

# RI-Loss: A Learnable Residual-Informed Loss for Time Series Forecasting

JiETING WANG<sup>1,2\*</sup>, XIAOLEI SHANG<sup>1</sup>, FEIJIANG LI<sup>1,2</sup>, FURONG PENG<sup>1,2</sup>

<sup>1</sup>Institute of Big Data Science and Industry, Shanxi University

<sup>2</sup>Key Laboratory of Evolutionary Science Intelligence of Shanxi Province, Taiyuan, Shanxi, China  
jtwang@sxu.edu.cn

## Abstract

Time series forecasting relies on predicting future values from historical data, yet most state-of-the-art approaches—including transformer and multilayer perceptron-based models—optimize using Mean Squared Error (MSE), which has two fundamental weaknesses: its point-wise error computation fails to capture temporal relationships, and it does not account for inherent noise in the data. To overcome these limitations, we introduce the Residual-Informed Loss (RI-Loss), a novel objective function based on the Hilbert-Schmidt Independence Criterion (HSIC). RI-Loss explicitly models noise structure by enforcing dependence between the residual sequence and a random time series, enabling more robust, noise-aware representations. Theoretically, we derive the first non-asymptotic HSIC bound with explicit double-sample complexity terms, achieving optimal convergence rates through Bernstein-type concentration inequalities and Rademacher complexity analysis. This provides rigorous guarantees for RI-Loss optimization while precisely quantifying kernel space interactions. Empirically, experiments across eight real-world benchmarks and five leading forecasting models demonstrate improvements in predictive performance, validating the effectiveness of our approach.

**Code** — <https://github.com/shang-xl/RI-Loss>

**Appendix** — <https://arxiv.org/abs/2511.10130>

## Introduction

Time series data consist of chronologically ordered observations, where forecasting involves identifying latent patterns and trends from historical data to predict future values. In practice, temporal datasets typically demonstrate considerable complexity, manifesting as nonlinear relationships, non-stationary distributions, long-term dependencies, and significant noise contamination. Deep learning architectures excel at modeling these intricate temporal structures due to their exceptional capacity for nonlinear function approximation, allowing for more precise characterization of dynamic temporal patterns (Cao et al. 2018; Liu et al. 2022; Deng and Hooi 2021). These advanced modeling capabilities have enabled transformative applications across multiple domains, including financial market prediction, indus-

trial process monitoring, disease spread forecasting, and meteorological condition estimation.

Recent advances in time series analysis have witnessed the emergence of transformer-based and multilayer perceptron (MLP)-based architectures as dominant approaches, with numerous studies demonstrating their state-of-the-art performance (Wu et al. 2021; Zeng et al. 2022; Liu et al. 2024). However, despite their widespread adoption of MSE as the standard loss function, this conventional choice presents two critical limitations for long-term forecasting tasks.

The first limitation of MSE lies in its point-wise error computation, which makes it susceptible to overfitting observational noise, thereby impairing model generalization. More critically, MSE’s formulation disregards the temporal dynamics of sequential data, neglecting essential interdependencies across time steps. These inherent weaknesses hinder the model’s capacity to capture and preserve long-range dependencies effectively.

To better capture temporal dependencies, we propose a novel loss function inspired by a fundamental principle from (Yi and Wang 2019): an ideal model should produce residuals that are statistically indistinguishable from random noise, indicating complete extraction of all predictable patterns. While directly quantifying the difference between residuals and random sequences is problematic due to noise’s inherent stochasticity, we reformulate this intuition through statistical independence measures. Specifically, our approach maximizes the dependence between model residuals and random noise sequences, thereby compelling the model to extract more informative temporal features from the input data.

Building on these insights, we propose a novel forecasting framework that incorporates kernel-based nonparametric independence testing via a residual-informed loss function. Our approach centers on the Hilbert-Schmidt Independence Criterion (HSIC), which provides a theoretically sound measure of statistical dependence while remaining computationally efficient. The framework achieves two key innovations: it effectively captures nonlinear temporal dependencies between model residuals and random noise through kernel methods, and it establishes a new dependency-aware modeling paradigm that unifies statistical learning principles with deep neural forecasting architectures.

Theoretically, we establish rigorous learning guarantees

by analyzing the self-bounding properties of HSIC and deriving novel generalization bounds through double-sample Rademacher complexity analysis. These theoretical contributions achieve dual significance: they formally validate the learnability of HSIC-Loss, and they provide a generalizable analytical framework for assessing the generalization capacity of dependence-aware loss functions. Our comprehensive empirical evaluation demonstrates consistent performance gains across diverse real-world datasets.

The main contributions of this paper are as follows:

- **Methodological Innovation:** We propose RI-Loss, a novel residual-informed loss function that simultaneously addresses two fundamental limitations of MSE: susceptibility to observational noise and failure to capture temporal dependencies, establishing a new kernel-based learning paradigm for time series forecasting.
- **Theoretical Advancement:** We develop the first second-order Rademacher complexity bounds for the Hilbert-Schmidt Independence Criterion (HSIC) in time series analysis, providing rigorous generalization guarantees and creating a theoretical framework for analyzing dependence-aware loss functions.
- **Empirical Validation:** Through extensive experiments on eight diverse real-world benchmarks, we demonstrate that RI-Loss consistently enhances performance across state-of-the-art architectures (including Transformers and MLPs).

Complete proofs and additional experimental results are provided in the Appendix.

## Related Work

### Time Series Forecasting

Time series forecasting has undergone significant evolution, progressing from simple moving averages to classical statistical methods like Autoregressive Integrated Moving Average (ARIMA) (Box et al. 1970), through machine learning approaches such as Support Vector Regression (SVR) (Smola and Scholkopf 2004), to contemporary deep learning architectures. The deep learning revolution has introduced specialized models for temporal data processing: Recurrent Neural Networks (RNNs) effectively capture sequential dependencies through iterative time-step processing (Cao et al. 2018); Convolutional Neural Networks (CNNs) extract localized temporal patterns via convolutional kernels (Hewage et al. 2020; Liu et al. 2022); Graph Neural Networks (GNNs) model complex multivariate interactions (Deng and Hooi 2021; Cao et al. 2020); Transformers utilize self-attention mechanisms to identify both long-range dependencies and cross-variable correlations (Zhou et al. 2021; Wu et al. 2021; Liu et al. 2024); while Multilayer Perceptrons (MLPs) provide computationally efficient solutions for resource-constrained scenarios (Zeng et al. 2022; Xu, Zeng, and Xu 2024). Despite these architectural advances, most deep time series models continue to rely on Mean Squared Error (MSE) loss for parameter optimization, presenting limitations in handling noise and temporal relationships.

Recent work has proposed MSE alternatives via shape alignment (Cuturi and Blondel 2017; Le Guen and Thome 2019) and dependency modeling (Wang et al. 2025), but these face computational complexity or residual pattern neglect. Our RI-Loss addresses both issues through statistical independence, offering noise-robust temporal modeling via kernel dependency measurement, and provable learning guarantees through HSIC-based theory.

### HSIC-based Learning Method

The Hilbert-Schmidt Independence Criterion (HSIC) is a powerful kernel-based independence measure widely used in machine learning. In representation learning, it enhances latent variable interpretability through deep generative model regularization (Li et al. 2021). For causal inference, HSIC enables nonlinear conditional independence testing, overcoming linearity constraints in causal discovery (Li et al. 2024; Hu, Sejdinovic, and Evans 2024). Transfer learning benefits from its synergistic use with maximum mean discrepancy for cross-domain feature alignment (Wang et al. 2020), while time series analysis leverages HSIC for nonlinear Granger causality detection in financial modeling (Ren, Li, and Han 2020). Despite these successful applications, HSIC’s potential for time series forecasting remains largely untapped, offering significant opportunities to advance temporal dependency modeling.

## Preliminaries

### Time Series Forecasting Task

Time series forecasting aims to predict future values based on historical observations. In our framework, given a look-back window  $\mathbf{X}_t = (X_t, X_{t-1}, \dots, X_{t-w+1})^\top \in \mathbb{R}^w$  of  $w$  historical observations, we predict the  $H$ -step future trajectory  $\mathbf{Y} = (Y_{t+1}, \dots, Y_{t+H})^\top \in \mathbb{R}^H$  through  $\hat{\mathbf{Y}} = f(\mathbf{X}_t)$ , where  $f : \mathbb{R}^w \rightarrow \mathbb{R}^H$  is the forecasting function. For multivariate series with  $d$ -dimensional observations  $\mathbf{x}_t \in \mathbb{R}^d$ , this extends to predicting  $\mathbf{Y} = \mathbf{y}_{t+1}, \dots, \mathbf{y}_{t+H} \in \mathbb{R}^{H \times d}$  from input  $\mathbf{X}_t = \mathbf{x}_t, \mathbf{x}_{t-1}, \dots, \mathbf{x}_{t-w+1} \in \mathbb{R}^{w \times d}$ .

Generally, the forecasting quality is evaluated via the MSE:

$$\text{MSE} = \frac{1}{H} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2 = \frac{1}{H} \sum_{k=1}^H \|\mathbf{y}_{t+k} - \hat{\mathbf{y}}_{t+k}\|_2^2, \quad (1)$$

where  $\|\cdot\|_F^2$  represents the square of the Frobenius norm for matrices. The optimal solution of MSE corresponds to the conditional expectation of the true data, which means minimizing MSE guarantees that the mean of the predicted sequence matches that of the true sequence.

### Additive Noise Model for Time Series Observations

Conventional time series prediction frameworks typically adopt an additive model structure of the form:  $\mathbf{Y} = h(\mathbf{X}_t) + \boldsymbol{\epsilon}$ , where  $\mathbf{Y}$  denotes the observation sequence,  $h : \mathbb{R}^w \rightarrow \mathbb{R}^H$  represents the true underlying system dynamics that maps historical patterns to future values, and  $\boldsymbol{\epsilon}$  is a zero-mean additive noise vector satisfying  $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$ . The noise process  $\boldsymbol{\epsilon} = (\epsilon_{t+1}, \dots, \epsilon_{t+H})^\top$  further exhibits two key

properties: (i) temporal uncorrelation with  $\mathbb{E}[\epsilon_{t+i}\epsilon_{t+j}] = 0$  for all  $i, j = 1, \dots, H$  where  $i \neq j$ , and (ii) signal independence with  $\epsilon_{t+i} \perp h(\mathbf{X}_t)$  for all prediction horizons  $i$ .

Although we define the loss within the above basic framework, it naturally generalizes to more complex scenarios, including heteroscedastic noise structures  $\epsilon_{t+i} \sim \mathcal{N}(0, \sigma_{t+i}^2)$  that accommodate time-varying volatility patterns commonly observed in financial and environmental time series, as well as state-dependent noise models of the form  $Y_{t+i} = h(\mathbf{X}_t) + g(\mathbf{X}_t)\epsilon_{t+i}$  where the noise amplitude  $g(\mathbf{X}_t)$  varies with the system state.

### Hilbert-Schmidt Independence Criterion

The Hilbert-Schmidt Independence Criterion (HSIC) provides a kernel-based measure of statistical dependence between random variables. Formally, for variables  $R \in \mathcal{R}$  and  $S \in \mathcal{S}$  with joint distribution  $\mathbb{P}_{R,S}$ , using reproducing kernel Hilbert spaces (RKHS)  $\mathcal{F}$  on  $\mathcal{R}$  and  $\mathcal{G}$  on  $\mathcal{S}$ , the HSIC is defined as the squared Hilbert-Schmidt norm of their cross-covariance operator.

As established by (Gretton et al. 2005; Greenfeld and Shalit 2020), the HSIC admits an important probabilistic interpretation: it represents the maximum achievable covariance between normalized functions in the respective RKHS. Mathematically, for  $f \in \mathcal{F}$  and  $g \in \mathcal{G}$ , this is expressed as:

$$\mathbb{E}(\text{HSIC}(R, S)) = \sup_{\|f\|_{\mathcal{F}} \leq 1, \|g\|_{\mathcal{G}} \leq 1} (\text{COV}[f(R), g(S)])^2, \quad (2)$$

where COV denotes covariance. The HSIC possesses two key properties:  $\text{HSIC}(R, S) = 0$  if and only if  $R$  and  $S$  are independent, and the magnitude of HSIC reflects the strength of dependence. This formulation makes HSIC particularly powerful for detecting nonlinear dependencies that conventional correlation measures might miss, as it operates in high-dimensional feature spaces through the kernel trick.

By expanding the Hilbert-Schmidt norm through kernel-based inner products, we derive the following equivalent formulation of HSIC:

**Definition 1** (Population HSIC). *For random variables  $R \in \mathcal{R}$  and  $S \in \mathcal{S}$  with joint distribution  $\mathbb{P}_{R,S}$ , let  $k : \mathcal{R} \times \mathcal{R} \rightarrow \mathbb{R}$  and  $l : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$  be characteristic kernels, and let  $(R', S')$  be an independent copy of  $(R, S)$ . The Hilbert-Schmidt Independence Criterion admits the following equivalent expression:*

$$\begin{aligned} \text{HSIC}(R, S) &= \mathbb{E}_{RR'SS'} [k(R, R')l(S, S')] \\ &+ \mathbb{E}_{RR'} [k(R, R')] \mathbb{E}_{SS'} [l(S, S')] \\ &- 2\mathbb{E}_{RS} [\mathbb{E}_{R'} [k(R, R')] \mathbb{E}_{S'} [l(S, S')]]. \end{aligned} \quad (3)$$

The HSIC statistic has range  $[0, \infty)$ . This formulation reveals that HSIC measures the discrepancy between the joint embedding and the product of marginal embeddings.

The kernel-based approach provides three key advantages: it is nonparametric (making no assumptions about the dependence structure), flexible (accommodating mixed data types through kernel choice), and universal (detecting any measurable dependence when using characteristic kernels).

We now present the finite-sample estimator of HSIC, which operationalizes the population measure for practical computation.

**Definition 2** (Empirical HSIC Estimator). *Given an i.i.d. sample  $\{(r_i, s_i)\}_{i=1}^n$ , the U-statistic estimator of HSIC is given by:*

$$\begin{aligned} \text{HSIC} &= \frac{1}{\binom{n}{2}} \sum_{i < j} k(r_i, r_j) l(s_i, s_j) \\ &+ \frac{1}{\binom{n}{4}} \sum_{i < j < k < l} k(r_i, r_j) l(s_k, s_l) \\ &- \frac{2}{\binom{n}{3}} \sum_{i < j < k} [k(r_i, r_j) l(s_j, s_k) + k(r_j, s_k) l(s_i, s_j)], \end{aligned} \quad (4)$$

where all sums are over distinct index tuples, the binomial coefficient  $\binom{n}{m}$  represents the number of ways to choose  $m$  elements from a set of  $n$  distinct items, and  $k, l$  are characteristic kernels.

### Residual-Informed Loss

#### Decomposition of MSE Loss for Additive Model

Under the additive noise model described earlier, where the observed outputs  $\mathbf{Y}$  are corrupted by additive noise  $\epsilon$  with mean zero such that  $\mathbf{Y} = h(\mathbf{X}_t) + \epsilon$ , we can analyze the relationship between the true MSE and the observed MSE. The true MSE measures the prediction error with respect to the noise-free target  $h(\mathbf{X}_t)$ :  $\text{MSE}_{\text{true}} = \frac{1}{H} \|h(\mathbf{X}_t) - \hat{\mathbf{Y}}\|_F^2$ . The observed MSE quantifies the discrepancy between predictions and noisy measurements:  $\text{MSE}_{\text{obs}} = \frac{1}{H} \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2$ .

Using the observation model  $h(\mathbf{X}_t) = \mathbf{Y} - \epsilon$ , we can express the true MSE as:

$$\begin{aligned} \text{MSE}_{\text{true}} &= \frac{1}{H} \|(\mathbf{Y} - \epsilon) - \hat{\mathbf{Y}}\|_F^2 \\ &= \frac{1}{H} \left( \|\mathbf{Y} - \hat{\mathbf{Y}}\|_F^2 - 2\langle \mathbf{Y} - \hat{\mathbf{Y}}, \epsilon \rangle + \|\epsilon\|_F^2 \right) \\ &= \text{MSE}_{\text{obs}} - \frac{2}{H} \langle \mathbf{Y} - \hat{\mathbf{Y}}, \epsilon \rangle + \frac{1}{H} \|\epsilon\|_F^2. \end{aligned} \quad (5)$$

Assuming the noise  $\epsilon$  is zero-mean and uncorrelated with the prediction error, taking expectations on both sides yields:

$$\mathbb{E}[\text{MSE}_{\text{true}}] = \mathbb{E}[\text{MSE}_{\text{obs}}] - \frac{1}{H} \mathbb{E}[\|\epsilon\|_F^2], \quad (6)$$

where  $\mathbb{E}[\|\epsilon\|_F^2]$  denotes the noise variance like term. This shows that the observed MSE systematically overestimates the true prediction error by exactly the noise power. However, in general cases, since the prediction  $\hat{\mathbf{Y}}$  is learned from noisy observations  $\mathbf{Y}$ , the prediction error  $\hat{\mathbf{Y}} - \mathbf{Y}$  depends on the noisy  $\epsilon$ , making the cross-term in Eq. (5) non-zero:

**Theorem 1** (Cross-Term Expectation for Linear Projection). *Let  $\mathbf{Y} = (Y_{t+1}, \dots, Y_{t+H})^\top \in \mathbb{R}^H$  be a noisy observation vector generated by:  $\mathbf{Y} = h(\mathbf{X}_t) + \epsilon$ , where  $h(\mathbf{X}_t) : \mathcal{X} \rightarrow \mathbb{R}^H$  is a deterministic mapping, and the noise vector  $\epsilon \in \mathbb{R}^H$  satisfies:  $\mathbb{E}[\epsilon] = \mathbf{0}$ ,  $\mathbb{E}[\epsilon\epsilon^\top] = \sigma^2 \mathbf{I}_H$ ,  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_H)$ . Here,  $\mathbf{I}_H$  denotes the  $H \times H$  identity matrix and  $\sigma^2 > 0$*

is the noise variance. For any linear estimator  $\hat{\mathbf{Y}} = \mathbf{P}\mathbf{Y}$  with projection matrix  $\mathbf{P} \in \mathbb{R}^{H \times H}$ , the expected normalized cross-term evaluates to:

$$\mathbb{E} \left[ \frac{1}{H} \langle \mathbf{Y} - \hat{\mathbf{Y}}, \boldsymbol{\epsilon} \rangle \right] = \frac{\sigma^2}{H} \text{tr}(\mathbf{I}_H - \mathbf{P}), \quad (7)$$

where  $\text{tr}(\cdot)$  represents the trace of a matrix, defined as the sum of its diagonal elements.

The decomposition of MSE in Eq.(5) and Theorem 1 motivate our design of a loss function that explicitly accounts for the dependency between prediction residuals and inherent noise. To explicitly penalize both prediction accuracy and residual-noise dependency, we first propose a naive loss combining the empirical MSE with a linear residual-noise alignment term:  $\mathcal{L}_{\text{naive}} = \text{MSE}_{\text{obs}} - \lambda \langle \mathbf{Y} - \hat{\mathbf{Y}}, \boldsymbol{\epsilon} \rangle$ , where  $\lambda$  is the trade-off parameter. However, this linear term only captures first-order dependencies.

### Definition of Residual-Informed Loss

To generalize  $\langle \mathbf{Y} - \hat{\mathbf{Y}}, \boldsymbol{\epsilon} \rangle$  to nonlinear relationships, we replace the inner product with the HSIC. We analyze their relationship. If  $R$  and  $S$  are centered random variables, that is,  $\mathbb{E}[R] = 0$  and  $\mathbb{E}[S] = 0$ , then their covariance simplifies to  $\text{COV}(R, S) = \mathbb{E}[RS]$ . At this time, we can observe that  $\mathbb{E}[RS] \leq \mathbb{E}(\text{HSIC}(R, S))$ . This inequality holds because if  $\mathcal{F}$  and  $\mathcal{G}$  contain the identity functions (i.e.,  $f(R) = R \in \mathcal{F}$ ,  $g(S) = S \in \mathcal{G}$ ), then  $\mathbb{E}(\text{HSIC}(R, S)) = \sup_{f \in \mathcal{F}, g \in \mathcal{G}} \text{COV}[f(R), g(S)] \geq \text{COV}[R, S]$ . Thus, assuming the residual expectation is zero ( $\mathbb{E}[\mathbf{Y} - \hat{\mathbf{Y}}] = \mathbf{0}$ ) and noise expectation is zero ( $\mathbb{E}[\boldsymbol{\epsilon}] = \mathbf{0}$ ), we have  $\mathbb{E}(\langle \mathbf{Y} - \hat{\mathbf{Y}}, \boldsymbol{\epsilon} \rangle) \leq \mathbb{E}(\text{HSIC}(\mathbf{Y} - \hat{\mathbf{Y}}, \boldsymbol{\epsilon}))$ . This yields the Residual-Informed (RI) loss:

**Definition 3.** For the observed output series  $\mathbf{Y}$ , the predicted series  $\hat{\mathbf{Y}}$  and a noise series  $\boldsymbol{\epsilon}$ , RI loss is defined as:

$$\text{RI}(\mathbf{Y}, \hat{\mathbf{Y}}) = \text{MSE}_{\text{obs}} + \lambda \exp(-\tau \cdot \text{HSIC}(\mathbf{Y} - \hat{\mathbf{Y}}, \boldsymbol{\epsilon})), \quad (8)$$

where the exponential transformation is applied to ensure that the loss function remains positive-valued and  $\lambda, \tau$  are the trade-off parameter and transformation parameter, respectively.

RI-Loss in Eq. (8) requires the model to preserve interpretable signals by minimizing the MSE to ensure prediction accuracy, while explicitly separating the noise structure by maximizing the dependence between the residuals and the noise, thereby forcing the model to push unexplainable variation (noise structure) into the residuals. Mathematically, when the residual  $\mathbf{Y} - \hat{\mathbf{Y}}$  is highly correlated with the noise  $\boldsymbol{\epsilon}$ , it indicates that the model has excluded as much of the noise-distributed components as possible from the prediction  $\hat{\mathbf{Y}}$ , ensuring that the prediction result  $\hat{\mathbf{Y}}$  contains only the true signal that is independent of the noise.

### RI Advantage

We investigate how RI loss varies with increasing noise ratios  $\rho \in [0, 1]$ , comparing its behavior to conventional MSE.

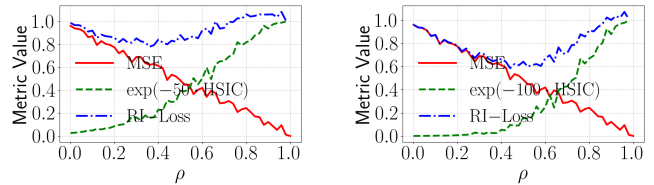


Figure 1: The RI-Loss varies with the noise ratio.

The results is demonstrated in Figure 1. Our experiment generates a sinusoidal signal with noise  $\mathbf{y}_{\text{true}} = \sin(\mathbf{x}) + \boldsymbol{\epsilon}$  where  $\mathbf{x} \in [0, 2\pi]$  contains 500 points and  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, 1)$  represents baseline Gaussian noise. Crucially, for each  $\rho$  value (51 linearly spaced steps), we corrupt the underlying noise-free sinusoidal component by adding  $\mathcal{N}(0, 1)$  noise to randomly selected  $\rho \times 100\%$  of the  $\sin(\mathbf{x})$  values, producing  $\mathbf{y}_{\text{noisy}}$ . The metrics include  $\text{MSE} = \mathbb{E}[(\mathbf{y}_{\text{true}} - \mathbf{y}_{\text{noisy}})^2]$  and  $\text{RI} = \text{MSE} + \exp(-\tau \cdot \text{HSIC}(\mathbf{y}_{\Delta}, \boldsymbol{\epsilon}))$  ( $\tau = 50, 100, \lambda = 1$ ), where  $\mathbf{y}_{\Delta} = \mathbf{y}_{\text{true}} - \mathbf{y}_{\text{noisy}}$ . For HSIC computation, we employ Gaussian kernel matrices  $k_{ij} = \exp(-|\mathbf{y}_{\Delta,i} - \mathbf{y}_{\Delta,j}|^2)$  for residual and  $l_{ij} = \exp(-|\boldsymbol{\epsilon}_i - \boldsymbol{\epsilon}_j|^2)$  for noise.

As shown in Figure 1, the RI-Loss demonstrates a convex dependence on the noise ratio  $\rho$ , reflecting the fundamental trade-off between signal reconstruction and noise suppression. For about  $\rho < 0.5$  in the left panel ( $\rho < 0.6$  in the right panel), the system prioritizes signal fitting, where controlled noise retention improves reconstruction fidelity (manifested in decreasing MSE). When about  $\rho > 0.5$  in the left panel ( $\rho > 0.6$  in the right panel), the exponential HSIC penalty ( $\exp(-\tau \cdot \text{HSIC}) \rightarrow 1$ ) becomes dominant, enforcing noise rejection at the expense of elevated RI-Loss. The characteristic minimum at  $\rho \approx 0.5$  in the left panel ( $\rho \approx 0.6$  in the right panel) naturally emerges as the optimal operating point that balances these competing requirements.

### Theoretical Analysis of RI

The major purpose of this paper is presenting a novel approach that enhances standard MSE optimization by incorporating a HSIC regularization term between model residuals and random noise. Hence, we establish a theoretical framework to investigate HSIC's generalization properties.

The current theoretical analysis operates under the baseline assumption of independent and identically distributed (i.i.d.) residuals and noise terms. This initial configuration provides crucial insights into HSIC's generalization properties. In future work, we will investigate the generalization properties of both HSIC and MSE under more complex dependency structures.

The generalization gap measures how well empirical training error approximates the true expected error (Bartlett and Mendelson 2003). A smaller gap indicates better generalization performance from training data to the underlying distribution. Theoretical bounds on this gap ensure reliable model behavior on unseen data. Our theorem shows that the empirical HSIC converges to the population HSIC as  $n \rightarrow \infty$ , providing generalization guarantees.

**Theorem 2.** Let  $\sigma_1, \dots, \sigma_n$  are independent Rademacher random

variables ( $\mathbb{P}(\sigma_i = \pm 1) = \frac{1}{2}$ ), where  $n$  is the number of samples. For random variables  $R \in \mathcal{R}$  and  $S \in \mathcal{S}$  with joint distribution  $\mathbb{P}_{RS}$ , let  $k : \mathcal{R} \times \mathcal{R} \rightarrow \mathbb{R}$  and  $l : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$  be characteristic kernels. Let  $(R', S')$  be an independent and identically distributed (i.i.d.) copy of  $(R, S)$ . Suppose that  $c_1(k) \leq k \leq c_2(k)$  and  $c_1(l) \leq l \leq c_2(l)$ , for  $\delta > 0$ , then the following holds:

$$\begin{aligned} & |\text{HSIC}(\{(r_i, s_i)\}_{i=1}^n) - \mathbb{E}[\text{HSIC}(R, S)]| \\ & \leq 3 \sum_{j=k,l} c_2(j) \sum_{m=1}^3 \gamma_m(j). \end{aligned} \quad (9)$$

Specifically, the above key functions for  $f \in \{k, l\}$  are defined as:

$$\begin{aligned} & \gamma_1(n, \delta; f) \\ & = \frac{10(c_2(f) - c_1(f))}{n} \ln \frac{2}{\delta} \\ & \gamma_2(n, \delta, R_\sigma; f) \\ & = 2(c_2(f) - c_1(f)) \sqrt{2 \ln \frac{2}{\delta} \left( \frac{R_\sigma}{2n} + \frac{1}{n^2} \ln \frac{2}{\delta} \right)} \\ & \gamma_3(n, \delta, W_{\sigma,\sigma}, W_{\sigma,\alpha}, W_\sigma, F; f) \\ & = 4 \left( \ln \frac{2}{\delta} \left( \frac{C_0}{n-1} (\mathbb{E}W_{\sigma,\sigma} + \sqrt{2}\mathbb{E}W_{\sigma,\alpha} \right. \right. \\ & \left. \left. + \frac{2(\mathbb{E}W_\sigma + F)}{n} + \frac{\sqrt{8}Fn^{1/2}}{n^2} + \frac{4F}{n^2} \right) + \frac{1}{n^2} \ln \frac{2}{\delta} \right)^{\frac{1}{2}}. \end{aligned} \quad (10)$$

In the above functions, the first-order and second-order sample Rademacher complexities are:

$$\mathcal{R}_\sigma = \mathbb{E}_\sigma \left[ \sup_{f \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i f_1(X_i) \right], \quad (11)$$

$$W_{\sigma,\sigma} = \sup_{f \in \mathcal{F}} \frac{1}{n^2} \left| \sum_{i,j} \sigma_i \sigma_j f_2(X_i, X_j) \right|, \quad (12)$$

$$W_{\sigma,\alpha} = \sup_{f \in \mathcal{F}} \sup_{\alpha: \|\alpha\|_2 \leq 1} \frac{1}{n^2} \sum_{i,j} \sigma_i \alpha_j f_2(X_i, X_j), \quad (13)$$

$$W_\sigma = \sup_{f \in \mathcal{F}, k=1, \dots, n} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i f_2(X_i, X_k) \right|, \quad (14)$$

$$F = \sup_{f \in \mathcal{F}} \|f_2\|_\infty, \quad (15)$$

where  $f_1$  and  $f_2$  respectively represent the centered kernels:

$$\begin{aligned} f_1(X_i) &= \mathbb{E}[f(X_i, X) | X_i] - \mathbb{E}[f(X_i, X_j)], \\ f_2(X_i, X_j) &= f(X_i, X_j) - f_1(X_i) - f_1(X_j) \\ & \quad - \mathbb{E}[f(X_i, X_j)]. \end{aligned} \quad (16)$$

Theorem 2 makes a significant breakthrough in the non-asymptotic analysis of HSIC statistics, establishing for the first time a convergence upper bound that explicitly incorporates double-sample Rademacher complexities. Through an innovative decomposition of the error into three terms with clear statistical significance: the  $\gamma_1$  term controls the deviation with typical  $\mathcal{O}(1/n)$  convergence rate ( $\mathcal{O}(\cdot)$  characterizes the asymptotic order); the  $\gamma_2$  term's square root structure reflects classical sub-exponential convergence properties; and the most innovative  $\gamma_3$  term employs sophisticated probability inequality techniques to constrain the influence of double-sample complexities (including  $W_{\sigma,\sigma}$  and  $W_{\sigma,\alpha}$ ) at the  $\mathcal{O}(1/\sqrt{n})$  level.

This theoretical achievement not only maintains optimal convergence rates but, more importantly, rigorously quantifies higher-order interaction effects in kernel function spaces in HSIC analysis, providing a new theoretical framework for high-dimensional nonparametric independence tests and offering important guidance for kernel method selection and sample complexity estimation. In the RI loss framework,  $r$  denotes the residual vector  $\mathbf{Y} - \hat{\mathbf{Y}}$ ,  $s$  denotes the noise vector  $\epsilon$ , and  $n$  is the sequence length. Theorem 2 establishes that the empirical HSIC( $r, s$ ) estimator exhibits finite-sample convergence with error decaying as  $\mathcal{O}(1/\sqrt{n})$ , dependent on kernel complexities  $\mathcal{R}_\sigma, W_{\sigma,\sigma}, W_{\sigma,\alpha}, W_\sigma$ , and achieves asymptotic consistency:  $\lim_{n \rightarrow \infty} \text{HSIC}_n(r, s) \stackrel{a.s.}{=} \mathbb{E}_{R,S}[\text{HSIC}(R, S)]$ . This result provides fundamental guarantees for RI loss optimization, ensuring both theoretical soundness and practical reliability in large-sample regimes.

## RI-based Time Series Forecasting

Our proposed RI-Loss demonstrates broad compatibility with existing deep time series architectures, as illustrated in Figure 2. The modular design enables seamless integration with both Transformer-based (e.g., Autoformer, Informer) and MLP-based (e.g., DLinear) forecasting frameworks. We employ a uniform distribution over the interval  $[-1, 1]$  as the noise distribution. Complete implementation details are provided in Algorithm 1 of Appendix B.

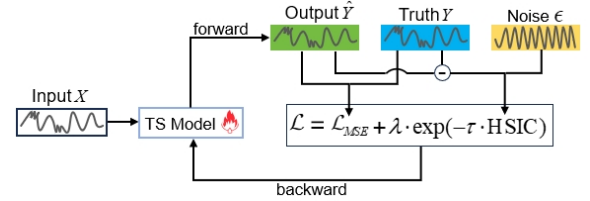


Figure 2: RI-Loss Based Time Series Model.

## Experimental Analysis

We conduct a comprehensive evaluation of the proposed RI-Loss. Additional experimental results, including the effects of the lookback window size, visualization of correlation structures, and hyperparameter sensitivity analysis, are provided in the appendix.

### Experiments Settings

**Datasets** We conduct extensive experiments on eight widely used real-world datasets, including ETT (ETTth1, ETTth2, ETTm1, ETTm2), Exchange, ILI, Weather and Electricity. Please refer to Appendix C for data descriptions.

**Backbones** We use five long-term time series forecasting models as our backbone models to comprehensively evaluate the effectiveness of RI-Loss. These include three Transformer-based models: Informer (Zhou et al. 2021), Autoformer (Wu et al. 2021), and iTransformer (Liu et al. 2024) and two MLP-based models: Decomposition-Linear (DLinear) (Zeng et al. 2022) and Retrieval Augmented Time Series Forecasting (RAFT) (Han et al. 2025).

Models		Informer(2021)		Autoformer(2021)		DLinear(2022)		iTransformer(2024)		RAFT(2025)	
Metric		MSE	Ours	MSE	Ours	MSE	Ours	MSE	Ours	MSE	Ours
ETTh1	96	0.878	<b>0.875</b>	0.456	<b>0.440</b>	0.384	<b>0.346</b>	0.387	<b>0.382</b>	0.387	<b>0.378</b>
	192	1.013	<b>0.970</b>	0.492	<b>0.474</b>	0.443	<b>0.404</b>	0.441	<b>0.437</b>	0.423	<b>0.422</b>
	336	<b>1.172</b>	1.225	<b>0.506</b>	0.522	0.447	<b>0.443</b>	0.491	<b>0.490</b>	<b>0.458</b>	0.460
	720	1.175	<b>1.169</b>	<b>0.496</b>	0.544	0.472	<b>0.471</b>	0.508	<b>0.490</b>	0.463	<b>0.462</b>
ETTh2	96	3.340	<b>2.299</b>	0.375	<b>0.336</b>	0.297	<b>0.279</b>	0.301	<b>0.296</b>	0.296	<b>0.293</b>
	192	5.684	<b>3.949</b>	0.451	<b>0.417</b>	0.373	<b>0.355</b>	0.380	<b>0.374</b>	0.384	<b>0.377</b>
	336	4.418	<b>3.412</b>	0.477	<b>0.444</b>	0.450	<b>0.412</b>	0.424	<b>0.419</b>	0.426	<b>0.415</b>
	720	3.249	<b>3.094</b>	0.476	<b>0.445</b>	0.696	<b>0.582</b>	0.430	<b>0.422</b>	0.434	<b>0.417</b>
ETTm1	96	0.678	<b>0.536</b>	0.478	<b>0.466</b>	0.302	<b>0.294</b>	0.342	<b>0.322</b>	0.329	<b>0.316</b>
	192	0.784	<b>0.610</b>	0.549	<b>0.534</b>	0.337	<b>0.333</b>	0.383	<b>0.376</b>	0.363	<b>0.357</b>
	336	1.011	<b>0.807</b>	0.516	<b>0.480</b>	0.371	<b>0.369</b>	0.418	<b>0.412</b>	0.391	<b>0.385</b>
	720	1.037	<b>0.978</b>	0.526	<b>0.523</b>	0.426	<b>0.423</b>	0.487	<b>0.478</b>	<b>0.444</b>	0.447
ETTm2	96	0.464	<b>0.343</b>	0.242	<b>0.216</b>	0.170	<b>0.164</b>	0.186	<b>0.175</b>	0.177	<b>0.174</b>
	192	0.811	<b>0.504</b>	0.300	<b>0.269</b>	0.226	<b>0.220</b>	0.254	<b>0.243</b>	0.243	<b>0.238</b>
	336	1.286	<b>1.124</b>	0.336	<b>0.325</b>	0.292	<b>0.276</b>	0.316	<b>0.306</b>	0.302	<b>0.297</b>
	720	3.949	<b>3.269</b>	0.432	<b>0.415</b>	0.406	<b>0.375</b>	0.414	<b>0.407</b>	0.402	<b>0.396</b>
Exchange	96	1.022	<b>0.901</b>	0.154	<b>0.138</b>	0.086	<b>0.084</b>	0.086	<b>0.085</b>	0.089	<b>0.084</b>
	192	1.204	<b>1.176</b>	0.297	<b>0.297</b>	0.199	<b>0.163</b>	0.181	<b>0.180</b>	0.192	<b>0.177</b>
	336	<b>1.577</b>	1.674	<b>0.457</b>	0.511	0.345	<b>0.265</b>	0.338	<b>0.337</b>	0.370	<b>0.343</b>
	720	<b>2.206</b>	2.574	<b>1.079</b>	1.188	<b>0.872</b>	1.128	<b>0.853</b>	0.863	1.133	<b>1.058</b>
ILI	24	6.158	<b>5.972</b>	3.427	<b>3.322</b>	2.205	<b>2.199</b>	<b>2.695</b>	2.739	2.040	<b>2.015</b>
	36	5.930	<b>5.739</b>	3.596	<b>3.425</b>	2.392	<b>2.106</b>	2.563	<b>2.531</b>	2.085	<b>1.962</b>
	48	<b>5.185</b>	5.237	3.493	<b>3.178</b>	2.298	<b>1.984</b>	2.567	<b>2.455</b>	2.069	<b>1.884</b>
	60	5.296	<b>5.091</b>	<b>2.847</b>	3.019	2.410	<b>1.970</b>	2.625	<b>2.480</b>	2.065	<b>1.859</b>
Weather	96	0.496	<b>0.479</b>	<b>0.224</b>	<b>0.224</b>	0.144	<b>0.143</b>	0.174	<b>0.169</b>	<b>0.189</b>	<b>0.189</b>
	192	0.582	<b>0.575</b>	0.305	<b>0.283</b>	0.188	<b>0.184</b>	0.224	<b>0.219</b>	0.239	<b>0.234</b>
	336	0.643	<b>0.609</b>	0.353	<b>0.347</b>	0.240	<b>0.234</b>	0.283	<b>0.278</b>	0.291	<b>0.283</b>
	720	0.651	<b>0.637</b>	0.456	<b>0.419</b>	0.317	<b>0.309</b>	0.359	<b>0.357</b>	0.366	<b>0.358</b>
Electricity	96	0.289	<b>0.281</b>	0.202	<b>0.197</b>	<b>0.140</b>	0.141	<b>0.148</b>	0.149	<b>0.174</b>	0.176
	192	0.296	<b>0.295</b>	0.218	<b>0.210</b>	<b>0.154</b>	0.155	0.166	<b>0.165</b>	0.172	<b>0.170</b>
	336	<b>0.297</b>	0.302	<b>0.231</b>	<b>0.231</b>	<b>0.169</b>	0.170	<b>0.178</b>	0.180	0.184	<b>0.181</b>
	720	<b>0.341</b>	0.367	0.257	<b>0.244</b>	<b>0.204</b>	<b>0.204</b>	<b>0.209</b>	<b>0.209</b>	<b>0.218</b>	0.220
1 <sup>st</sup> Count		6	26	7	27	5	28	7	26	5	28

Table 1: The table reports multivariate forecasting results on eight datasets and compares them with several baseline models. All experiments use the same prediction length settings as the corresponding backbone models. Ours indicates models trained with RI-Loss under identical architectures. The best results are shown in bold. The last row indicates the number of top-1 times of each model.

ID	-3dB		3dB		10dB	
	MSE	Ours	MSE	Ours	MSE	Ours
1	0.485	<b>0.473</b>	0.481	<b>0.457</b>	0.481	<b>0.456</b>
2	0.693	<b>0.579</b>	0.649	<b>0.558</b>	0.613	<b>0.551</b>
3	0.437	<b>0.433</b>	0.428	<b>0.423</b>	0.425	<b>0.422</b>
4	0.471	<b>0.391</b>	0.459	<b>0.380</b>	0.455	<b>0.375</b>
6	0.340	<b>0.329</b>	0.326	<b>0.319</b>	0.325	<b>0.320</b>

Table 2: MSE comparison of DLinear at 720 steps under different noise levels. The dataset IDs follow the same order as in Table 1.

**Evaluation Metrics** Two widely used evaluation metrics in long-term time series forecasting, MSE and Mean Absolute Error (MAE), are employed to validate the proposed loss function. Lower values indicate better performance.

Method	Dlinear	iTransformer	RAFT
MSE	16.3	32.5	1386
RI-Loss	30.5	40.2	1393

Table 3: Training time comparison (ms/iter) between MSE and RI-Loss at forecasting horizon 720 on the Weather dataset.

**Implementation Details** We conduct our experiments using the official implementations of each backbone model from their respective repositories. To ensure a fair comparison, we adopt the original experimental setups and hyperparameter configurations when incorporating RI-Loss to enhance the performance of the backbone models. The hyperparameter  $\lambda$  in RI-Loss is set to 10, and  $\tau$  is set to 1. The kernel function is implemented as a Gaussian kernel with a bandwidth of 1. To mitigate the impact of randomness, all

Models	Loss	ETTh2			
		96	192	336	720
Informer	RI-Loss	<b>2.299</b>	<b>3.949</b>	<b>3.412</b>	3.094
	MAE	2.444	3.995	3.713	<b>3.083</b>
	MSE	3.340	5.684	4.418	3.249
	Pearson+MSE	2.569	4.193	4.406	3.123
Autoformer	RI-Loss	<b>0.336</b>	<b>0.417</b>	<b>0.444</b>	0.445
	MAE	0.337	0.420	0.449	<b>0.443</b>
	MSE	0.375	0.451	0.477	0.476
	Pearson+MSE	0.364	0.421	0.449	0.448
DLinear	RI-Loss	<b>0.279</b>	<b>0.355</b>	<b>0.412</b>	<b>0.582</b>
	MAE	<b>0.279</b>	0.358	0.420	0.590
	MSE	0.297	0.373	0.450	0.696
	Pearson+MSE	0.305	0.404	0.494	0.714
iTrans.	RI-Loss	<b>0.296</b>	<b>0.374</b>	<b>0.419</b>	<b>0.422</b>
	MAE	0.298	0.377	0.421	0.427
	MSE	0.301	0.380	0.424	0.430
	Pearson+MSE	0.299	0.381	0.424	0.435
RAFT	RI-Loss	<b>0.293</b>	<b>0.377</b>	<b>0.415</b>	<b>0.417</b>
	MAE	0.294	0.379	0.420	0.422
	MSE	0.296	0.384	0.426	0.434
	Pearson+MSE	0.300	0.386	0.426	0.431

Table 4: Comparison of different loss functions across five backbone models on the ETTh2 dataset. "iTrans." denotes the iTransformer model.

reported results are averaged over five independent runs.

## Main Results

Table 1 presents a comparison of model performance using RI-Loss versus the traditional MSE Loss function. Overall, RI-Loss yields better results. We analyzed the prediction error reduction across all datasets and found that RI-Loss leads to significant improvements in both the Informer and DLinear models. Specifically, it achieves an average reduction of 9.4% in MSE for Informer, and 5.2% in MSE for DLinear. Other models also show similar gains. Across 160 test cases, RI-Loss outperforms MSE loss in 130 instances, further demonstrating its robustness and broad applicability. Results under the MAE metric are provided in Appendix F.1.

The Friedman test is conducted to evaluate the statistical significance of performance differences between the two loss functions, showing that all five models achieve varying degrees of ranking improvement after incorporating RI-Loss. Further details are given in Appendix E.

## Noise Robustness

To evaluate the robustness of our method under noisy conditions, we add zero-mean Gaussian noise to the input data at varying Signal-to-Noise Ratios (SNRs) and conduct experiments using DLinear as the backbone model. The results, summarized in Table 2, show that our loss consistently outperforms the conventional MSE across all SNR levels. While MSE suffers significant performance degradation as noise increases, our method effectively mitigates the impact

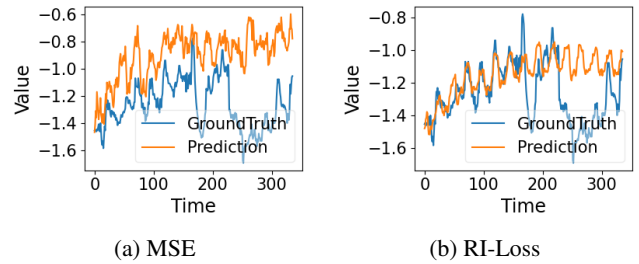


Figure 3: Forecasting results on ETTh1 with input-336-predict-336, comparing RI-Loss and MSE objectives. Blue denotes ground truth, and orange denotes predictions.

of noise, demonstrating better generalization and practical applicability for noisy time series forecasting.

## Running Cost

To evaluate the computational overhead of RI-Loss, we adopt DLinear, iTransformer, and RAFT as backbone models and measure per-iteration training time (ms/iter) on the Weather (forecasting horizon 720). As shown in Table 3, RI-Loss adds a small computational overhead because HSIC is costlier than MSE, but it consistently boosts backbone accuracy and remains effective for long-term forecasting on large datasets. Additional results are provided in Appendix F.2.

## Ablation Study

To analyze the impact of different components in our loss, we conducted four ablation experiments on ETTh2 dataset: RI-Loss, MAE, MSE, and RI-Loss with HSIC replaced by the Pearson Correlation (PC). As shown in Table 4, replacing HSIC with PC reduces performance, since PC captures only linear relationships and cannot model nonlinear dependencies in complex time series. RI-Loss also outperforms MSE and MAE, confirming the effectiveness of our design. Further details are in Appendix F.3.

## Visualization of Prediction Results

We use DLinear as the backbone model to visualize the prediction results. As shown in Figure 3, compared with the MSE loss, RI-Loss significantly improves the alignment with the ground truth. Predictions with RI-Loss align more closely with the ground truth's statistical properties, yielding results with higher structural similarity. More visualizations are provided in Appendix F.4.

## Conclusion

Time series forecasting is crucial across many domains; however, most existing models rely on MSE, which struggles with temporal dependencies and inherent noise. To address this, we propose RI-Loss, an HSIC-based loss function that enhances temporal pattern extraction while mitigating noise interference. We further establish a learning bound for HSIC to support its learnability. Experiments on diverse benchmarks demonstrate the superior performance of RI-Loss for time series forecasting. Future work will explore more complex noise correlations and residual dependencies.

## Acknowledgements

This work was supported by the National Natural Science Foundation of China (Nos. T2495251, 62306170, U24A20253, 62136005, 62476160, 62441239, 62276162), the Science and Technology Major Project of Shanxi (No. 202201020101006), the Special Fund for Science and Technology Innovation Teams of Shanxi Province (No. 202304051001001), and the Key Research and Development Project of Shanxi Province (No. 202402020101004).

## References

- Bartlett, P. L.; and Mendelson, S. 2003. Rademacher and gaussian complexities: risk bounds and structural results. *Journal of Machine Learning Research*, 3: 463–482.
- Box, G. E. P.; Jenkins, G. M.; Reinsel, G. C.; and Ljung, G. M. 1970. *Time Series Analysis: Forecasting and Control*.
- Cao, D.; Wang, Y.; Duan, J.; Zhang, C.; Zhu, X.; Huang, C.; Tong, Y.; Xu, B.; Bai, J.; Tong, J.; and Zhang, Q. 2020. Spectral temporal graph neural network for multivariate time-series forecasting. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, 17766–17778.
- Cao, W.; Wang, D.; Li, J.; Zhou, H.; Li, L.; and Li, Y. 2018. BRITS: Bidirectional Recurrent Imputation for Time Series. In *Neural Information Processing Systems*, 6776–6786.
- Cuturi, M.; and Blondel, M. 2017. Soft-DTW: a differentiable loss function for time-series. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, 894–903.
- Deng, A.; and Hooi, B. 2021. Graph Neural Network-Based Anomaly Detection in Multivariate Time Series. In *AAAI Conference on Artificial Intelligence*.
- Greenfeld, D.; and Shalit, U. 2020. Robust Learning with the Hilbert-Schmidt Independence Criterion. In III, H. D.; and Singh, A., eds., *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, 3759–3768. PMLR.
- Gretton, A.; Bousquet, O.; Smola, A.; and Schölkopf, B. 2005. Measuring Statistical Dependence with Hilbert-Schmidt Norms. In Jain, S.; Simon, H. U.; and Tomita, E., eds., *Algorithmic Learning Theory*, 63–77. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Han, S.; Lee, S.; Cha, M.; Arik, S. O.; and Yoon, J. 2025. Retrieval Augmented Time Series Forecasting. In *Forty-second International Conference on Machine Learning*.
- Hewage, P.; Behera, A.; Trovati, M.; Pereira, E.; Ghahremani, M.; Palmieri, F.; and Liu, Y. 2020. Temporal convolutional neural (TCN) network for an effective weather forecasting using time-series data from the local weather station. *Soft Comput.*, 24(21): 16453–16482.
- Hu, R.; Sejdinovic, D.; and Evans, R. J. 2024. A Kernel Test for Causal Association via Noise Contrastive Backdoor Adjustment. *Journal of Machine Learning Research*, 1–56.
- Le Guen, V.; and Thome, N. 2019. Shape and Time Distortion Loss for Training Deep Time Series Forecasting Models. In *Advances in Neural Information Processing Systems*, 4191–4203.
- Li, J.; Qian, Y.; Wang, J.; and Liu, S. 2024. PHSIC against Random Consistency and Its Application in Causal Inference. In *In Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, 2108–2116.
- Li, Y.; Pogodin, R.; Sutherland, D. J.; and Gretton, A. 2021. Self-Supervised Learning with Kernel Dependence Maximization. In *Neural Information Processing Systems*, volume 34, 15543–15556.
- Liu, M.; Zeng, A.; Chen, M.-H.; Xu, Z.; Lai, Q.; Ma, L.; and Xu, Q. 2022. SCINet: Time Series Modeling and Forecasting with Sample Convolution and Interaction. In *Neural Information Processing Systems*, volume 35, 5816–5828.
- Liu, Y.; Hu, T.; Zhang, H.; Wu, H.; Wang, S.; Ma, L.; and Long, M. 2024. iTransformer: Inverted Transformers Are Effective for Time Series Forecasting. In *The Twelfth International Conference on Learning Representations*.
- Ren, W.; Li, B.; and Han, M. 2020. A novel Granger causality method based on HSIC-Lasso for revealing nonlinear relationship between multivariate time series. *Physica A: Statistical Mechanics and its Applications*.
- Smola, A.; and Schölkopf, B. 2004. A tutorial on support vector regression. *Statistics and Computing*, 14: 199–222.
- Wang, H.; Pan, L.; Shen, Y.; Chen, Z.; Yang, D.; Yang, Y.; Zhang, S.; Liu, X.; Li, H.; and Tao, D. 2025. FreDF: Learning to Forecast in the Frequency Domain. In *The Thirteenth International Conference on Learning Representations*.
- Wang, S.; Zhang, L.; Zuo, W.; and Zhang, B. 2020. Class-Specific Reconstruction Transfer Learning for Visual Recognition Across Domains. *IEEE Transactions on Image Processing*, 29: 2424–2438.
- Wu, H.; Xu, J.; Wang, J.; and Long, M. 2021. Autoformer: decomposition transformers with auto-correlation for long-term series forecasting. In *Proceedings of the 35th International Conference on Neural Information Processing Systems*, 22419–22430.
- Xu, Z.; Zeng, A.; and Xu, Q. 2024. FITS: Modeling Time Series with 10k Parameters. In *The Twelfth International Conference on Learning Representations*.
- Yi, D.; and Wang, Y. 2019. *Applied Time Series Analysis*. Beijing: Renmin University of China Press, 5nd edition.
- Zeng, A.; Chen, M.-H.; Zhang, L.; and Xu, Q. 2022. Are Transformers Effective for Time Series Forecasting? In *AAAI Conference on Artificial Intelligence*, 11121–11128.
- Zhou, H.; Zhang, S.; Peng, J.; Zhang, S.; Li, J.; Xiong, H.; and Zhang, W. 2021. Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting. In *AAAI Conference on Artificial Intelligence*, 11106–11115.