on MLP-based models  
 219 including TimeMixer, SOFTS, and TSMixer. For instance, on the ETTm2 dataset, the MSE of  
 220 PatchTST model decreases from 0.211 to 0.196 (a reduction of approximately 7.11%), while the  
 221 MSE of SOFTS model drops from 0.210 to 0.198 (approximately 5.71% reduction). These results  
 222 demonstrate the strong architecture-agnostic generalization ability of AMRC, highlighting its broad  
 223 applicability across a wide range of time series forecasting models.

Table 3: Ablation Study Results on Different Model Components

Model		ETTh1		ETTh2		Solar-Energy		Weather	
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
SOFTS	AML only	0.401	0.405	0.322	0.358	0.297	0.309	0.192	0.228
	ESP only	0.393	0.398	0.318	0.351	0.295	0.318	0.208	0.241
	AMRC	<b>0.389</b>	<b>0.393</b>	<b>0.311</b>	<b>0.362</b>	<b>0.290</b>	<b>0.309</b>	<b>0.196</b>	<b>0.220</b>
iTransformer	AML only	0.410	0.413	0.328	0.363	0.398	0.347	0.205	0.230
	ESP only	0.407	0.408	0.326	0.359	0.402	0.351	0.210	0.248
	AMRC	<b>0.402</b>	<b>0.399</b>	<b>0.324</b>	<b>0.356</b>	<b>0.392</b>	<b>0.342</b>	<b>0.201</b>	<b>0.221</b>
TimeMixer	AML only	0.395	0.412	0.319	0.351	0.287	0.319	0.189	0.232
	ESP only	0.391	0.406	0.317	0.347	0.293	0.325	0.202	0.248
	AMRC	<b>0.388</b>	<b>0.401</b>	<b>0.316</b>	<b>0.339</b>	<b>0.284</b>	<b>0.317</b>	<b>0.186</b>	<b>0.228</b>
PatchTST	AML only	0.419	0.420	0.325	0.361	0.369	0.379	0.214	0.274
	ESP only	0.417	0.418	0.323	0.357	0.375	0.384	0.217	0.281
	AMRC	<b>0.411</b>	<b>0.415</b>	<b>0.319</b>	<b>0.356</b>	<b>0.361</b>	<b>0.376</b>	<b>0.210</b>	<b>0.264</b>
TSMixer	AML only	0.396	0.404	0.324	0.356	0.285	0.317	0.216	0.283
	ESP only	0.390	0.399	0.322	0.352	0.291	0.323	0.224	0.292
	AMRC	<b>0.386</b>	<b>0.397</b>	<b>0.319</b>	<b>0.340</b>	<b>0.280</b>	<b>0.313</b>	<b>0.212</b>	<b>0.281</b>

224 **Generalization on Low-Channel Datasets.** On datasets with fewer input channels (ETTh1, ETTh2,  
225 ETTm1, ETTm2), AMRC effectively enhances model performance. For instance, on ETTm1, the  
226 MSE of iTransformer decreases from 0.517 to 0.502, and that of TSMixer drops from 0.440 to 0.432.  
227 These results demonstrate AMRC’s ability to mitigate overfitting and improve prediction accuracy in  
228 low-dimensional time series forecasting tasks.

229 **Robustness on High-Channel Datasets.** For high-dimensional datasets such as Weather (21 chan-  
230 nels) and Solar-Energy (137 channels) see in Appendix F, AMRC consistently improves robustness  
231 by reducing the impact of signal noise and inter-channel redundancy. On the Weather dataset,  
232 TimeMixer’s MSE decreases from 0.197 to 0.186 and MAE from 0.237 to 0.228, while iTransformer  
233 sees an MAE drop from 0.237 to 0.221. On Solar-Energy, PatchTST’s MSE drops from 0.374 to  
234 0.361, and SOFTS sees a slight MAE reduction from 0.314 to 0.309. These enhancements highlight  
235 AMRC’s effectiveness in managing complexity in multivariate time series with high channel counts.

236 **Generalizable Training Framework.** The consistent performance improvements observed across all  
237 models validate the strong scalability and integrability of AMRC. As a constraint-based optimization  
238 strategy, AMRC does not rely on any specific model architecture, making it highly generalizable. It  
239 serves as a versatile training framework for enhancing both the efficiency and accuracy of time series  
240 forecasting models.

### 241 4.3 Ablation Study

242 **Setup.** We evaluate ablation variants on four diverse datasets: ETTh1 and ETTh2, representing  
243 hourly electricity load with varying degrees of seasonality; Solar-Energy, which exhibits weather-  
244 driven variability and periodicity; and Weather, a multivariate meteorological dataset with complex  
245 inter-variable dependencies. We adopt a fixed input horizon following standard benchmarks.

246 **Evaluation protocol.** For each dataset, we apply the ablation study to five baseline models SOFTS,  
247 iTransformer, TimeMixer, PatchTST, and TSMixer under four configurations: 1) baseline + AML, 2)  
248 baseline + ESP, and 3) baseline + both AML and ESP. This design allows us to assess the standalone  
249 effectiveness of each module as well as their combined synergy.

250 **Findings.** We evaluate the individual and joint effects of the AML and ESP components using five  
251 representative forecasting architectures across four datasets. As shown in Table 3, both components  
252 contribute measurable performance gains in isolation, while their combination AMRC consistently  
253 leads to the best forecasting accuracy in terms of MSE and MAE. AML provides stronger improve-  
254 ments across most settings, supporting its role in suppressing redundant prefixes during training.  
255 ESP, while often delivering smaller standalone gains, remains beneficial by promoting geometric  
256 alignment between embedding and output spaces. Together, these findings demonstrate that each  
257 component addresses a distinct source of generalization error.

258 **Component impact across architectures.** The benefits of AML and ESP are consistently ob-  
259 served across all backbone models, regardless of architectural differences. For instance, models

260 with strong expressiveness, such as iTransformer and TimeMixer, benefit significantly from AML,  
 261 achieving notable MSE reductions on datasets like Weather and ETTh2. Even architectures with-  
 262 out attention mechanisms, such as SOFTS and TSMixer, exhibit consistent gains, highlighting the  
 263 broad applicability of adaptive prefix masking. In contrast, the improvements from ESP are often  
 264 more dataset-dependent, being particularly effective on high-dimensional multivariate inputs where  
 265 representation alignment plays a critical role. For example, ESP yields non-trivial reductions in MAE  
 266 on Weather, where multiple variables evolve under shared dynamics. Notably, we observe relatively  
 267 smaller improvements on the Solar-Energy dataset for transformer-based models such as PatchTST  
 268 and iTransformer, which may be attributed to their reliance on longer input sequences for stable  
 269 attention computation.

270 **Complementarity and synergy.** The AMRC configuration, which jointly applies AML and ESP,  
 271 consistently outperforms its ablated variants across all benchmarks. The performance improvement  
 272 from combining both components generally exceeds the stronger of the two individual effects,  
 273 indicating synergistic interaction. This complementarity can be attributed to their distinct operational  
 274 scopes: AML operates on the input level by learning to suppress non-informative temporal segments,  
 275 while ESP regularizes the latent space to align representations across semantically related inputs.  
 276 As a result, AMRC improves both the quality of features learned from the data and the consistency  
 277 of their usage in prediction. The robust gains observed across datasets and architectures suggest  
 278 that jointly addressing input redundancy and representation inconsistency is critical for improving  
 generalization in time series forecasting.

Table 4: AMRC Effectiveness Across Datasets and Models. Ratio is the percentage of training samples with reduced MSE under prefix masking. Ratio\* is the same metric after training with AMRC, reflecting improved robustness. Results are from the ablation setting with input length set to 48. Detailed results are provided in Appendix E.1.

Model	ETTh1		ETTh2		Solar-Energy		Weather		
	Metric	Ratio	Ratio*	Ratio	Ratio*	Ratio	Ratio*	Ratio	Ratio*
SOFTS		64%	<b>57.33%</b>	28.72%	<b>20.28%</b>	41.58%	<b>33.49%</b>	54.93%	<b>47.12%</b>
iTransformer		60.07%	<b>49.95%</b>	32.16%	<b>23.28%</b>	68.43%	<b>63.21%</b>	80.26%	<b>70.29%</b>
TimeMixer		58.04%	<b>46.29%</b>	44.52%	<b>34.17%</b>	36.25%	<b>27.90%</b>	66.13%	<b>52.28%</b>
PatchTST		65.51%	<b>51.63%</b>	42.46%	<b>26.19%</b>	51.66%	<b>47.64%</b>	42.43%	<b>30.78%</b>
TSMixer		59.19%	<b>46.62%</b>	42.13%	<b>27.98%</b>	40.12%	<b>28.36%</b>	70.88%	<b>58.23%</b>

279  
 280 **Effectiveness of AMRC in Reducing Redundant Features** We evaluate the model’s robustness to  
 281 redundant input by computing the proportion of training samples with improved MSE under prefix  
 282 masking Ratio and compare it to the value after applying AMRC Ratio\*. As shown in Table 7,  
 283 AMRC consistently improves or maintains this ratio, indicating its effectiveness in suppressing the  
 284 impact of redundant temporal information.

## 285 5 Conclusion

286 This study pioneers the investigation into the negative effects of redundant feature learning in time  
 287 series forecasting and introduces AMRC, a plug-and-play solution that suppresses such learning  
 288 without requiring architectural modifications. Unlike prior work focused on enhancing predictive  
 289 features, AMRC improves accuracy by reducing reliance on redundant features while maintaining  
 290 model flexibility. Its key advantages include: 1) seamless integration with existing models, 2) effective  
 291 suppression of feature redundancy, and 3) strong generalization performance across benchmark tests.  
 292 By addressing the long-overlooked issue of redundant learning, this research provides a novel and  
 293 practical methodology for optimizing forecasting models.

## 294 References

- 295 [1] Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational  
 296 information bottleneck, 2019. URL <https://arxiv.org/abs/1612.00410>.

- 297 [2] Peng Chen, Yingying Zhang, Yunyao Cheng, Yang Shu, Yihang Wang, Qingsong Wen, Bin  
 298 Yang, and Chenjuan Guo. Pathformer: Multi-scale transformers with adaptive pathways for  
 299 time series forecasting. *arXiv preprint arXiv:2402.05956*, 2024.
- 300 [3] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov.  
 301 Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint  
 302 arXiv:1901.02860*, 2019.
- 303 [4] Abhimanyu Das, Weihao Kong, Andrew Leach, Shaan Mathur, Rajat Sen, and Rose Yu. Long-  
 304 term forecasting with tide: Time-series dense encoder, 2024. URL [https://arxiv.org/abs/  
 305 2304.08424](https://arxiv.org/abs/2304.08424).
- 306 [5] Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, Nanning  
 307 Zheng, and Furu Wei. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint  
 308 arXiv:2307.02486*, 2023.
- 309 [6] Alexandre Drouin, Étienne Marcotte, and Nicolas Chapados. Tactis: Transformer-attentional  
 310 copulas for time series. In *International Conference on Machine Learning*, pages 5447–5493.  
 311 PMLR, 2022.
- 312 [7] Vijay Ekambaram, Arindam Jati, Nam Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam.  
 313 Tsmixer: Lightweight mlp-mixer model for multivariate time series forecasting. In *Proceedings  
 314 of the 29th ACM SIGKDD conference on knowledge discovery and data mining*, pages 459–469,  
 315 2023.
- 316 [8] Lu Han, Xu-Yang Chen, Han-Jia Ye, and De-Chuan Zhan. Softs: Efficient multivariate time  
 317 series forecasting with series-core fusion. *arXiv preprint arXiv:2404.14197*, 2024.
- 318 [9] Hansika Hewamalage, Christoph Bergmeir, and Kasun Bandara. Recurrent neural networks  
 319 for time series forecasting: Current status and future directions. *International Journal of  
 320 Forecasting*, 37(1):388–427, 2021.
- 321 [10] Shizhe Hu, Zhengzheng Lou, Xiaoqiang Yan, and Yangdong Ye. A survey on information  
 322 bottleneck. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- 323 [11] Natalia Kashpruk, Cezary Piskor-Ignatowicz, and Jerzy Baranowski. Time series prediction in  
 324 industry 4.0: a comprehensive review and prospects for future advancements. *Applied Sciences*,  
 325 13(22):12374, 2023.
- 326 [12] Bryan Lim and Stefan Zohren. Time-series forecasting with deep learning: a survey. *Philosophical  
 327 Transactions of the Royal Society A*, 379(2194):20200209, 2021.
- 328 [13] Shengsheng Lin, Weiwei Lin, Wentai Wu, Feiyu Zhao, Ruichao Mo, and Haotong Zhang.  
 329 Segrrnn: Segment recurrent neural network for long-term time series forecasting. *arXiv preprint  
 330 arXiv:2308.11200*, 2023.
- 331 [14] Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.  
 332 itransformer: Inverted transformers are effective for time series forecasting. *arXiv preprint  
 333 arXiv:2310.06625*, 2023.
- 334 [15] Chao Ma, Yikai Hou, Xiang Li, Yinggang Sun, and Haining Yu. Long input sequence network  
 335 for long time series forecasting. *arXiv preprint arXiv:2407.15869*, 2024.
- 336 [16] Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is  
 337 worth 64 words: Long-term forecasting with transformers. In *International Conference on  
 338 Learning Representations*, 2023.
- 339 [17] Boris N Oreshkin, Dmitri Carpov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis  
 340 expansion analysis for interpretable time series forecasting. *arXiv preprint arXiv:1905.10437*,  
 341 2019.
- 342 [18] Asiya K Ozcanli, Fatma Yaprakdal, and Mustafa Baysal. Deep learning methods and applications  
 343 for electrical power systems: A comprehensive review. *International Journal of Energy  
 344 Research*, 44(9):7136–7157, 2020.
- 345 [19] Fotios Petropoulos, Daniele Apiletti, Vassilios Assimakopoulos, Mohamed Zied Babai, Devon K  
 346 Barrow, Souhaib Ben Taieb, Christoph Bergmeir, Ricardo J Bessa, Jakub Bijak, John E Boylan,  
 347 et al. Forecasting: theory and practice. *International Journal of forecasting*, 38(3):705–871,  
 348 2022.

- 349 [20] Yankun Ren, Longfei Li, Xinxing Yang, and Jun Zhou. Autotransformer: Automatic transformer  
 350 architecture design for time series classification. In *Pacific-Asia Conference on Knowledge  
 351 Discovery and Data Mining*, pages 143–155. Springer, 2022.
- 352 [21] Koushik Roy, Abtahi Ishmam, and Kazi Abu Taher. Demand forecasting in smart grid using long  
 353 short-term memory. In *2021 International Conference on Automation, Control and Mechatronics  
 354 for Industry 4.0 (ACMI)*, pages 1–5. IEEE, 2021.
- 355 [22] Zezhi Shao, Fei Wang, Yongjun Xu, Wei Wei, Chengqing Yu, Zhao Zhang, Di Yao, Tao  
 356 Sun, Guangyin Jin, Xin Cao, et al. Exploring progress in multivariate time series forecasting:  
 357 Comprehensive benchmarking and heterogeneity analysis. *IEEE Transactions on Knowledge  
 358 and Data Engineering*, 2024.
- 359 [23] Naom Slonim and Naftali Tishby. Agglomerative information bottleneck. *Advances in neural  
 360 information processing systems*, 12, 1999.
- 361 [24] Naftali Tishby and Noga Zaslavsky. Deep learning and the information bottleneck principle. In  
 362 *2015 ieee information theory workshop (itw)*, pages 1–5. Ieee, 2015.
- 363 [25] Naftali Tishby, Fernando C. Pereira, and William Bialek. The information bottleneck method,  
 364 2000. URL <https://arxiv.org/abs/physics/0004057>.
- 365 [26] Shiyu Wang, Haixu Wu, Xiaoming Shi, Tengge Hu, Huakun Luo, Lintao Ma, James Y Zhang,  
 366 and Jun Zhou. Timemixer: Decomposable multiscale mixing for time series forecasting. *arXiv  
 367 preprint arXiv:2405.14616*, 2024.
- 368 [27] Zhijian Xu, Ailing Zeng, and Qiang Xu. Fits: Modeling time series with 10 $k$  parameters. *arXiv  
 369 preprint arXiv:2307.03756*, 2023.
- 370 [28] Kun Yi, Qi Zhang, Wei Fan, Shoujin Wang, Pengyang Wang, Hui He, Ning An, Defu Lian,  
 371 Longbing Cao, and Zhendong Niu. Frequency-domain mlps are more effective learners in time  
 372 series forecasting. *Advances in Neural Information Processing Systems*, 36:76656–76679, 2023.
- 373 [29] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti,  
 374 Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird:  
 375 Transformers for longer sequences. *Advances in neural information processing systems*, 33:  
 376 17283–17297, 2020.
- 377 [30] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series  
 378 forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages  
 379 11121–11128, 2023.
- 380 [31] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series  
 381 forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages  
 382 11121–11128, 2023.
- 383 [32] Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency  
 384 for multivariate time series forecasting. In *The eleventh international conference on learning  
 385 representations*, 2023.
- 386 [33] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wanhai  
 387 Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In  
 388 *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115,  
 389 2021.
- 390 [34] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer:  
 391 Frequency enhanced decomposed transformer for long-term series forecasting. In *International  
 392 conference on machine learning*, pages 27268–27286. PMLR, 2022.

393 **A Related Work**

394 The Information Bottleneck (IB) method was first introduced by Tishby et al. [25] as an information-  
395 theoretic framework that aims to compress input signals while preserving as much relevant information  
396 as possible about the target output. In the field of machine learning, IB theory has been widely adopted  
397 as a regularization technique. For instance, Alemi et al.[1] proposed the Variational Information  
398 Bottleneck (VIB), which leverages variational inference to construct a tractable lower bound on  
399 the IB objective. Building upon this, Tishby and Zaslavsky[24]further explored the applicability of  
400 information-theoretic objectives to deep neural networks. Research on IB has also extended into  
401 the domain of clustering. Slonim et al.[23] developed a distributional clustering algorithm based on  
402 mutual information maximization and demonstrated its effectiveness on the 20 Newsgroups dataset,  
403 achieving substantial compression with minimal loss of relevant information. More recently, Hu  
404 et al.[10] conducted a comprehensive survey of the IB literature, reviewing over two decades of  
405 theoretical developments, methodological advances, and practical applications in both traditional  
406 machine learning and deep learning settings, aiming to provide a unified perspective for future  
407 research in this area.

408 Time series forecasting has attracted significant attention due to its wide-ranging applications in  
409 fields such as energy, transportation, and economics. In the context of deep learning, time series  
410 forecasting methods can be broadly categorized into MLP-based, RNN-based, and Transformer-based  
411 approaches. Among MLP-based models, DLinear[30] and TSMixer[7] are representative examples,  
412 featuring relatively simple architectures while achieving strong performance across multiple datasets.  
413 RNN-based methods, such as SegRnn[6] and LSTMlong[21], focus on structural modifications to  
414 address challenges related to parallel prediction and long-sequence modeling. Transformer-based  
415 models include Informer[33], Autoformer[20], and iTransformer[14]. Informer introduces a sparse  
416 attention mechanism to improve the scalability of traditional attention for time series modeling;  
417 Autoformer incorporates frequency-domain information to enhance attention; and iTransformer  
418 further extends attention across channels by embedding multivariate sequences for variable-aware  
419 representation.

420 **B Limitations**

421 Despite the demonstrated effectiveness (Table 2) of our Adaptive Masking Loss with Representation  
422 Consistency (AMRC) in suppressing redundant feature learning, several limitations remain and  
423 warrant further investigation and refinement in future work.

424 **1. High Computational Overhead**

425 In AMRC, the Adaptive Masking Loss (AML) requires the generation of  $m$  masked variants  
426 per training batch, each undergoing a forward pass to select the optimal representation. This  
427 mechanism results in a per-step training cost approximately  $m$  times higher than standard  
428 training. In particular, as the input sequence length  $L$  increases, a larger  $m$  is typically  
429 needed to ensure effective mask selection, which further exacerbates the computational  
430 burden. While this overhead is acceptable in offline training scenarios, it may become a  
431 bottleneck in latency-sensitive applications such as high-frequency trading or real-time fault  
432 detection.

433 **2. Distance Degeneration in High-Dimensional Spaces**

434 The Embedding Similarity Penalty (ESP) relies on normalized Frobenius distances ( $\Delta E$  and  
435  $\Delta O$ ) to quantify geometric discrepancies between embedding and output spaces. However,  
436 in high-dimensional embedding spaces (e.g., when  $D > 64$ ), euclidean distances tend  
437 to concentrate, causing inter-sample distances to shrink and resulting in  $P_{ij} \approx 0$ , which  
438 weakens the regularization signal. This effect is especially pronounced when the output  
439 length  $P$  is large. Moreover, as shown in Table3, the benefit of ESP is highly dataset-  
440 dependent: significant performance gains are observed in datasets with strong multivariate  
441 interactions (e.g., *Weather*), while improvements are more limited in datasets with simpler  
442 structure or weaker inter-variable coupling. This suggests that the effectiveness of ESP may  
443 be constrained by the complexity of the task-specific data distribution.

444 **3. Residual Redundant Features**

445 Although AMRC effectively reduces the influence of redundant features through AML and

446       ESP, it does not eliminate them entirely. This limitation mainly arises from the approximation  
 447       nature of AML: the optimal mask  $k^*$  is selected from a finite set of  $m$  randomly sampled  
 448       variants, introducing the possibility of missing the true theoretical optimum. Nonetheless,  
 449       this strategy strikes a practical balance between computational efficiency and redundancy  
 450       suppression. Future work may explore dynamic sampling strategies or learning-based  
 451       mask generation mechanisms to better approximate the optimal mask and further enhance  
 452       redundancy elimination.

## 453     C Details of the Baseline Model

454     All models are reproduced based on their official open-source implementations:

- 455       1. **SOFTS** from <https://github.com/Secilia-Cxy/SOFTS>.
- 456       2. **TimeMixer** from <https://github.com/kwuking/TimeMixer>.
- 457       3. **iTransformer** from <https://github.com/thuml/iTransformer>.
- 458       4. **PatchTST** from <https://github.com/yuqinie98/PatchTST>.
- 459       5. **TSMixer** from <https://github.com/ditschuk/pytorch-tsmixer>.

460     The hyperparameters for each model on different datasets follow the official configurations provided in  
 461     their corresponding GitHub repositories. For the PatchTST model on the Solar-Energy dataset, since  
 462     no official configuration was provided, we adopted the hyperparameter settings from iTransformer.

## 463     D Model Detail

### 464     D.1 ESP

---

**Algorithm 1:** Embedding-Similarity Penalty (ESP) for Time Series Forecasting

---

**Input:** Mini-batch  $\mathcal{B} = \{(X_i, Y_i)\}_{i=1}^n$ , encoder  $f_{\text{enc}}$ , predictor  $f_{\text{pred}}$   
**Output:** Penalty loss  $\mathcal{L}_{\text{ESP}}$

```

1 1. Forward pass to compute encoder outputs
2 for  $i \leftarrow 1$  to  $n$  do
3    $Z_i \leftarrow f_{\text{enc}}(X_i)$                                      #encoder output  $\in \mathbb{R}^{L \times D}$ 
4 end
5 2. Compute pairwise Frobenius distances
6 Initialize  $\Delta^E, \Delta^O \in \mathbb{R}^{n \times n}$ 
7 for  $i \leftarrow 1$  to  $n$  do
8   for  $j \leftarrow i$  to  $n$  do
9      $\Delta_{ij}^E \leftarrow \frac{1}{L \times D} \|Z_i - Z_j\|_F^2$            #embedding similarity
10     $\Delta_{ij}^O \leftarrow \frac{1}{P \times D} \|Y_i - Y_j\|_F^2$            #output similarity
11     $\Delta_{ji}^E \leftarrow \Delta_{ij}^E,$ 
12     $\Delta_{ji}^O \leftarrow \Delta_{ij}^O$                                      #symmetry
13  end
14 end
15 3. Compute pairwise penalties
16 for  $i \leftarrow 1$  to  $n$  do
17   for  $j \leftarrow 1$  to  $n$  do
18      $P_{ij} \leftarrow |\Delta_{ij}^E - \Delta_{ij}^O|_+$                  #element-wise consistency penalty
19   end
20 end
21 4. Compute final regularization loss
22  $\mathcal{L}_{\text{ESP}} \leftarrow \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P_{ij}$ 
23 5. Backward pass and update
24 Update  $\theta$  using forecasting loss  $+ \lambda_{\text{ESP}} \cdot \mathcal{L}_{\text{ESP}}$ 

```

---

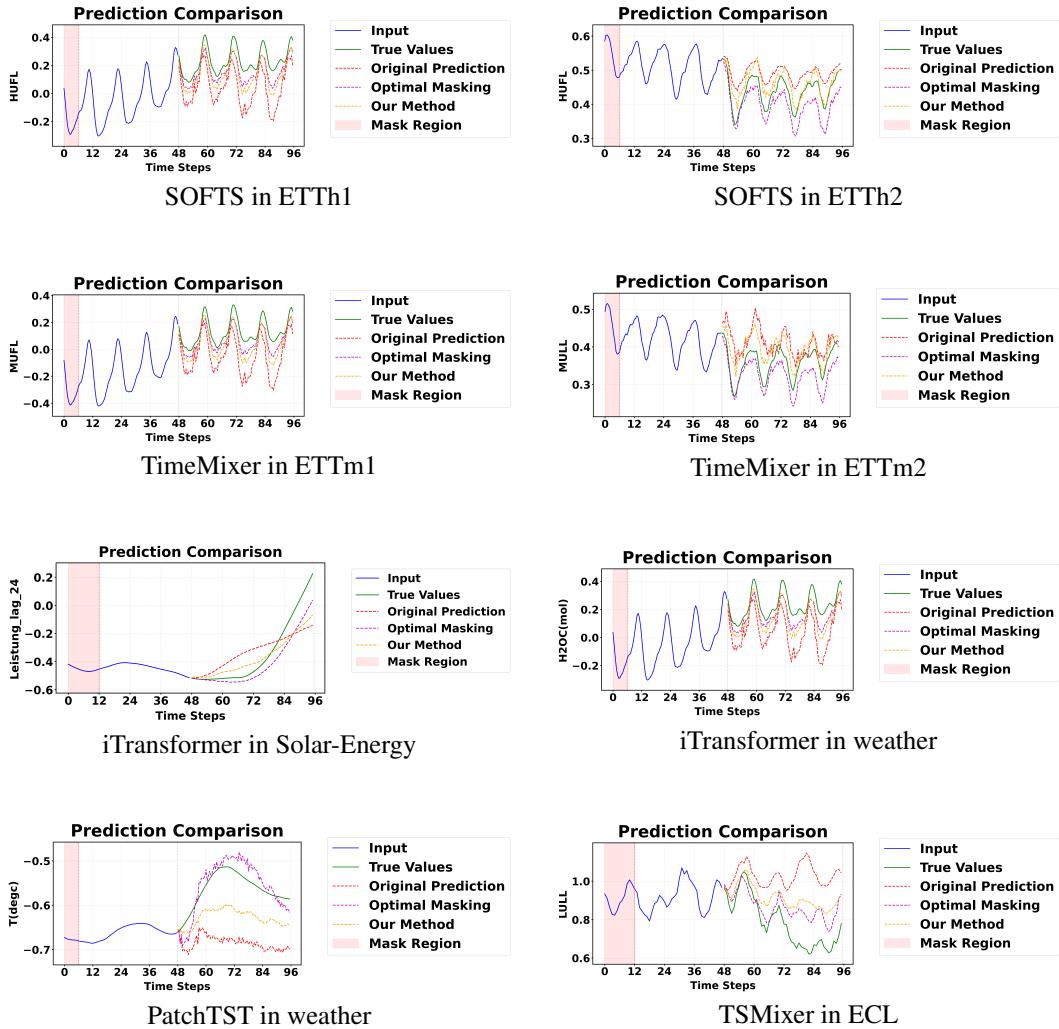


Table 7: Average AMRC Effectiveness Across Datasets and Models. Ratio is the percentage of samples with reduced MSE under ideal masking. Ratio\* is the same metric after training with AMRC, reflecting improved robustness.

Models	SOFTS		TimeMixer		iTransformer		PatchTST		TSMixer			
	Metric	Ratio	Ratio*	Metric	Ratio	Ratio*	Metric	Ratio	Ratio*	Metric	Ratio	Ratio*
ETTh1	57.14%	48.7%	49.69%	37.81%	51.88%	42.93%	57.81%	42.91%	52.92%	59.1%		
ETTh2	30.99%	21.09%	47.16%	34.89%	33.28%	24.11%	43.54%	27.91%	44.29%	29.63%		
Solar-Energy	44.87%	33.83%	37.52%	29.61%	71.18%	67.72%	53.26%	48.02%	41.66%	30.11%		
Weather	54.63%	48.33%	67.39%	52.78%	79.4%	69.36%	41.98%	29.86%	69.32%	56.26%		

## 468 E.2 Visualized Prediction Comparison Chart

Figure 4



469 **F Dataset description**

470 Here we provide detailed descriptions along with download links for each dataset:

- 471 1. **ETT (Electricity Transformer Temperature)** [33]<sup>1</sup>: This collection includes two hourly-resolution datasets (ETTh) and two 15-minute-resolution datasets (ETTm). Each dataset captures seven key operational metrics (including oil and load measurements) from electricity transformers, spanning from July 2016 to July 2018.
- 472 2. **Electricity**<sup>2</sup>: Comprising hourly power consumption records from 321 customers, this dataset covers the period from 2012 to 2014.
- 473 3. **Weather**: Featuring 21 meteorological indicators (such as air temperature and humidity), this dataset provides 10-minute-interval recordings throughout 2020, sourced from weather stations in Germany.
- 474 4. **Solar-Energy**: Documents the solar power generation output of 137 photovoltaic plants in 2006, with measurements taken at 10-minute intervals.

Table 8: Detailed Dataset Descriptions. The table summarizes key characteristics of the time series datasets, including the number of channels, prediction lengths, dataset splits, temporal granularity, and application domains.

Dataset	Channels	Prediction Length	Dataset Split (Train, Val, Test)	Granularity	Domain
ETTh1, ETTh2	7	{48, 72, 96, 120, 144, 168, 192}	(8545, 2881, 2881)	Hourly	Electricity
ETTm1, ETTm2	7	{48, 72, 96, 120, 144, 168, 192}	(34465, 11521, 11521)	15min	Electricity
Weather	21	{48, 72, 96, 120, 144, 168, 192}	(36792, 5271, 10540)	10min	Weather
ECL	321	{48, 72, 96, 120, 144, 168, 192}	(18317, 2633, 5261)	Hourly	Electricity
Solar-Energy	137	{48, 72, 96, 120, 144, 168, 192}	(36601, 5161, 10417)	10min	Energy

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

<sup>2</sup><https://archive.ics.uci.edu/dataset/321/electricityloaddiagrams20112014>

482 **NeurIPS Paper Checklist**

483 The checklist is designed to encourage best practices for responsible machine learning research,  
484 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
485 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
486 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
487 towards the page limit.

488 Please read the checklist guidelines carefully for information on how to answer these questions. For  
489 each question in the checklist:

- 490 • You should answer [Yes] , [No] , or [NA] .  
491 • [NA] means either that the question is Not Applicable for that particular paper or the  
492 relevant information is Not Available.  
493 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

494 **The checklist answers are an integral part of your paper submission.** They are visible to the  
495 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
496 (after eventual revisions) with the final version of your paper, and its final version will be published  
497 with the paper.

498 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
499 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
500 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
501 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
502 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
503 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
504 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
505 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
506 please point to the section(s) where related material for the question can be found.

507 **IMPORTANT**, please:

- 508 • **Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",**  
509 • **Keep the checklist subsection headings, questions/answers and guidelines below.**  
510 • **Do not modify the questions and only use the provided macros for your answers.**

511 **1. Claims**

512 Question: Do the main claims made in the abstract and introduction accurately reflect the  
513 paper's contributions and scope?

514 Answer: [Yes]

515 Justification: The main claims are clearly written in the abstract and introduction.

516 Guidelines:

- 517 • The answer NA means that the abstract and introduction do not include the claims  
518 made in the paper.  
519 • The abstract and/or introduction should clearly state the claims made, including the  
520 contributions made in the paper and important assumptions and limitations. A No or  
521 NA answer to this question will not be perceived well by the reviewers.  
522 • The claims made should match theoretical and experimental results, and reflect how  
523 much the results can be expected to generalize to other settings.  
524 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
525 are not attained by the paper.

526 **2. Limitations**

527 Question: Does the paper discuss the limitations of the work performed by the authors?

528 Answer: [Yes]

529 Justification: We discussed the limitation of our method in Appendix A.

530 Guidelines:

- 531 • The answer NA means that the paper has no limitation while the answer No means that  
532 the paper has limitations, but those are not discussed in the paper.
- 533 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 534 • The paper should point out any strong assumptions and how robust the results are to  
535 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
536 model well-specification, asymptotic approximations only holding locally). The authors  
537 should reflect on how these assumptions might be violated in practice and what the  
538 implications would be.
- 539 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
540 only tested on a few datasets or with a few runs. In general, empirical results often  
541 depend on implicit assumptions, which should be articulated.
- 542 • The authors should reflect on the factors that influence the performance of the approach.  
543 For example, a facial recognition algorithm may perform poorly when image resolution  
544 is low or images are taken in low lighting. Or a speech-to-text system might not be  
545 used reliably to provide closed captions for online lectures because it fails to handle  
546 technical jargon.
- 547 • The authors should discuss the computational efficiency of the proposed algorithms  
548 and how they scale with dataset size.
- 549 • If applicable, the authors should discuss possible limitations of their approach to  
550 address problems of privacy and fairness.
- 551 • While the authors might fear that complete honesty about limitations might be used by  
552 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
553 limitations that aren't acknowledged in the paper. The authors should use their best  
554 judgment and recognize that individual actions in favor of transparency play an impor-  
555 tant role in developing norms that preserve the integrity of the community. Reviewers  
556 will be specifically instructed to not penalize honesty concerning limitations.

557 **3. Theory assumptions and proofs**

558 Question: For each theoretical result, does the paper provide the full set of assumptions and  
559 a complete (and correct) proof?

560 Answer: [Yes]

561 Justification: All the theories and hypotheses we proposed are supported by experimental  
562 and mathematical derivations.

563 Guidelines:

- 564 • The answer NA means that the paper does not include theoretical results.
- 565 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
566 referenced.
- 567 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 568 • The proofs can either appear in the main paper or the supplemental material, but if  
569 they appear in the supplemental material, the authors are encouraged to provide a short  
570 proof sketch to provide intuition.
- 571 • Inversely, any informal proof provided in the core of the paper should be complemented  
572 by formal proofs provided in appendix or supplemental material.
- 573 • Theorems and Lemmas that the proof relies upon should be properly referenced.

574 **4. Experimental result reproducibility**

575 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
576 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
577 of the paper (regardless of whether the code and data are provided or not)?

578 Answer: [Yes]

579 Justification: We provide detailed descriptions of the hyperparameters in the paper and  
580 appendices, along with an anonymous link to the experimental demo in the abstract.

581 Guidelines:

- 582 • The answer NA means that the paper does not include experiments.

- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 613 5. Open access to data and code

614 Question: Does the paper provide open access to the data and code, with sufficient instructions  
 615 to faithfully reproduce the main experimental results, as described in supplemental  
 616 material?

617 Answer: [Yes]

618 Justification: We have included a link to an anonymous demo of our experiments in the  
 619 abstract.

620 Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

- 638           • Providing as much information as possible in supplemental material (appended to the  
639           paper) is recommended, but including URLs to data and code is permitted.

640           **6. Experimental setting/details**

641           Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
642           parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
643           results?

644           Answer: [Yes]

645           Justification: We provide the experimental setup details in both the main text and appendices.

646           Guidelines:

- 647           • The answer NA means that the paper does not include experiments.  
648           • The experimental setting should be presented in the core of the paper to a level of detail  
649           that is necessary to appreciate the results and make sense of them.  
650           • The full details can be provided either with the code, in appendix, or as supplemental  
651           material.

652           **7. Experiment statistical significance**

653           Question: Does the paper report error bars suitably and correctly defined or other appropriate  
654           information about the statistical significance of the experiments?

655           Answer: [Yes]

656           Justification: The margin of error is reported in the appendix.

657           Guidelines:

- 658           • The answer NA means that the paper does not include experiments.  
659           • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
660           dence intervals, or statistical significance tests, at least for the experiments that support  
661           the main claims of the paper.  
662           • The factors of variability that the error bars are capturing should be clearly stated (for  
663           example, train/test split, initialization, random drawing of some parameter, or overall  
664           run with given experimental conditions).  
665           • The method for calculating the error bars should be explained (closed form formula,  
666           call to a library function, bootstrap, etc.)  
667           • The assumptions made should be given (e.g., Normally distributed errors).  
668           • It should be clear whether the error bar is the standard deviation or the standard error  
669           of the mean.  
670           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
671           preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
672           of Normality of errors is not verified.  
673           • For asymmetric distributions, the authors should be careful not to show in tables or  
674           figures symmetric error bars that would yield results that are out of range (e.g. negative  
675           error rates).  
676           • If error bars are reported in tables or plots, The authors should explain in the text how  
677           they were calculated and reference the corresponding figures or tables in the text.

678           **8. Experiments compute resources**

679           Question: For each experiment, does the paper provide sufficient information on the com-  
680           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
681           the experiments?

682           Answer: [Yes]

683           Justification: We provide sufficient computational resource details for each experiment in  
684           both the main text and appendices.

685           Guidelines:

- 686           • The answer NA means that the paper does not include experiments.  
687           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
688           or cloud provider, including relevant memory and storage.

- 689           • The paper should provide the amount of compute required for each of the individual  
690           experimental runs as well as estimate the total compute.  
691           • The paper should disclose whether the full research project required more compute  
692           than the experiments reported in the paper (e.g., preliminary or failed experiments that  
693           didn't make it into the paper).

694           **9. Code of ethics**

695           Question: Does the research conducted in the paper conform, in every respect, with the  
696           NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

697           Answer: [Yes]

698           Justification: Our methodology and implementation fully adhere to the ethical code standards  
699           set forth by NeurIPS.

700           Guidelines:

- 701           • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
702           • If the authors answer No, they should explain the special circumstances that require a  
703           deviation from the Code of Ethics.  
704           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
705           eration due to laws or regulations in their jurisdiction).

706           **10. Broader impacts**

707           Question: Does the paper discuss both potential positive societal impacts and negative  
708           societal impacts of the work performed?

709           Answer: [NA]

710           Justification: We have discussed the broader impact of time series forecasting in both abstract  
711           and introduction.

712           Guidelines:

- 713           • The answer NA means that there is no societal impact of the work performed.  
714           • If the authors answer NA or No, they should explain why their work has no societal  
715           impact or why the paper does not address societal impact.  
716           • Examples of negative societal impacts include potential malicious or unintended uses  
717           (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
718           (e.g., deployment of technologies that could make decisions that unfairly impact specific  
719           groups), privacy considerations, and security considerations.  
720           • The conference expects that many papers will be foundational research and not tied  
721           to particular applications, let alone deployments. However, if there is a direct path to  
722           any negative applications, the authors should point it out. For example, it is legitimate  
723           to point out that an improvement in the quality of generative models could be used to  
724           generate deepfakes for disinformation. On the other hand, it is not needed to point out  
725           that a generic algorithm for optimizing neural networks could enable people to train  
726           models that generate Deepfakes faster.  
727           • The authors should consider possible harms that could arise when the technology is  
728           being used as intended and functioning correctly, harms that could arise when the  
729           technology is being used as intended but gives incorrect results, and harms following  
730           from (intentional or unintentional) misuse of the technology.  
731           • If there are negative societal impacts, the authors could also discuss possible mitigation  
732           strategies (e.g., gated release of models, providing defenses in addition to attacks,  
733           mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
734           feedback over time, improving the efficiency and accessibility of ML).

735           **11. Safeguards**

736           Question: Does the paper describe safeguards that have been put in place for responsible  
737           release of data or models that have a high risk for misuse (e.g., pretrained language models,  
738           image generators, or scraped datasets)?

739           Answer: [No]

740           Justification: This paper does not have this risk.

741 Guidelines:

- 742 • The answer NA means that the paper poses no such risks.
- 743 • Released models that have a high risk for misuse or dual-use should be released with
- 744 necessary safeguards to allow for controlled use of the model, for example by requiring
- 745 that users adhere to usage guidelines or restrictions to access the model or implementing
- 746 safety filters.
- 747 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 748 should describe how they avoided releasing unsafe images.
- 749 • We recognize that providing effective safeguards is challenging, and many papers do
- 750 not require this, but we encourage authors to take this into account and make a best
- 751 faith effort.

752 **12. Licenses for existing assets**

753 Question: Are the creators or original owners of assets (e.g., code, data, models), used in

754 the paper, properly credited and are the license and terms of use explicitly mentioned and

755 properly respected?

756 Answer: [Yes]

757 Justification: We included it in implementation details and appendix.

758 Guidelines:

- 759 • The answer NA means that the paper does not use existing assets.
- 760 • The authors should cite the original paper that produced the code package or dataset.
- 761 • The authors should state which version of the asset is used and, if possible, include a
- 762 URL.
- 763 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 764 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 765 service of that source should be provided.
- 766 • If assets are released, the license, copyright information, and terms of use in the
- 767 package should be provided. For popular datasets, [paperswithcode.com/datasets](http://paperswithcode.com/datasets)
- 768 has curated licenses for some datasets. Their licensing guide can help determine the
- 769 license of a dataset.
- 770 • For existing datasets that are re-packaged, both the original license and the license of
- 771 the derived asset (if it has changed) should be provided.
- 772 • If this information is not available online, the authors are encouraged to reach out to
- 773 the asset's creators.

774 **13. New assets**

775 Question: Are new assets introduced in the paper well documented and is the documentation

776 provided alongside the assets?

777 Answer: [NA]

778 Justification: N/A.

779 Guidelines:

- 780 • The answer NA means that the paper does not release new assets.
- 781 • Researchers should communicate the details of the dataset/code/model as part of their
- 782 submissions via structured templates. This includes details about training, license,
- 783 limitations, etc.
- 784 • The paper should discuss whether and how consent was obtained from people whose
- 785 asset is used.
- 786 • At submission time, remember to anonymize your assets (if applicable). You can either
- 787 create an anonymized URL or include an anonymized zip file.

788 **14. Crowdsourcing and research with human subjects**

789 Question: For crowdsourcing experiments and research with human subjects, does the paper

790 include the full text of instructions given to participants and screenshots, if applicable, as

791 well as details about compensation (if any)?

792                  Answer: [NA]

793                  Justification: N/A.

794                  Guidelines:

- 795                  • The answer NA means that the paper does not involve crowdsourcing nor research with  
796                  human subjects.
- 797                  • Including this information in the supplemental material is fine, but if the main contribu-  
798                  tion of the paper involves human subjects, then as much detail as possible should be  
799                  included in the main paper.
- 800                  • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
801                  or other labor should be paid at least the minimum wage in the country of the data  
802                  collector.

803                  **15. Institutional review board (IRB) approvals or equivalent for research with human  
804                  subjects**

805                  Question: Does the paper describe potential risks incurred by study participants, whether  
806                  such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
807                  approvals (or an equivalent approval/review based on the requirements of your country or  
808                  institution) were obtained?

809                  Answer: [NA]

810                  Justification: N/A.

811                  Guidelines:

- 812                  • The answer NA means that the paper does not involve crowdsourcing nor research with  
813                  human subjects.
- 814                  • Depending on the country in which research is conducted, IRB approval (or equivalent)  
815                  may be required for any human subjects research. If you obtained IRB approval, you  
816                  should clearly state this in the paper.
- 817                  • We recognize that the procedures for this may vary significantly between institutions  
818                  and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
819                  guidelines for their institution.
- 820                  • For initial submissions, do not include any information that would break anonymity (if  
821                  applicable), such as the institution conducting the review.

822                  **16. Declaration of LLM usage**

823                  Question: Does the paper describe the usage of LLMs if it is an important, original, or  
824                  non-standard component of the core methods in this research? Note that if the LLM is used  
825                  only for writing, editing, or formatting purposes and does not impact the core methodology,  
826                  scientific rigorosity, or originality of the research, declaration is not required.

827                  Answer: [No]

828                  Justification: We did not use any large language models (LLMs) in this work.

829                  Guidelines:

- 830                  • The answer NA means that the core method development in this research does not  
831                  involve LLMs as any important, original, or non-standard components.
- 832                  • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
833                  for what should or should not be described.