

# BEYOND MODEL RANKING: PREDICTABILITY-ALIGNED EVALUATION FOR TIME SERIES FORECASTING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In the era of increasingly complex AI models for time series forecasting, progress is often measured by marginal improvements on benchmark leaderboards. However, this approach suffers from a fundamental flaw: standard evaluation metrics conflate a model’s performance with the data’s intrinsic unpredictability. To address this pressing challenge, we introduce a novel, predictability-aligned diagnostic framework grounded in spectral coherence. Our framework makes two primary contributions: the **Spectral Coherence Predictability (SCP)**, a computationally efficient ( $O(N \log N)$ ) and task-aligned score that quantifies the inherent difficulty of a given forecasting instance, and the **Linear Utilization Ratio (LUR)**, a frequency-resolved diagnostic tool that precisely measures how effectively a model exploits the linearly predictable information within the data. We validate our framework’s effectiveness and leverage it to reveal two core insights. First, we provide the first systematic evidence of “predictability drift”, demonstrating that a task’s forecasting difficulty varies sharply over time. Second, our evaluation reveals a key architectural trade-off: complex models are superior for low-predictability data, whereas linear models are highly effective on more predictable tasks. We advocate for a paradigm shift, moving beyond simplistic aggregate scores toward a more insightful, predictability-aware evaluation that fosters fairer model comparisons and a deeper understanding of model behavior. Codes and data are available at [https://anonymous.4open.science/r/TS\\_Predictability-C8B7](https://anonymous.4open.science/r/TS_Predictability-C8B7).

## 1 INTRODUCTION

Despite the proliferation of ever-more-complex models for time-series forecasting, true progress in the field remains notoriously difficult to measure (Bergmeir, 2024). The community relies on standard metrics, such as Mean Squared Error (MSE) and Mean Absolute Error (MAE), which quantify the magnitude of prediction errors but fail to reveal their origin. These metrics obscure a critical distinction: is a model failing because of its own limitations, or has it hit the fundamental ceiling of what is predictable in the data (Ehrenberg & Bound, 1993)? This ambiguity creates a severe evaluation dilemma, where a sophisticated model on a chaotic series may appear worse than a simple one. It also stalls scientific progress by wasting significant computational resources on problems where gains are impossible. Disentangling model inadequacy from intrinsic data unpredictability is therefore one of the most pressing challenges in the field today (Erkintalo, 2015).

The ideal solution to this dilemma is to first quantify a time series’ intrinsic predictability and then evaluate a model’s performance relative to this established baseline. However, designing a predictability metric that is fit for the modern deep learning forecasting paradigm presents a series of formidable challenges (Pennekamp et al., 2019). First, such a metric must be task-aligned, meaning its theoretical foundation should cohere with multi-horizon forecasting under a squared-error loss, rather than traditional single-step classification accuracy (Mishra & Palanisamy, 2018). Second, it must be computationally efficient to handle the massive, high-dimensional time series prevalent today (Fiecas et al., 2019). Finally, a single, global predictability score is insufficient; a truly useful tool must be diagnostic, offering insights to reveal where a model succeeds or fails in capturing predictable patterns.

Viewed through the lens of these challenges, existing tools are ill-suited for this purpose. Traditional proxies for predictability, such as entropy-rate estimators and Lempel-Ziv complexity, suffer from a fundamental paradigm mismatch (Aboy et al., 2006). These information-theoretic methods were developed for symbolic dynamics and classification tasks, designed to estimate an upper bound on next-symbol accuracy under a 0-1 loss, not a lower bound on multi-step regression error (Zhao et al., 2021). Furthermore, they often assume stationarity and are computationally expensive, which typically entails quadratic-to-cubic time complexity, rendering them impractical for the large-scale, non-stationary datasets common in modern applications (Kontoyiannis et al., 2002; Wyner & Ziv, 2002). An entirely new framework is thus urgently needed to bridge the significant gap between classical predictability theory and contemporary forecasting practice.

To bridge this gap, we introduce a novel diagnostic framework grounded in spectral coherence that is computationally efficient, directly aligned with the MSE objective, and provides rich, multi-scale insights. Our framework consists of two core components: 1) **Spectral Coherence Predictability (SCP)**, a per-instance predictability score  $\mathcal{P}$  derived from a linear MSE lower bound,  $\text{MSE}_{\text{lb}}$ . It is computed in  $O(N \log N)$  time and is intrinsically consistent with the MSE objective. 2) **Linear Utilization Ratio (LUR)**, a frequency-resolved diagnostic that quantifies how effectively a model exploits linearly predictable information across different spectral bands, enabling fine-grained assessments of under-use, saturation, and beyond-linear gains. With this framework, we aim to shift the time-series evaluation paradigm from simple “model ranking” toward a more profound analysis of “model-data diagnostics,” thereby providing clear, actionable guidance for developing more powerful and reliable forecasting models. Through extensive experiments on synthetic and real-world benchmarks, we validate our proposed framework. We demonstrate that our predictability score is well-calibrated and strongly correlates with the empirical performance of state-of-the-art models. Furthermore, our diagnostics uncover the highly time-varying nature of predictability, enabling a fairer, stratified evaluation that uncovers the differential strengths of different architectures, moving beyond the limitations of aggregate scores.

In summary, our contributions are as follows:

- We are the first to systematically address the evaluation ambiguity in modern time-series forecasting by introducing a framework to decouple model error from a time series’ intrinsic, irreducible unpredictability.
- We propose a computationally efficient and task-aligned Spectral Coherence Predictability (SCP). It provides a per-instance predictability score with a corresponding error lower bound and a frequency-resolved diagnostic to analyze a model’s information utilization.
- Extensive experiments validates this framework’s alignment with state-of-the-art models. We then use it to reveal novel insights like “predictability drift” and introduce a stratified evaluation that uncovers the complementary strengths of different models.

## 2 RELATED WORK

Our work is positioned at the intersection of two key research areas: the quantification of sequence predictability and the use of spectral methods for time-series analysis.

**Predictability of Time Series.** Entropy-based notions have long been used to proxy sequence predictability, from Shannon’s entropy and entropy rate to variants usable on continuous data (approximate, sample, fuzzy, and permutation entropy) (Shannon, 1948; Pincus, 1991; Richman & Moorman, 2000; Bandt & Pompe, 2002; Garland et al., 2014). Compression-driven estimators (e.g., Lempel-Ziv) provide nonparametric estimates of entropy rate for symbolic, stationary sources (Ziv & Lempel, 1977). These approaches have also been popular in human mobility, where spatio-temporal regularity supports predictability limits under coarse symbolizations (González et al., 2008; Song et al., 2010; Wang et al., 2021). However, they face three key limitations for general forecasting: (i) computational burden the need for discretization of continuous data; (ii) theoretical misalignment with multi-step squared-loss objectives; and (iii) sensitivity of differential entropy to reparameterization and divergence issues in non-stationary settings (Mohammed et al., 2024). Consequently, existing predictability proxies fail to provide a task-aligned and reliable estimate of the achievable error lower bound for multi-step time-series forecasting.

**Spectral Analysis in Time Series Forecasting.** Spectral analysis is a cornerstone of time-series modeling, inspiring many recent deep learning architectures. For instance, Autoformer was designed with an auto-correlation mechanism to discover period-based dependencies efficiently (Wu et al., 2021). FEDformer directly integrates Fourier/Wavelet transforms into attention for frequency-domain computation with reduced complexity (Zhou et al., 2022). TimesNet captures complex multi-periodicity by transforming the 1D time series into a 2D representation for analysis (Wu et al., 2023). These methods all leverage spectral properties to build better models. In contrast, our work uses spectral coherence to build a novel diagnostic framework for analyzing data predictability and evaluating the utilization of existing models.

### 3 PRELIMINARIES

**Problem setup and notation.** We focus on a setting in which an observed sequence is decomposed into a past (history) portion used as input and a future portion serving as ground-truth for evaluation. Formally, a sample consists of a history  $\mathbf{x} \in \mathbb{R}^N$  and a future  $\mathbf{y} \in \mathbb{R}^N$  drawn from a distribution  $\mathbb{D}$ . The goal is to learn a measurable predictor  $f : \mathbb{R}^N \rightarrow \mathbb{R}^N$  that produces a forecast  $\hat{\mathbf{y}} = f(\mathbf{x})$ . We evaluate predictions with the mean squared error (MSE) per forecast step:

$$\text{MSE}(f; \mathbf{x}, \mathbf{y}) = \frac{1}{N} \|\mathbf{f}(\mathbf{x}) - \mathbf{y}\|_2^2, \quad (1)$$

where  $\|\cdot\|_2$  denotes the Euclidean norm. The learning objective is to minimize the expected risk  $\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[\text{MSE}(f; \mathbf{x}, \mathbf{y})]$ .

**Intrinsic predictability via Bayes risk.** Under the MSE metric, the risk-minimizing predictor is the conditional expectation  $f^*(\mathbf{x}) = \mathbb{E}[\mathbf{y} | \mathbf{x}]$  (Chen et al., 2016). The corresponding minimum achievable risk (Bayes risk) is

$$\text{MSE}^* = \mathbb{E} \left[ \frac{1}{N} \|\mathbf{y} - \mathbb{E}[\mathbf{y} | \mathbf{x}]\|_2^2 \right]. \quad (2)$$

Using the unconditional variance  $\text{Var}(\mathbf{y})$  as a baseline, we define intrinsic predictability as the normalized reduction of uncertainty:

$$\mathcal{P}_{xy}^* = 1 - \frac{\text{MSE}^*}{\text{Var}(\mathbf{y})}, \quad (3)$$

which coincides with the theoretical coefficient of determination ( $R^2$ ). By the law of total variance,  $\text{Var}(\mathbf{y}) = \text{MSE}^* + \text{Var}(\mathbb{E}[\mathbf{y} | \mathbf{x}])$ , hence  $\mathcal{P}^* \in [0, 1]$ . At the extremes,  $\mathcal{P}^* = 1$  if and only if  $\text{Var}(\mathbf{y} | \mathbf{x}) = 0$  almost surely, i.e.,  $\mathbf{y}$  is a deterministic function of  $\mathbf{x}$ . Conversely,  $\mathcal{P}^* = 0$  if and only if  $\text{Var}(\mathbb{E}[\mathbf{y} | \mathbf{x}]) = 0$ , i.e., the conditional mean  $\mathbb{E}[\mathbf{y} | \mathbf{x}]$  is constant and  $\mathbf{x}$  conveys no information for predicting the mean of  $\mathbf{y}$ .

## 4 METHOD

### 4.1 SPECTRAL COHERENCE PREDICTABILITY

Building on the intrinsic predictability in Eq. (3), we seek a computable surrogate from a single realization by leveraging frequency-domain structure. The resulting Spectral Coherence Predictability (SCP) quantifies how much of the future segment  $\mathbf{y}$  is linearly explainable by the history segment  $\mathbf{x}$  across frequencies, and then aggregates the explained and unexplained power.

We operate in the frequency domain using Welch’s method. Let  $\hat{S}_{yy}(f)$  and  $\hat{S}_{xx}(f)$  denote the power spectral densities (PSD) of  $\mathbf{y}$  and  $\mathbf{x}$ , and let  $\hat{S}_{xy}(f)$  denote their cross-power spectral density (CPSD). All spectra are computed on the same discrete Fourier transform (DFT) grid with identical Welch parameters after mean removal. The squared coherence between  $\mathbf{y}$  and  $\mathbf{x}$  is

$$\gamma_{xy}^2(f) = \frac{|\hat{S}_{xy}(f)|^2}{(\hat{S}_{xx}(f) + \varepsilon)(\hat{S}_{yy}(f) + \varepsilon)} \in [0, 1], \quad (4)$$

**Algorithm 1** Spectral Coherence Predictability (SCP)

**Require:** History  $\mathbf{x} \in \mathbb{R}^N$ , future  $\mathbf{y} \in \mathbb{R}^N$ ; Welch parameters (window, segment length, overlap); stability constant  $\varepsilon > 0$ ; optional frequency band  $\mathcal{F}_b$ .

**Ensure:** MSE linear lower bound  $\text{MSE}_{\text{lb}}$  and predictability  $\mathcal{P}_{xy}$ .

- 1: **Mean removal:**  $m_x \leftarrow \text{mean}(x)$ ,  $m_y \leftarrow \text{mean}(y)$ ;  $\Delta^2 \leftarrow (m_y - m_x)^2$ ;  $x \leftarrow x - m_x$ ,  $y \leftarrow y - m_y$ .
- 2: **Welch spectra:** Compute the PSD  $\widehat{S}_{xx}(f)$ ,  $\widehat{S}_{yy}(f)$  and the CPSD  $\widehat{S}_{xy}(f)$  on the discrete frequency domain  $\mathcal{F}$ .
- 3: **Squared coherence:**

$$\gamma^2(f) \leftarrow \frac{|\widehat{S}_{xy}(f)|^2}{(\widehat{S}_{xx}(f) + \varepsilon)(\widehat{S}_{yy}(f) + \varepsilon)} \in [0, 1], \quad \forall f \in \mathcal{F}.$$

- 4: **Residual spectrum:**  $\widehat{S}_e(f) \leftarrow \widehat{S}_{yy}(f)(1 - \gamma^2(f))$ ,  $\forall f \in \mathcal{F}$ .
- 5: **Frequency set:**  $\mathcal{F}_* \leftarrow \mathcal{F}_b$  if a band  $\mathcal{F}_b$  is provided; otherwise  $\mathcal{F}_* \leftarrow \mathcal{F}$ .
- 6: **Aggregate:**

$$\widehat{\text{Var}}(y) \leftarrow \sum_{f \in \mathcal{F}_*} \widehat{S}_{yy}(f), \quad \text{MSE}_{\text{lb}} \leftarrow \Delta^2 + \sum_{f \in \mathcal{F}_*} \widehat{S}_e(f).$$

- 7: **Predictability:**  $\mathcal{P}_{xy} \leftarrow 1 - \text{MSE}_{\text{lb}}/\widehat{\text{Var}}(y)$ .
- 8: **return**  $\text{MSE}_{\text{lb}}$ ,  $\mathcal{P}_{xy}$ .

where  $\varepsilon > 0$  is a small term for numerical stability (Mandel & Wolf, 1976; Wang et al., 2019). Interpreting  $\gamma_{xy}^2(f)$  as a linearly explained–power ratio, the unexplained (residual) spectrum is

$$\widehat{S}_e(f) = \widehat{S}_{yy}(f)(1 - \gamma_{xy}^2(f)). \quad (5)$$

Let  $\mathcal{F}$  denote the discrete frequency domain under our normalization, so that the total spectral power equals the sample variance, i.e.,  $\sum_{f \in \mathcal{F}} \widehat{S}_{yy}(f) = \widehat{\text{Var}}(\mathbf{y})$ . To correct a boundary mean mismatch, we optionally add a boundary mean–shift term  $\Delta^2 = (\text{mean}(\mathbf{y}) - \text{mean}(\mathbf{x}))^2$ . This yields a lower bound on the mean squared error of any linear time–invariant predictor that uses  $\mathbf{x}$ :

$$\text{MSE}_{\text{lb}} = \Delta^2 + \sum_{f \in \mathcal{F}} \widehat{S}_e(f), \quad \widehat{\text{Var}}(\mathbf{y}) = \sum_{f \in \mathcal{F}} \widehat{S}_{yy}(f). \quad (6)$$

The SCP estimate of predictability is then

$$\mathcal{P}_{xy} = 1 - \frac{\text{MSE}_{\text{lb}}}{\widehat{\text{Var}}(\mathbf{y})} \in [0, 1]. \quad (7)$$

Algorithm 1 summarizes the steps. Computationally, with fast Fourier transform and Welch estimation, SCP costs  $\mathcal{O}(N \log N)$  per sample. This is substantially lower than matching–based Lempel–Ziv–style predictability estimators, which typically entail at least quadratic–to–cubic time in sequence length (e.g.,  $\mathcal{O}(N^3)$  in naive implementations) and usually target single–step predictability, whereas SCP yields a multi–step estimate aligned with the evaluation horizon.

**Theoretical interpretation.** If  $(\mathbf{x}, \mathbf{y})$  is jointly Gaussian and wide–sense stationary around the boundary, the Bayes predictor is linear (Ko & Fox, 2009). In this case, Eq. (7) is a consistent estimator of the intrinsic predictability  $\mathcal{P}_{xy}^*$  as the effective sample size grows. For general (possibly non-linear) processes,  $\widehat{\text{MSE}}_{\text{lb}}$  lower–bounds the error of any linear time–invariant forecaster using  $\mathbf{x}$ , hence  $\mathcal{P}_{xy}$  is a conservative estimate that isolates the explanatory power of stable linear dynamics across frequencies and sets a meaningful baseline for non-linear forecasters.

## 4.2 SLINEAR UTILIZATION RATIO

Instead of relying on standard metrics like MSE or MAE, which only tell us if a model is accurate, our approach analyzes why it is accurate by comparing its performance to the data’s theoretical predictability limit in the frequency domain.

**Algorithm 2** Linear Utilization Ratio (LUR)

**Require:** History  $\mathbf{x} \in \mathbb{R}^N$ , future  $\mathbf{y} \in \mathbb{R}^N$ , model prediction  $\hat{\mathbf{y}} \in \mathbb{R}^N$ ; Welch parameters (window, segment length, overlap); stability  $\varepsilon > 0$ ; optional frequency band  $\mathcal{F}_b$ .

**Ensure:** Model-explained power  $P_{\text{model}}$ ; linear-explainable power  $P_{\text{linear}}$ ; evaluation ratio LUR.

1: **Mean removal:**  $\mathbf{x} \leftarrow \mathbf{x} - \text{mean}(\mathbf{x})$ ;  $\mathbf{y} \leftarrow \mathbf{y} - \text{mean}(\mathbf{y})$ ;  $\hat{\mathbf{y}} \leftarrow \hat{\mathbf{y}} - \text{mean}(\hat{\mathbf{y}})$ .

2: **Welch spectra (same parameters for all):**

$$\hat{S}_{xx}(f), \hat{S}_{yy}(f), \hat{S}_{\hat{y}\hat{y}}(f), \hat{S}_{xy}(f), \hat{S}_{y\hat{y}}(f), \quad \forall f \in \mathcal{F}.$$

3: **Coherences:**

$$\gamma_{yx}^2(f) \leftarrow \frac{|\hat{S}_{yx}(f)|^2}{(\hat{S}_{yy}(f) + \varepsilon)(\hat{S}_{xx}(f) + \varepsilon)}, \quad \gamma_{y\hat{y}}^2(f) \leftarrow \frac{|\hat{S}_{y\hat{y}}(f)|^2}{(\hat{S}_{yy}(f) + \varepsilon)(\hat{S}_{\hat{y}\hat{y}}(f) + \varepsilon)}.$$

▷ Eq. (4) and Eq. (8).

4: **Frequency set:**  $\mathcal{F}_* \leftarrow \mathcal{F}_b$  if a band  $\mathcal{F}_b$  is provided; otherwise  $\mathcal{F}_* \leftarrow \mathcal{F}$ .

5: **Power-weighted aggregation:**

$$P_{\text{model}} \leftarrow \sum_{f \in \mathcal{F}_*} \gamma_{y\hat{y}}^2(f) \hat{S}_{yy}(f), \quad P_{\text{linear}} \leftarrow \sum_{f \in \mathcal{F}_*} \gamma_{yx}^2(f) \hat{S}_{yy}(f).$$

6: **LUR ratio:**  $\text{LUR} \leftarrow P_{\text{model}}/P_{\text{linear}}$ .

7: **return**  $P_{\text{model}}$ ,  $P_{\text{linear}}$ , LUR.

Our method, detailed in Algorithm 2, is built on two key quantities: The first key quantity is the linear limit, given by the squared coherence between the future  $\mathbf{y}$  and the history  $\mathbf{x}$ ,  $\gamma_{yx}^2(f)$  in Eq. (4), which upper-bounds, at frequency  $f$ , the fraction of  $\mathbf{y}$ 's power that any linear predictor can explain using  $\mathbf{x}$ . The second quantity is the prediction-target coherence, measuring how much of  $\mathbf{y}$ 's power is actually captured by the model prediction  $\hat{\mathbf{y}}$  at frequency  $f$ :

$$\gamma_{y\hat{y}}^2(f) = \frac{|\hat{S}_{y\hat{y}}(f)|^2}{(\hat{S}_{yy}(f) + \varepsilon)(\hat{S}_{\hat{y}\hat{y}}(f) + \varepsilon)} \in [0, 1]. \quad (8)$$

Comparing these two coherences yields a per-frequency diagnosis:

- If  $\gamma_{y\hat{y}}^2(f) < \gamma_{yx}^2(f)$ , the model under-utilizes linearly predictable information in the history  $\mathbf{x}$ .
- If  $\gamma_{y\hat{y}}^2(f) \approx \gamma_{yx}^2(f)$ , the model saturates the linear limit of the history  $\mathbf{x}$ .
- If  $\gamma_{y\hat{y}}^2(f) > \gamma_{yx}^2(f)$ , the model surpasses the linear limit of the history  $\mathbf{x}$ , indicating exploitation of nonlinear patterns or additional, correlated data sources.

To summarize over frequency while prioritizing strong signal regions, we weight coherence by the target power spectrum  $\hat{S}_{yy}(f)$  on the discrete frequency domain  $\mathcal{F}$ . Define the aggregated powers

$$P_{\text{model}} = \sum_{f \in \mathcal{F}} \gamma_{y\hat{y}}^2(f) \hat{S}_{yy}(f), \quad P_{\text{linear}} = \sum_{f \in \mathcal{F}} \gamma_{yx}^2(f) \hat{S}_{yy}(f), \quad (9)$$

The fraction of target energy that is linearly predictable is

$$\eta_{\text{linear}} = \frac{P_{\text{linear}}}{\sum_{f \in \mathcal{F}} \hat{S}_{yy}(f)} = \frac{P_{\text{linear}}}{\widehat{\text{Var}}(\mathbf{y})}. \quad (10)$$

To further assess how much of this predictable energy the model actually captures, we use the Linear Utilization Ratio (LUR)

$$\text{LUR} = \frac{P_{\text{model}}}{P_{\text{linear}}} \geq 0. \quad (11)$$

When  $\text{LUR} < 1$ , the model under-exploits linear information available in the history  $\mathbf{x}$ ; when  $\text{LUR} \approx 1$ , it saturates the linear limit; when  $\text{LUR} > 1$ , it achieves gains beyond linear structure (e.g., nonlinear dynamics or cross-instance inductive patterns learned from the training set).

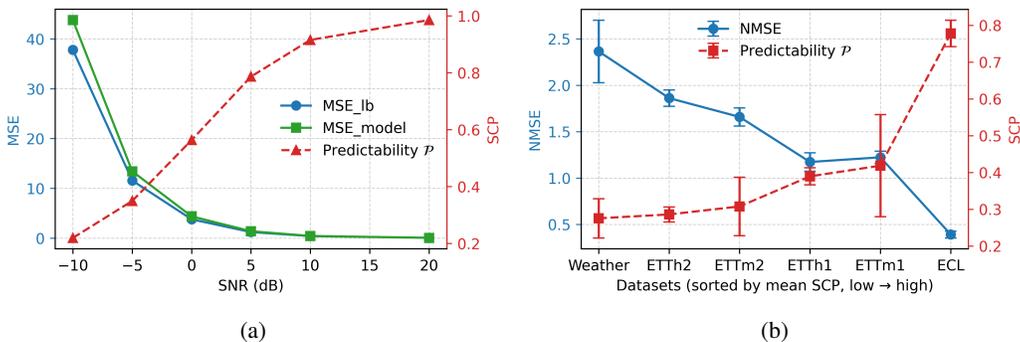


Figure 1: Relationship between model error (MSE) and predictability  $\mathcal{P}$ . (a) MSE of the best linear predictor on a synthetic Gaussian process with varying SNR. (b) Average performance of state-of-the-art prediction models on real datasets. We report normalized MSE (NMSE), obtained by dividing MSE by the corresponding variance.

To analyze behavior across scales, we additionally partition the discrete frequency domain into disjoint bands  $\{\mathcal{F}_b\}_{b=1}^B$  (e.g., low/mid/high), using the band partition as in Algorithms 1 and 2. This yields band-limited counterparts  $MSE_{lb,b}$ ,  $\mathcal{P}_{xy,b}$ , and  $LUR_b$ , which enable localized diagnosis of under-use, saturation, or beyond-linear gains within each frequency band.

## 5 EXPERIMENTS

This section empirically evaluates our metric and diagnostics across synthetic and real-world settings, guided by three questions: **Q1:** Does our proposed metric accurately reflect the intrinsic predictability of time series data (Secs. 5.1 and 5.2)? **Q2:** What novel insights and observations can be derived from this new metric (Sec. 5.3)? **Q3:** How can our metric be leveraged to enable more comprehensive evaluation methodologies for forecasting models (Secs. 5.4 and 5.5)?

### 5.1 TOY STUDY

We first validate our proposed predictability score  $\mathcal{P}$  and theoretical error lower bound  $MSE_{lb}$  in a controlled, synthetic environment. Our setup consists of a Gaussian process with additive noise at varying SNRs, on which we evaluate an optimal linear forecaster (Fig. 1a). For each SNR, we report the model’s test MSE, the estimated linear lower bound  $MSE_{lb}$ , and the corresponding predictability  $\mathcal{P}$ . As noise decreases (higher SNR),  $\mathcal{P}$  increases monotonically toward one, and the model MSE approaches  $MSE_{lb}$ . Across all SNRs,  $MSE_{lb}$  remains strictly below the realized MSE, confirming it is a valid lower bound; the shrinking gap at high SNR indicates that the linear model saturates the data-implied limit when the process is nearly noise-free. It validates both the calibration (monotonic response to noise) and the tightness (small bound–error gap in the linear regime) of our estimates.

### 5.2 ALIGNING PREDICTABILITY AND FORECASTING PERFORMANCE

To ensure a fair and rigorous comparison, we evaluate all baselines under a strictly controlled setup: (i) the forecast horizon and the history length are fixed and identical across models; (ii) the common “drop–last” heuristic is disabled. For correlation analyses, we report the Pearson coefficient  $R$  between each model’s empirical prediction MSE and the theoretically estimated  $MSE_{lb}$ , averaged over all variables and samples in the test set to provide a holistic summary.

As shown in Table 1,  $MSE_{lb}$  strongly correlates with the empirical MSE of SOTA models across datasets, with Pearson correlations typically  $R \geq 0.8$ . This alignment indicates that the bound accurately predicts where forecasting is intrinsically easier or harder, in agreement with realized model accuracy. To visualize this relationship, Fig. 2 plots iTransformer’s MSE against  $MSE_{lb}$  for each channel in ECL and Weather datasets, revealing a near-linear trend. Collectively, these dataset- and channel-level results validate the effectiveness of our  $MSE_{lb}$  estimator. They also highlight substantial within-dataset heterogeneity: predictability and MSE can vary markedly across channels.

Table 1: Long-term multivariate forecasting results. We report MSE, MAE, NMSE for forecasting lengths equal to history length  $N \in \{96, 192, 336, 720\}$  under an identical protocol (same preprocessing and no drop-last). **Bold** marks the best (lowest) MSE/MAE per column across models. *Average* rows give the column-wise mean across models. Predictability reports the per-task linear MSE lower bound ( $MSE_{lb}$ ) and SCP  $\mathcal{P}$  (higher is easier).

Models	Metric	ETTh1				ETTh2				ETTm1				ETTm2				ECL				Weather			
		96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720				
iTransformer (Liu et al., 2024)	MSE	0.387	0.441	0.471	0.700	0.301	0.381	0.426	0.425	0.342	0.345	0.379	0.448	0.186	0.254	0.289	0.382	<b>0.148</b>	0.156	0.170	<b>0.194</b>	0.176	0.214	0.255	0.353
	MAE	0.405	0.440	0.464	0.608	0.350	0.405	0.438	0.455	0.377	0.378	0.403	0.449	0.272	0.319	0.341	0.407	<b>0.239</b>	0.250	0.266	<b>0.287</b>	0.216	0.255	0.290	0.387
	NMSE	1.121	1.095	1.056	1.292	1.598	1.829	1.563	1.432	1.340	1.167	1.140	1.159	1.609	1.696	1.595	1.917	0.288	0.268	0.280	0.319	2.629	2.353	2.032	2.150
	R	0.844	0.876	0.899	0.747	0.907	0.883	0.877	0.842	0.845	0.782	0.803	0.869	0.868	0.834	0.898	0.840	0.723	0.778	0.821	0.826	0.900	0.876	0.801	0.824
TimeMixer (Wang et al., 2024)	MSE	<b>0.381</b>	0.440	0.482	0.631	<b>0.289</b>	0.377	0.390	0.435	<b>0.322</b>	0.337	0.380	0.484	<b>0.176</b>	0.231	0.280	0.376	0.153	0.155	0.172	0.214	<b>0.169</b>	<b>0.198</b>	0.249	0.347
	MAE	0.400	0.434	0.460	0.561	<b>0.340</b>	0.406	0.423	0.458	0.359	0.372	0.396	0.469	0.259	0.296	0.332	0.390	0.245	0.244	0.264	0.310	<b>0.215</b>	<b>0.242</b>	0.291	0.355
	NMSE	1.131	1.069	1.062	1.155	1.540	1.724	1.524	1.446	1.303	1.105	1.135	1.202	1.493	1.485	1.498	1.689	0.282	0.274	0.286	0.334	2.602	2.161	2.107	2.162
	R	0.815	0.889	0.848	0.793	0.916	0.801	0.909	0.906	0.829	0.781	0.752	0.798	0.843	0.867	0.910	0.732	0.706	0.607	0.682	0.887	0.911	0.862	0.885	0.852
DLinear (Zeng et al., 2023)	MSE	0.383	<b>0.422</b>	0.447	0.507	0.329	0.375	0.463	0.740	0.346	0.342	0.372	<b>0.415</b>	0.187	0.242	0.278	0.374	0.195	0.163	0.169	0.197	0.197	0.225	0.263	0.315
	MAE	<b>0.396</b>	<b>0.421</b>	0.448	0.517	0.380	0.410	0.472	0.609	0.374	0.369	<b>0.389</b>	<b>0.415</b>	0.281	0.315	0.338	0.406	0.277	0.259	0.268	0.295	0.255	0.282	0.314	0.354
	NMSE	1.214	1.143	1.310	1.782	2.927	2.067	2.728	3.896	1.327	1.205	1.208	1.121	1.676	1.722	1.620	1.793	0.868	0.678	0.574	0.684	3.507	2.899	2.474	2.069
	R	0.869	0.878	0.872	0.804	0.845	0.880	0.798	0.439	0.868	0.887	0.819	0.884	0.833	0.813	0.910	0.902	0.880	0.867	0.909	0.864	0.924	0.931	0.911	0.923
PatchTST (Nie et al., 2023)	MSE	0.391	0.429	<b>0.436</b>	<b>0.465</b>	0.293	<b>0.357</b>	<b>0.363</b>	<b>0.406</b>	<b>0.322</b>	<b>0.328</b>	0.395	0.417	0.177	<b>0.230</b>	<b>0.276</b>	<b>0.356</b>	0.167	<b>0.151</b>	<b>0.167</b>	0.212	0.176	0.202	<b>0.247</b>	<b>0.309</b>
	MAE	0.403	0.426	<b>0.440</b>	<b>0.482</b>	0.342	<b>0.387</b>	<b>0.402</b>	<b>0.442</b>	<b>0.358</b>	<b>0.364</b>	0.390	0.419	<b>0.258</b>	<b>0.294</b>	<b>0.329</b>	0.385	0.252	<b>0.242</b>	<b>0.258</b>	0.304	0.217	0.243	<b>0.281</b>	<b>0.331</b>
	NMSE	1.205	1.227	1.113	1.341	1.952	2.115	1.556	1.527	1.391	1.424	1.352	1.361	1.714	1.621	1.933	2.012	0.316	0.345	0.350	0.482	2.555	2.531	2.294	2.055
	R	0.869	0.881	0.915	0.876	0.886	0.920	0.911	0.864	0.813	0.652	0.741	0.835	0.904	0.908	0.863	0.909	0.735	0.809	0.737	0.969	0.912	0.903	0.868	0.865
TimesNet (Wu et al., 2023)	MSE	0.389	0.460	0.487	0.641	0.337	0.405	0.399	0.447	0.334	0.414	0.429	0.482	0.189	0.239	0.320	0.383	0.168	0.189	0.209	0.305	<b>0.169</b>	0.220	0.272	0.334
	MAE	0.412	0.456	0.477	0.582	0.371	0.424	0.433	0.463	0.375	0.414	0.434	0.477	0.266	0.306	0.357	0.408	0.272	0.291	0.308	0.382	0.219	0.265	0.301	0.350
	NMSE	1.205	1.227	1.113	1.341	1.952	2.115	1.556	1.527	1.391	1.424	1.352	1.361	1.714	1.621	1.933	2.012	0.316	0.345	0.350	0.482	2.555	2.531	2.294	2.055
	R	0.869	0.881	0.915	0.876	0.886	0.920	0.911	0.864	0.813	0.652	0.741	0.835	0.904	0.908	0.863	0.909	0.735	0.809	0.737	0.969	0.912	0.903	0.868	0.865
Average	MSE	0.386	0.438	0.465	0.589	0.310	0.379	0.408	0.491	0.333	0.353	0.385	0.449	0.183	0.239	0.289	0.374	0.166	0.163	0.177	0.224	0.177	0.212	0.257	0.332
	MAE	0.403	0.435	0.458	0.550	0.357	0.406	0.434	0.485	0.369	0.379	0.402	0.446	0.267	0.306	0.339	0.399	0.257	0.257	0.273	0.316	0.224	0.257	0.295	0.349
	NMSE	1.149	1.120	1.105	1.321	1.912	1.863	1.738	1.936	1.321	1.199	1.193	1.184	1.608	1.599	1.626	1.805	0.415	0.366	0.356	0.432	2.814	2.421	2.179	2.047
Predictability	$MSE_{lb}$	0.354	0.417	0.404	0.412	0.298	0.360	0.309	0.356	0.228	0.307	0.513	0.436	0.175	0.248	0.295	0.361	0.239	0.219	0.167	0.241	0.185	0.244	0.278	0.317
	$\mathcal{P}$	0.422	0.379	0.368	0.389	0.305	0.270	0.302	0.267	0.590	0.460	0.268	0.356	0.415	0.315	0.230	0.271	0.751	0.755	0.829	0.777	0.345	0.240	0.228	0.289

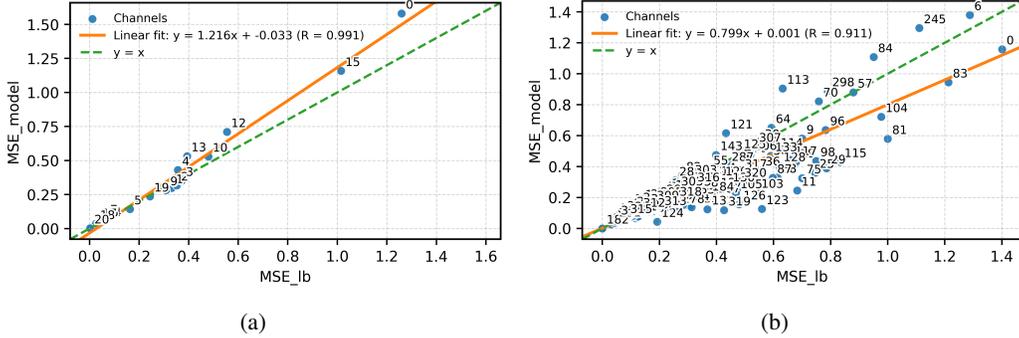


Figure 2: Per-variable scatter plots on Weather (a) and ECL (b) comparing the estimated MSE lower bound ( $MSE_{lb}$ ) with iTransformer’s prediction error ( $MSE_{model}$ ).

To compare predictability across datasets, we aggregate over horizons and report, for each dataset, the mean±std of SCP  $\mathcal{P}$  and realized NMSE (Fig. 1b). The relationship is strongly inverse: datasets with higher  $\mathcal{P}$  tend to exhibit lower NMSE on average. In particular, ECL shows the highest  $\mathcal{P}$  and the lowest NMSE (easiest to forecast), whereas Weather exhibits the lowest  $\mathcal{P}$  and the highest NMSE (hardest). This ordering indicates that  $\mathcal{P}$  ranks dataset difficulty in a manner that aligns with realized errors.

Overall, no single architecture dominates across all datasets and horizons—an expected outcome in predictability-limited settings that, as shown later, arises from time- and band-dependent variability in the exploitable history  $x$  that challenges both linear and nonlinear representations.

### 5.3 TIME-VARYING PREDICTABILITY

Standard evaluation metrics, which average performance over an entire test set, implicitly assume that the forecasting task is statistically stationary. However, this assumption often fails for time series tasks. Focusing on a single channel, we move sample by sample, plotting two key metrics in parallel: the model’s instantaneous error (MSE) and our measure of the data’s predictable energy, decomposed by frequency.

The results, shown in Fig. 3, are striking. We find that a model’s performance is not random but is tightly coupled to the instantaneous predictable energy in the data. This reveals that predictability is not a fixed property of a dataset; it varies sharply from one model to the next. This phenomenon,

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390

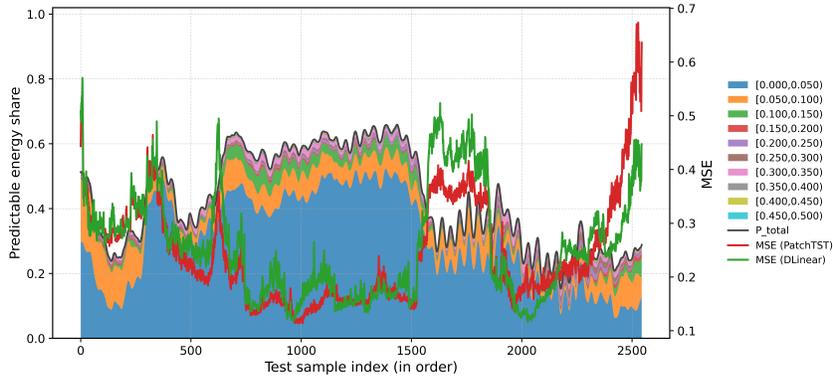


Figure 3: ETTh1 test set, channel 1, horizon  $N = 336$ . Relationship between per-sample linearly predictable energy (decomposed by frequency band as a share of total) and the corresponding MSE of DLinear and PatchTST.

391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407

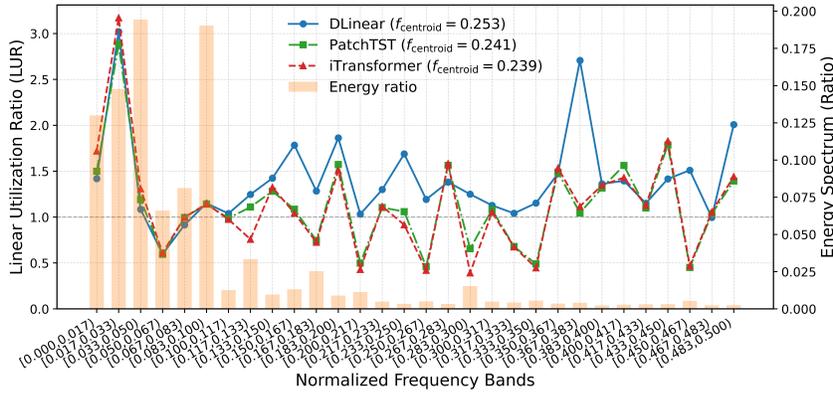


Figure 4: Band-wise analysis on ETTh1, representative channel. Normalized energy shares and LUR across frequency bands for three models.

412 which we term predictability drift, means the forecasting task should be viewed not as a single  
413 problem, but as a time-varying mixture of easy and hard regimes.

414 When the total predictable share is low, or when the dominant predictable bands shift rapidly, MSE  
415 spikes. These observations explain why aggregate test-set statistics often fail to discriminate mod-  
416 els. Fluctuations in sample-level predictability always dominate realized error. They also motivate  
417 finer-grained evaluation protocols, such as band-wise diagnostics and predictability-aware report-  
418 ing, which enable grouping evaluations conditioned on data characteristics and provide actionable  
419 guidance for model development.

420  
421 **5.4 BAND-WISE EVALUATION**

422 To gain a more granular understanding, we use the Linear Utilization Ratio (LUR) to analyze model  
423 behavior in the frequency domain. We decompose the test signal into distinct frequency bands and,  
424 for each band, measure both the proportion of the signal’s total energy and the LUR for several  
425 representative models.  
426

427 Figure 4 reveals the distinct strategies employed by different model architectures. In the five low-  
428 frequency bands, which contain the majority of the signal’s energy, all three models exhibit a similar  
429 pattern: their LURs closely track the energy distribution. This indicates that all models correctly  
430 identify and attempt to capture these dominant, high-energy components of the signal. Within these  
431 critical bands, however, the PatchTST and iTransformer achieve a slightly higher LUR than DLinear,  
suggesting they are more efficient at extracting information from these primary signal components.

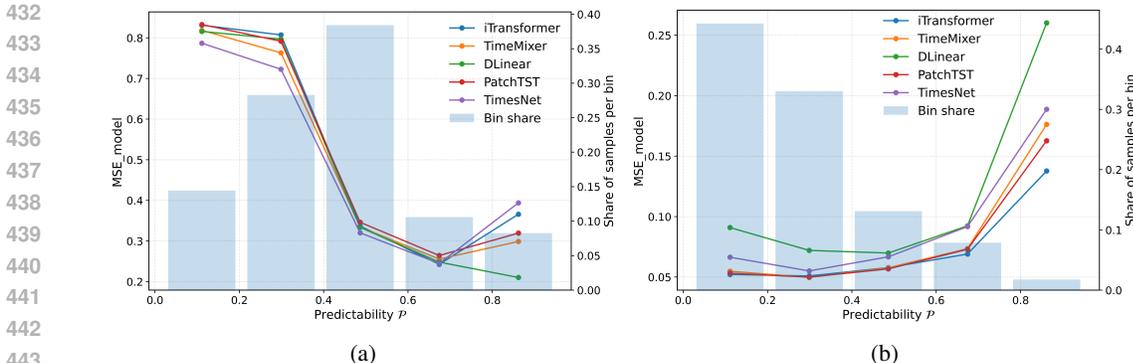


Figure 5: ETTh1 dataset with forecasting length  $N = 96$ : per-channel evaluation stratified by predictability  $\mathcal{P}$ . Samples are grouped into equal-width  $\mathcal{P}$  bins; each point reports the mean MSE within the bin.

A more significant divergence appears in the higher-frequency bands. Here, DLinear maintains a significantly higher LUR than both PatchTST and iTransformer. This finding exposes the fundamental difference in their underlying mechanisms. DLinear, being a simple linear model, acts like a broadband filter, attempting to capture information indiscriminately across the entire spectrum. In contrast, Transformer-based models like PatchTST and iTransformer function as frequency-selective filters. They intelligently focus their capacity on the high-energy, low-frequency bands where the most predictable signal resides, while largely ignoring the less informative, high-frequency noise. This demonstrates their more sophisticated inductive bias for typical time-series data.

### 5.5 PREDICTABILITY-AWARE EVALUATION

Although most models achieve a similar MSE of  $\approx 0.38$  on the ETTh1 dataset (horizon  $N = 96$ ), this average score obscures critical performance differences. By stratifying the test set by our predictability score  $\mathcal{P}$ , we uncover these hidden behaviors (Fig. 5). In Fig. 5a, **nonlinear forecasters (e.g., TimesNet) perform better in the low- $\mathcal{P}$  bin (hard samples), whereas DLinear outperforms in the high- $\mathcal{P}$  (easy samples) bin.** In Fig. 5b, where samples concentrate in the low- $\mathcal{P}$  region, iTransformer, TimeMixer, and PatchTST exhibit convergent performance, while DLinear is capacity-limited on hard cases and tends to underfit easy cases due to their low ratio. These contrasts highlight complementary strengths across architectures: nonlinear models excel when the linear predictable signal is scarce, whereas simple linear models are highly competitive when predictability is high.

## 6 CONCLUSION AND FUTURE WORK

This work confronts a fundamental flaw in time series forecasting evaluation: the conflation of model error with intrinsic data difficulty. We introduced a predictability-aligned diagnostic framework, centered on the Spectral Coherence Predictability (SCP) score and the Linear Utilization Ratio (LUR), to disentangle these factors. Our analysis not only provided the first systematic evidence of “predictability drift” but also revealed that different architectures possess complementary strengths, a nuance lost in aggregate metrics. Ultimately, our contribution is a call to shift the community’s focus from simplistic model rankings to a more insightful, diagnostic approach. This paradigm change fosters a deeper understanding of model behavior, enables fairer comparisons, and paves the way for the development of more robust and targeted forecasting solutions. Looking ahead, our framework opens the door to new data-aware learning paradigms. For instance, the SCP score can guide curriculum learning strategies by ordering data by difficulty, while the LUR could serve as a novel regularizer, paving the way for more adaptive and robust forecasting systems.

486 ETHICS STATEMENT  
487

488 This work does not introduce new data collection involving human subjects. All experiments use  
489 publicly available datasets (ETT, Weather, ECL) under their original licenses and terms of use; no  
490 sensitive attributes (e.g., protected classes or personally identifiable information) are inferred or  
491 used. Our analyses operate on aggregated forecasting errors and frequency-domain summaries, not  
492 on individual-level predictions tied to identity or behavior.  
493

494 REPRODUCIBILITY STATEMENT  
495

496 We aim for full reproducibility and have released all code, configurations, and scripts to repro-  
497 duce our results, including data preparation, backbone training/evaluation, and the SCP/LUR diag-  
498 nostics. All code and data can be found in [https://anonymous.4open.science/r/TS\\_](https://anonymous.4open.science/r/TS_Predictability-C8B7)  
499 [Predictability-C8B7](https://anonymous.4open.science/r/TS_Predictability-C8B7).  
500

501 REFERENCES  
502

503 Mateo Aboy, Roberto Hornero, Daniel Abásolo, and Daniel Álvarez. Interpretation of the Lempel-  
504 Ziv complexity measure in the context of biomedical signal analysis. *IEEE transactions on*  
505 *biomedical engineering*, 53(11):2282–2288, 2006.  
506

507 Christoph Bandt and Bernd Pompe. Permutation Entropy: A Natural Complexity Measure for Time  
508 Series. *Physical Review Letters*, 88(17):174102, April 2002. ISSN 0031-9007, 1079-7114. doi:  
509 10.1103/PhysRevLett.88.174102.  
510

511 Christoph Bergmeir. Fundamental limitations of foundational forecasting models: The need for  
512 multimodality and rigorous evaluation. In *Proc. NeurIPS Workshop*, 2024.

513 Xi Chen, Adityanand Guntuboyina, and Yuchen Zhang. On Bayes risk lower bounds. *Journal of*  
514 *Machine Learning Research*, 17(218):1–58, 2016.  
515

516 Andrew SC Ehrenberg and John A. Bound. Predictability and prediction. *Journal of the Royal*  
517 *Statistical Society Series A: Statistics in Society*, 156(2):167–194, 1993.  
518

519 Miro Erkintalo. Predicting the unpredictable? *Nature Photonics*, 9(9):560–562, 2015.  
520

521 Mark Fiecas, Chenlei Leng, Weidong Liu, and Yi Yu. Spectral analysis of high-dimensional time  
522 series. 2019.

523 Joshua Garland, Ryan James, and Elizabeth Bradley. Model-free quantification of time-series pre-  
524 dictability. *Physical Review E*, 90(5):052910, November 2014. ISSN 1539-3755, 1550-2376. doi:  
525 10.1103/PhysRevE.90.052910.  
526

527 Marta C. González, César A. Hidalgo, and Albert-László Barabási. Understanding individual human  
528 mobility patterns. *Nature*, 453(7196):779–782, June 2008. ISSN 1476-4687. doi: 10.1038/  
529 nature06958.

530 Jonathan Ko and Dieter Fox. GP-BayesFilters: Bayesian filtering using Gaussian process prediction  
531 and observation models. *Autonomous Robots*, 27(1):75–90, July 2009. ISSN 0929-5593, 1573-  
532 7527. doi: 10.1007/s10514-009-9119-x.  
533

534 Ioannis Kontoyiannis, Paul H. Algoet, Yu M. Suhov, and Abraham J. Wyner. Nonparametric entropy  
535 estimation for stationary processes and random fields, with applications to English text. *IEEE*  
536 *transactions on information theory*, 44(3):1319–1327, 2002.  
537

538 Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long-and short-term  
539 temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference*  
*on research & development in information retrieval*, pp. 95–104, 2018.

- 540 Shiyang Li, Xiaoyong Jin, Yao Xuan, Xiyong Zhou, Wenhui Chen, Yu-Xiang Wang, and Xifeng  
541 Yan. Enhancing the Locality and Breaking the Memory Bottleneck of Transformer on Time  
542 Series Forecasting. In *Advances in Neural Information Processing Systems*, volume 32. Curran  
543 Associates, Inc., 2019.
- 544 Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.  
545 iTransformer: Inverted Transformers Are Effective for Time Series Forecasting. In *The Thirteenth  
546 International Conference on Learning Representations*. The Thirteenth International Conference  
547 on Learning Representations, March 2024.
- 548 L Mandel and E Wolf. Spectral coherence and the concept of cross-spectral purity. *Journal of the  
549 Optical Society of America*, 66(6):529–535, 1976.
- 550 Sakshi Mishra and Praveen Palanisamy. Multi-time-horizon solar forecasting using recurrent neural  
551 network. In *2018 IEEE Energy Conversion Congress and Exposition (ECCE)*, pp. 18–24. IEEE,  
552 2018.
- 553 Jamal Mohammed, Michael H. Böhlen, and Sven Helmer. Quantifying and Estimating the Pre-  
554 dictability Upper Bound of Univariate Numeric Time Series. In *Proceedings of the 30th ACM  
555 SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '24*, pp. 2236–2247, New  
556 York, NY, USA, August 2024. Association for Computing Machinery. ISBN 979-8-4007-0490-1.  
557 doi: 10.1145/3637528.3671995.
- 558 Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A Time Series is Worth  
559 64 Words: Long-term Forecasting with Transformers. arXiv, March 2023. doi: 10.48550/arXiv.  
560 2211.14730.
- 561 Frank Pennekamp, Alison C. Iles, Joshua Garland, Georgina Brennan, Ulrich Brose, Ursula Gaedke,  
562 Ute Jacob, Pavel Kratina, Blake Matthews, Stephan Munch, Mark Novak, Gian Marco Palamara,  
563 Björn C. Rall, Benjamin Rosenbaum, Andrea Tabi, Colette Ward, Richard Williams, Hao Ye, and  
564 Owen L. Petchey. The intrinsic predictability of ecological time series and its potential to guide  
565 forecasting. *Ecological Monographs*, 89(2):e01359, May 2019. ISSN 0012-9615, 1557-7015.  
566 doi: 10.1002/ecm.1359.
- 567 S M Pincus. Approximate entropy as a measure of system complexity. *Proceedings of the National  
568 Academy of Sciences*, 88(6):2297–2301, March 1991. ISSN 0027-8424, 1091-6490. doi: 10.  
569 1073/pnas.88.6.2297.
- 570 Joshua S. Richman and J. Randall Moorman. Physiological time-series analysis using approximate  
571 entropy and sample entropy. *American Journal of Physiology-Heart and Circulatory Physiology*,  
572 278(6):H2039–H2049, June 2000. ISSN 0363-6135, 1522-1539. doi: 10.1152/ajpheart.2000.  
573 278.6.H2039.
- 574 C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27  
575 (3):379–423, July 1948. ISSN 0005-8580. doi: 10.1002/j.1538-7305.1948.tb01338.x.
- 576 Chaoming Song, Zehui Qu, Nicholas Blumm, and Albert-László Barabási. Limits of predictability in  
577 human mobility. *Science*, 327(5968):1018–1021, February 2010. doi: 10.1126/science.1177170.
- 578 Dong Wang, Xuejun Zhao, Lin-Lin Kou, Yong Qin, Yang Zhao, and Kwok-Leung Tsui. A simple  
579 and fast guideline for generating enhanced/squared envelope spectra from spectral coherence for  
580 bearing fault diagnosis. *Mechanical Systems and Signal Processing*, 122:754–768, 2019.
- 581 Huandong Wang, Sihan Zeng, Yong Li, and Depeng Jin. Predictability and Prediction of Human  
582 Mobility Based on Application-Collected Location Data. *IEEE Transactions on Mobile Comput-  
583 ing*, 20(7):2457–2472, July 2021. ISSN 1558-0660. doi: 10.1109/TMC.2020.2981441.
- 584 Yuxuan Wang, Haixu Wu, Jiayang Dong, Guo Qin, Haoran Zhang, Yong Liu, Yunzhong Qiu, Jian-  
585 min Wang, and Mingsheng Long. TimeXer: Empowering Transformers for Time Series Forecast-  
586 ing with Exogenous Variables. In *The Thirty-eighth Annual Conference on Neural Information  
587 Processing Systems*, November 2024.

594 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition trans-  
595 formers with auto-correlation for long-term series forecasting. *Advances in neural information*  
596 *processing systems*, 34:22419–22430, 2021.

597 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. TimesNet:  
598 Temporal 2D-Variation Modeling for General Time Series Analysis, April 2023.

600 Aaron D. Wyner and Jacob Ziv. Some asymptotic properties of the entropy of a stationary ergodic  
601 data source with applications to data compression. *IEEE Transactions on Information Theory*, 35  
602 (6):1250–1258, 2002.

603 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are Transformers Effective for Time Series  
604 Forecasting? *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(9):11121–11128,  
605 June 2023. ISSN 2374-3468. doi: 10.1609/aaai.v37i9.26317.

607 Kai Zhao, Denis Khryashchev, and Huy Vo. Predicting Taxi and Uber Demand in Cities: Approach-  
608 ing the Limit of Predictability. *IEEE Transactions on Knowledge and Data Engineering*, 33(6):  
609 2723–2736, June 2021. ISSN 1558-2191. doi: 10.1109/TKDE.2019.2955686.

611 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.  
612 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings*  
613 *of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11106–11115, 2021.

614 Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. Fedformer: Frequency  
615 enhanced decomposed transformer for long-term series forecasting. In *International Conference*  
616 *on Machine Learning*, pp. 27268–27286. PMLR, 2022.

617 J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions*  
618 *on Information Theory*, 23(3):337–343, May 1977. ISSN 1557-9654. doi: 10.1109/TIT.1977.  
619 1055714.  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

## A USAGE OF LLMs

We utilized a large language model (LLM) to proofread and improve the grammatical clarity of this manuscript. All scientific ideas, methodologies, and conclusions presented are the original work of the authors.

## B EXPERIMENTAL SETUP

### B.1 TOY STUDY

We synthesize a multiband Gaussian process and evaluate linear forecasting under controlled band-limited noise. Signals are split into history/future with boundary-paired segments of length  $N_p = 512$  from a total length  $2N$  with  $N = 1024$ . Power spectra and coherences are estimated via Welch’s method (Hann window) with  $n_{\text{perseg}} = 256$  and  $n_{\text{overlap}} = 128$ . The forecaster is a causal FIR least-squares filter (Wiener approximation) of length  $L_{\text{FIR}} = 64$  with ridge  $10^{-6}$ . The base process has four spectral peaks at rFFT bins  $\{32, 96, 192, 384\}$  with widths  $\{6, 10, 14, 18\}$  and amplitudes  $\{3.0, 2.0, 1.5, 1.0\}$ . We sweep noise levels  $\{0, 0.25, 0.5, 1.0, 2.0, 4.0\}$  on a single band (index 1 in the plot) and average over 3 trials, reporting model MSE,  $\text{MSE}_{\text{lb}}$ , and SCP.

### B.2 BACKBONE

We evaluate seven state-of-the-art backbones spanning diverse architectures: Transformer-based (iTransformer (Liu et al., 2024), PatchTST (Nie et al., 2023)), MLP-based (DLinear (Zeng et al., 2023), TimeMixer (Wang et al., 2024)), and CNN-based (TimesNet (Wu et al., 2023)). We adopt the official implementations and recommended hyperparameters from their repositories. To ensure strict comparability, we fix the forecasting horizon and enforce equal input and output lengths for all backbones (no “drop-last”), using identical preprocessing and dataset splits across models.

### B.3 DATASETS

We conduct experiments on eight standard long-horizon multivariate forecasting benchmarks: ETTh1, ETTh2, ETTm1, ETTm2, ECL, Weather, Traffic, and ILI. These datasets cover electricity systems, meteorology, transportation, and epidemiology, and are widely used in recent long-horizon time series forecasting studies.

- **ETT (Electricity Transformer Temperature).** The ETT benchmarks record seven variables from electricity transformers collected in two regions between July 2016 and July 2018. Variants differ in sampling granularity: “h” denotes hourly data and “m” denotes 15-minute data; suffixes “1/2” indicate the two regions (ETTh1/ETTh2, ETTm1/ETTM2) (Zhou et al., 2021). Following standard practice, we treat each variant as a separate multivariate dataset.
- **ECL (Electricity Consuming Load).** ECL provides hourly electricity consumption for 321 clients from 2012 to 2014 (Li et al., 2019). It is a large-scale dataset with strong daily and weekly periodicities and substantial cross-series correlations.
- **Weather.** Weather contains 21 meteorological variables (for example, temperature, humidity, wind speed) measured every 10 minutes in 2020 at the Max Planck Biogeochemistry Institute’s weather station (Zhou et al., 2021). It provides a medium-scale, high-frequency benchmark with strong diurnal and seasonal patterns.
- **Traffic.** Traffic consists of hourly road occupancy rates collected from 862 loop sensors on San Francisco Bay Area freeways over two years (Wu et al., 2021). The series are bounded in  $[0, 1]$  and exhibit strong rush-hour and weekday/weekend patterns, making the dataset a high-dimensional transportation benchmark.
- **ILI (Influenza-Like Illness).** ILI contains weekly statistics on influenza-like illness in the United States, including the number of patients and the ILI ratio aggregated over multiple regions (Lai et al., 2018). Compared with the other datasets, ILI is much shorter and lower frequency, with pronounced annual seasonality and relatively high noise, and is therefore evaluated with shorter prediction horizons.

**Algorithm 3** Multivariate Spectral Coherence Predictability ( $\text{SCP}_{\text{multi}}$ )

**Require:** History  $\mathbf{x} \in \mathbb{R}^{d_x \times N}$ , future  $\mathbf{y} \in \mathbb{R}^{d_y \times N}$ ; Welch parameters (window, length, overlap); stability constant  $\varepsilon > 0$ ; optional frequency band  $\mathcal{F}_b$ .

**Ensure:** Multivariate MSE lower bound  $\text{MSE}_{\text{lb}}^{\text{multi}}$  and predictability  $\mathcal{P}_{xy}^{\text{multi}}$ .

1: **Mean removal:**  $\Delta^2 \leftarrow \|\boldsymbol{\mu}_y - \boldsymbol{\mu}_x\|_2^2$ ;  $\mathbf{x} \leftarrow \mathbf{x} - \boldsymbol{\mu}_x$ ,  $\mathbf{y} \leftarrow \mathbf{y} - \boldsymbol{\mu}_y$ .

2: **Welch spectra:** Compute matrix-valued PSDs  $\widehat{S}_{xx}(f)$ ,  $\widehat{S}_{yy}(f)$  and CPSD  $\widehat{S}_{xy}(f)$  on  $\mathcal{F}$ ; set  $\widehat{S}_{yx}(f) = \widehat{S}_{xy}(f)^H$ .

3: **Multichannel Wiener spectra:**

$$\widehat{S}_{\hat{y}\hat{y}}(f) = \widehat{S}_{yy}(f)(\widehat{S}_{xx}(f) + \varepsilon I_{d_x})^{-1}\widehat{S}_{xy}(f), \quad \widehat{S}_e(f) = \widehat{S}_{yy}(f) - \widehat{S}_{\hat{y}\hat{y}}(f).$$

4: **Frequency set:**  $\mathcal{F}_* \leftarrow \mathcal{F}_b$  if  $\mathcal{F}_b$  is provided; otherwise  $\mathcal{F}_* \leftarrow \mathcal{F}$ .

5: **Aggregate:**

$$\widehat{\text{Var}}(\mathbf{y}) \leftarrow \sum_{f \in \mathcal{F}_*} \text{tr} \widehat{S}_{yy}(f), \quad \text{MSE}_{\text{lb}}^{\text{multi}} \leftarrow \Delta^2 + \sum_{f \in \mathcal{F}_*} \text{tr} \widehat{S}_e(f).$$

6: **Predictability:**  $\mathcal{P}_{xy}^{\text{multi}} \leftarrow 1 - \text{MSE}_{\text{lb}}^{\text{multi}} / \widehat{\text{Var}}(\mathbf{y})$ .

7: **return**  $\text{MSE}_{\text{lb}}^{\text{multi}}$ ,  $\mathcal{P}_{xy}^{\text{multi}}$ .

## B.4 TIME-TO-FREQUENCY

For each test instance we split the sequence into history  $\mathbf{x}$  and future  $\mathbf{y}$  (equal lengths by default), remove sample means, and estimate power and cross-spectra with Welch’s method using identical settings for  $\mathbf{x}$ ,  $\mathbf{y}$ , and (when available)  $\hat{\mathbf{y}}$ : Hann window with length  $n_{\text{win}} = \lfloor 0.25N \rfloor$ , 50% overlap, and real FFT on the one-sided grid  $\mathcal{F}$  with variance-preserving normalization. We form squared coherences with a small ridge  $\varepsilon$  for stability, compute the residual spectrum to obtain the linear lower bound  $\text{MSE}_{\text{lb}}$  and predictability  $\mathcal{P} = 1 - \text{MSE}_{\text{lb}} / \widehat{\text{Var}}(\mathbf{y})$ , and derive utilization metrics (global or band-wise) via target-power-weighted aggregation of  $\gamma_{\hat{y}\hat{y}}^2$  and  $\gamma_{yx}^2$ .

## C METHOD EXTENSIONS

## C.1 MULTIVARIATE EXTENSION

## C.1.1 MULTIVARIATE SCP

We extend the univariate SCP in Sec. 4.1 to multivariate histories and futures with input dimensionality  $d_x$  and output dimensionality  $d_y$ . Let  $\mathbf{x}_t \in \mathbb{R}^{d_x}$  and  $\mathbf{y}_t \in \mathbb{R}^{d_y}$  denote a length- $N$  history–future pair, and let  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^{d_x \times N}$ ,  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_N) \in \mathbb{R}^{d_y \times N}$ . Using Welch’s method with shared parameters for all components, we compute multivariate power spectral density (PSD) and cross-power spectral density (CPSD) matrices on a discrete frequency grid  $\mathcal{F}$ :

$$\widehat{S}_{xx}(f) \in \mathbb{C}^{d_x \times d_x}, \quad \widehat{S}_{yy}(f) \in \mathbb{C}^{d_y \times d_y}, \quad \widehat{S}_{xy}(f) \in \mathbb{C}^{d_x \times d_y}, \quad (12)$$

and set  $\widehat{S}_{yx}(f) = \widehat{S}_{xy}(f)^H$ , where  $H$  denotes the Hermitian transpose.

At frequency  $f$ , the optimal linear time-invariant predictor from  $\mathbf{x}$  to  $\mathbf{y}$  in the least-squares sense has transfer matrix

$$H(f) = \widehat{S}_{yx}(f)(\widehat{S}_{xx}(f) + \varepsilon I_{d_x})^{-1}, \quad (13)$$

where  $\varepsilon > 0$  is the same Tikhonov regularization as in Eq. (4), and  $I_{d_x}$  is the  $d_x \times d_x$  identity matrix. The spectrum of the linearly predictable component of  $\mathbf{y}$  is then

$$\widehat{S}_{\hat{y}\hat{y}}(f) = H(f)\widehat{S}_{xx}(f)H(f)^H = \widehat{S}_{yx}(f)(\widehat{S}_{xx}(f) + \varepsilon I_{d_x})^{-1}\widehat{S}_{xy}(f) \in \mathbb{C}^{d_y \times d_y}. \quad (14)$$

In the scalar case  $d_x = d_y = 1$ , Eq. (14) reduces to  $\widehat{S}_{\hat{y}\hat{y}}(f) = |\widehat{S}_{xy}(f)|^2 / (\widehat{S}_{xx}(f) + \varepsilon)$ , which coincides with the univariate expression  $\gamma_{yx}^2(f) \widehat{S}_{yy}(f)$  in Eq. (4).

**Algorithm 4** Multivariate Linear Utilization Ratio (LUR<sub>multi</sub>)

**Require:** History  $\mathbf{x} \in \mathbb{R}^{d_x \times N}$ , future  $\mathbf{y} \in \mathbb{R}^{d_y \times N}$ , prediction  $\hat{\mathbf{y}} \in \mathbb{R}^{d_y \times N}$ ; Welch parameters (window, length, overlap); stability  $\varepsilon > 0$ ; optional band  $\mathcal{F}_b$ .

**Ensure:** Multivariate model-explained power  $P_{\text{model}}$ , linear-explainable power  $P_{\text{linear}}$ , and utilization ratio LUR<sup>multi</sup>.

1: **Mean removal:**  $\mathbf{x} \leftarrow \mathbf{x} - \text{mean}(\mathbf{x})$ ;  $\mathbf{y} \leftarrow \mathbf{y} - \text{mean}(\mathbf{y})$ ;  $\hat{\mathbf{y}} \leftarrow \hat{\mathbf{y}} - \text{mean}(\hat{\mathbf{y}})$ .

2: **Welch spectra:** Compute  $\hat{S}_{xx}(f)$ ,  $\hat{S}_{yy}(f)$ ,  $\hat{S}_{\hat{y}\hat{y}}^{\text{pred}}(f)$ ,  $\hat{S}_{xy}(f)$ ,  $\hat{S}_{y\hat{y}}(f)$  on  $\mathcal{F}$ ; set  $\hat{S}_{yx}(f) = \hat{S}_{xy}(f)^H$  and  $\hat{S}_{\hat{y}y}(f) = \hat{S}_{y\hat{y}}(f)^H$ .

3: **Linear limit (per frequency):**

$$\hat{S}_{\hat{y}\hat{y}}(f) \leftarrow \hat{S}_{yx}(f)(\hat{S}_{xx}(f) + \varepsilon I_{d_x})^{-1}\hat{S}_{xy}(f), \quad P_{\text{linear}}(f) \leftarrow \text{tr} \hat{S}_{\hat{y}\hat{y}}(f).$$

4: **Model-explained power (per frequency):**

$$P_{\text{model}}(f) \leftarrow \text{tr} \left( \hat{S}_{\hat{y}\hat{y}}(f)(\hat{S}_{\hat{y}\hat{y}}^{\text{pred}}(f) + \varepsilon I_{d_y})^{-1}\hat{S}_{\hat{y}\hat{y}}(f) \right).$$

5: **Frequency set:**  $\mathcal{F}_* \leftarrow \mathcal{F}_b$  if a band  $\mathcal{F}_b$  is provided; otherwise  $\mathcal{F}_* \leftarrow \mathcal{F}$ .

6: **Aggregation:**

$$P_{\text{linear}} \leftarrow \sum_{f \in \mathcal{F}_*} P_{\text{linear}}(f), \quad P_{\text{model}} \leftarrow \sum_{f \in \mathcal{F}_*} P_{\text{model}}(f).$$

7: **LUR ratio:** LUR<sup>multi</sup>  $\leftarrow P_{\text{model}}/P_{\text{linear}}$ .

8: **return**  $P_{\text{model}}$ ,  $P_{\text{linear}}$ , LUR<sup>multi</sup>.

The residual spectrum matrix is

$$\hat{S}_e(f) = \hat{S}_{yy}(f) - \hat{S}_{\hat{y}\hat{y}}(f), \quad \forall f \in \mathcal{F}. \quad (15)$$

Since  $\hat{S}_{\hat{y}\hat{y}}(f)$  is the least-squares projection of  $\hat{S}_{yy}(f)$  onto the subspace linearly spanned by  $\mathbf{x}$ , the true residual spectrum is positive semidefinite, and the regularization  $\varepsilon I_{d_x}$  stabilizes this property numerically. Let the estimated total variance (total power) of  $\mathbf{y}$  be the trace-aggregated spectrum

$$\widehat{\text{Var}}(\mathbf{y}) = \sum_{f \in \mathcal{F}} \text{tr} \hat{S}_{yy}(f), \quad (16)$$

where  $\text{tr}(\cdot)$  denotes the matrix trace. Using the same frequency grid, the multivariate MSE lower bound induced by linear time-invariant predictors is

$$\text{MSE}_{\text{lb}}^{\text{multi}} = \Delta^2 + \sum_{f \in \mathcal{F}} \text{tr} \hat{S}_e(f), \quad (17)$$

where  $\Delta^2$  is the same boundary mean-shift term as in the univariate case, generalized to the  $(d_x, d_y)$ -dimensional setting.

The multivariate SCP is defined by normalizing the residual energy as in Eq. (7):

$$\mathcal{P}_{xy}^{\text{multi}} = 1 - \frac{\text{MSE}_{\text{lb}}^{\text{multi}}}{\widehat{\text{Var}}(\mathbf{y})} \in [0, 1]. \quad (18)$$

When  $d_x = d_y = 1$ , Eq. (18) reduces exactly to the univariate SCP in Eq. (7).

### C.1.2 MULTIVARIATE LUR

The spectrum of the linearly predictable component in Eq. (14) induces the linear-explainable power

$$P_{\text{linear}}(f) = \text{tr} \hat{S}_{\hat{y}\hat{y}}(f) = \text{tr} \left( \hat{S}_{yx}(f)(\hat{S}_{xx}(f) + \varepsilon I_{d_x})^{-1}\hat{S}_{xy}(f) \right). \quad (19)$$

For the model, we form the auto- and cross-spectra of the prediction,

$$\hat{S}_{\hat{y}\hat{y}}^{\text{pred}}(f) \in \mathbb{C}^{d_y \times d_y}, \quad \hat{S}_{y\hat{y}}(f) \in \mathbb{C}^{d_y \times d_y}, \quad \hat{S}_{\hat{y}y}(f) = \hat{S}_{y\hat{y}}(f)^H, \quad (20)$$

and define the model–explained power via the optimal linear projection of  $\mathbf{y}$  onto the subspace spanned by  $\hat{\mathbf{y}}$ :

$$P_{\text{model}}(f) = \text{tr} \left( \widehat{S}_{\hat{\mathbf{y}}\hat{\mathbf{y}}}(f) \left( \widehat{S}_{\hat{\mathbf{y}}\hat{\mathbf{y}}}^{\text{pred}}(f) + \varepsilon I_{d_{\hat{\mathbf{y}}}} \right)^{-1} \widehat{S}_{\hat{\mathbf{y}}\mathbf{y}}(f) \right). \quad (21)$$

Aggregating over the discrete frequency domain  $\mathcal{F}$ ,

$$P_{\text{linear}} = \sum_{f \in \mathcal{F}} P_{\text{linear}}(f), \quad P_{\text{model}} = \sum_{f \in \mathcal{F}} P_{\text{model}}(f), \quad (22)$$

and normalizing as in Sec. 4.2 gives the multivariate linear utilization ratio

$$\text{LUR}^{\text{multi}} = \frac{P_{\text{model}}}{P_{\text{linear}}}. \quad (23)$$

When  $d_x = d_y = 1$ , these expressions reduce to the univariate definitions of  $P_{\text{linear}}$ ,  $P_{\text{model}}$ , and LUR.

## C.2 NONLINEAR EXTENSION

The SCP framework is linear by construction: it characterizes the best linear time–invariant (LTI) predictor in the original observation space. To relax this restriction while preserving the same spectral machinery, we introduce a nonlinear feature map

$$\phi : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_z}, \quad \mathbf{z}_t = \phi(\mathbf{x}_t) \in \mathbb{R}^{d_z}, \quad (24)$$

and apply multivariate SCP in the resulting feature space. We then form the feature sequence  $\mathbf{z} = (\mathbf{z}_1, \dots, \mathbf{z}_N) \in \mathbb{R}^{d_z \times N}$ . The map  $\phi$  can use explicit nonlinear features (e.g., polynomial expansions or a shallow encoder), or be defined implicitly by a kernel  $k(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$  in an RKHS. Using the same Welch configuration as before, we estimate the multivariate spectra

$$\widehat{S}_{\mathbf{z}\mathbf{z}}(f) \in \mathbb{C}^{d_z \times d_z}, \quad \widehat{S}_{\mathbf{y}\mathbf{y}}(f) \in \mathbb{C}^{d_y \times d_y}, \quad \widehat{S}_{\mathbf{y}\mathbf{z}}(f) \in \mathbb{C}^{d_y \times d_z}, \quad (25)$$

and set  $\widehat{S}_{\mathbf{z}\mathbf{y}}(f) = \widehat{S}_{\mathbf{y}\mathbf{z}}(f)^H$ .

In feature space, the optimal LTI predictor of  $\mathbf{y}$  from  $\mathbf{z}$  takes the same form as the multivariate Wiener filter in Eq. (13), but with  $(\mathbf{x}, \widehat{S}_{\mathbf{x}\mathbf{x}})$  replaced by  $(\mathbf{z}, \widehat{S}_{\mathbf{z}\mathbf{z}})$ :

$$H_\phi(f) = \widehat{S}_{\mathbf{y}\mathbf{z}}(f) \left( \widehat{S}_{\mathbf{z}\mathbf{z}}(f) + \varepsilon I_{d_z} \right)^{-1}, \quad (26)$$

where  $\varepsilon > 0$  is the same Tikhonov regularization as before and  $I_{d_z}$  is the  $d_z \times d_z$  identity. The spectrum of the component of  $\mathbf{y}$  that is linearly predictable from the nonlinear features is

$$\widehat{S}_{\hat{\mathbf{y}}\hat{\mathbf{y}}}^{\text{ker}}(f) = H_\phi(f) \widehat{S}_{\mathbf{z}\mathbf{z}}(f) H_\phi(f)^H = \widehat{S}_{\mathbf{y}\mathbf{z}}(f) \left( \widehat{S}_{\mathbf{z}\mathbf{z}}(f) + \varepsilon I_{d_z} \right)^{-1} \widehat{S}_{\mathbf{z}\mathbf{y}}(f) \in \mathbb{C}^{d_y \times d_y}. \quad (27)$$

When  $\phi$  is the identity map ( $d_z = d_x$  and  $\mathbf{z}_t = \mathbf{x}_t$ ), Eq. (27) reduces to the multivariate linear spectrum in Eq. (14).

The residual spectrum under the feature-space predictor is

$$\widehat{S}_e^{\text{ker}}(f) = \widehat{S}_{\mathbf{y}\mathbf{y}}(f) - \widehat{S}_{\hat{\mathbf{y}}\hat{\mathbf{y}}}^{\text{ker}}(f), \quad \forall f \in \mathcal{F}, \quad (28)$$

which is positive semidefinite in the ideal (population) setting. Aggregating as in Eq. (16), the total variance of  $\mathbf{y}$  and the corresponding nonlinear MSE lower bound are

$$\widehat{\text{Var}}(\mathbf{y}) = \sum_{f \in \mathcal{F}} \text{tr} \widehat{S}_{\mathbf{y}\mathbf{y}}(f), \quad \text{MSE}_{\text{lb}}^{\text{ker}} = \Delta^2 + \sum_{f \in \mathcal{F}} \text{tr} \widehat{S}_e^{\text{ker}}(f), \quad (29)$$

where  $\Delta^2$  is the same boundary mean–shift term used in Eq. (17), applied to the multivariate setting.

The nonlinear SCP is then obtained by normalizing the feature-space residual:

$$\mathcal{P}_{xy}^{\text{nonlin}} = 1 - \frac{\text{MSE}_{\text{lb}}^{\text{ker}}}{\widehat{\text{Var}}(\mathbf{y})}. \quad (30)$$

This quantity measures the fraction of future variance that is explainable by LTI predictors acting on the chosen nonlinear feature representation, providing a feature-dependent notion of nonlinear predictability.

### 864 C.3 BEYOND EVALUATION

865 The predictability scores from SCP (and its multivariate / nonlinear variants)  $\mathcal{P}_{xy} \in [0, 1]$  can be  
866 used not only for post-hoc analysis, but also to shape how data are selected and organized during  
867 training.

868  
869 **Hard-example mining** For a dataset  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^M$  with per-sample predictability  $\mathcal{P}^{(i)} \equiv$   
870  $\mathcal{P}_{xy}(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ , we can directly use  $\mathcal{P}^{(i)}$  to reweight the loss:

$$871 L = \sum_{i=1}^M w^{(i)} \ell(\hat{\mathbf{y}}^{(i)}, \mathbf{y}^{(i)}), \quad w^{(i)} \propto (\mathcal{P}^{(i)})^\alpha, \quad (31)$$

872 with  $\alpha < 0$ . This up-weights intrinsically predictable segments and down-weights near-  
873 unpredictable ones that mainly contain irreducible noise.

874  
875 **Curriculum learning** The same scores induce a simple curriculum over data difficulty. Let  
876  $\{\tau_s\}_{s=1}^S$  be a decreasing sequence of thresholds,  $\tau_1 > \tau_2 > \dots > \tau_S$ . At stage  $s$ , we restrict  
877 training to

$$878 \mathcal{D}_s = \{i : \mathcal{P}^{(i)} \geq \tau_s\}, \quad (32)$$

879 i.e., the model first sees highly predictable segments, and progressively incorporates samples with  
880 lower  $\mathcal{P}^{(i)}$  as  $s$  increases.

881  
882 **Anomaly detection and change points** On a time series stream, we compute predictability over  
883 a sliding window ending at time  $t$ , for example  $\mathcal{P}_t = \mathcal{P}_{xy}(\mathbf{x}_{t-N+1:t}, \mathbf{y}_{t+1:t+N})$ . Let  $\mu_{\mathcal{P}}, \sigma_{\mathcal{P}}$  be the  
884 mean and standard deviation of  $\mathcal{P}_t$  on a reference (normal) period. We flag  $t$  as anomalous when

$$885 |\mathcal{P}_t - \mu_{\mathcal{P}}| > \kappa \sigma_{\mathcal{P}}, \quad (33)$$

886 with a chosen threshold  $\kappa > 0$ . Sudden drops or spikes in  $\mathcal{P}_t$  indicate changes in intrinsic pre-  
887 dictability, and thus potential regime shifts or anomalous behavior.

### 888 C.4 COMPARISON WITH TIME-DOMAIN CORRELATION DIAGNOSTICS

889  
890 Classical time-domain tools such as the autocorrelation function (ACF) provide a convenient way  
891 to visualize second-order structure by plotting correlation as a function of lag. ACF is particularly  
892 useful for qualitatively assessing periodicity and dependence decay. However, it is primarily a de-  
893 scriptive tool for the self-correlation of a single series. In particular, ACF does not directly quantify  
894 how well a future window can be linearly predicted from a past window under the MSE objective,  
895 especially in the multi-horizon, multivariate setting we consider.

896 In contrast, our SCP/LUR framework is explicitly constructed around the past–future prediction  
897 task. SCP is derived from the cross-spectral density and coherence between the history and future  
898 segments, and measures the fraction of the future variance that is linearly explainable from the  
899 observed history, yielding an MSE-aligned notion of intrinsic predictability. LUR further decom-  
900 poses this explainable energy across frequency bands, revealing which parts of the spectrum are well  
901 captured or systematically missed by a given model.

902 A simple example illustrates the difference between naive time-domain correlation and spectral  
903 coherence. Consider two noiseless signals

$$904 x_t = \sin(\omega_0 t), \quad y_t = \cos(\omega_0 t).$$

905 Here  $y_t$  is a phase-shifted version of  $x_t$ , obtained by a linear time-invariant transformation. In other  
906 words,  $y$  is perfectly linearly predictable from  $x$ .

907 If we look only at the zero-lag Pearson correlation

$$908 \rho_{xy}(0) = \text{corr}(x_t, y_t),$$

909 and average over many periods by treating  $\theta = \omega_0 t$  as uniform on  $[0, 2\pi]$ , we obtain

$$910 \mathbb{E}[\sin \theta \cos \theta] = 0,$$

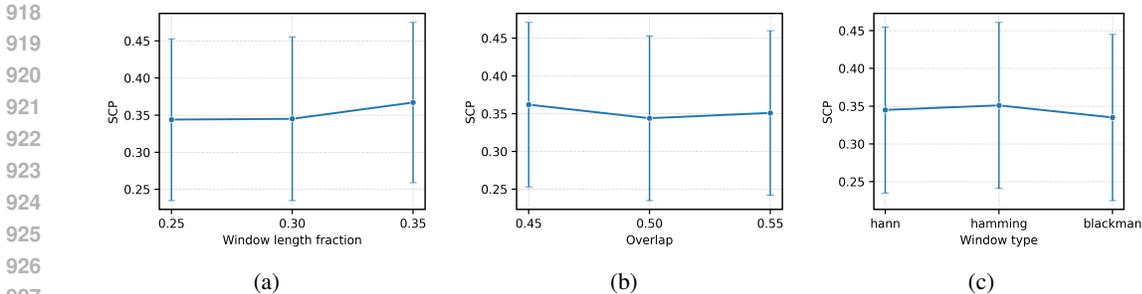


Figure 6: Sensitivity of SCP to Welch hyperparameters on the Weather dataset. Error bars indicate the mean and standard deviation computed over 10000 samples.

hence  $\rho_{xy}(0) = 0$ . A time-domain diagnostic based solely on zero-lag correlation would therefore suggest that  $x$  and  $y$  are “unrelated”, even though  $y$  is deterministically generated from  $x$  by a linear filter.

In the frequency domain, both  $x$  and  $y$  have all their energy concentrated at the same frequency  $\omega_0$ . Their cross-spectrum at  $\omega_0$  differs only by a constant phase factor, so the squared coherence

$$\gamma^2(\omega_0) = \frac{|S_{xy}(\omega_0)|^2}{S_{xx}(\omega_0) S_{yy}(\omega_0)}$$

evaluates to  $\gamma^2(\omega_0) = 1$ . In our framework, this implies a linear MSE lower bound of zero and SCP equal to one: the spectral diagnostic correctly recognizes that  $y$  is fully predictable from  $x$  despite the phase shift. This example highlights that simple time-domain summaries such as zero-lag correlation can miss strong linear predictability when phase shifts or distributed lags are present, whereas coherence (and thus SCP) aggregates information over all lags at each frequency and is invariant to such shifts.

## D SUPPLEMENTARY EXPERIMENTS

### D.1 SENSITIVITY ANALYSIS

This section presents two complementary sensitivity studies. First, we investigate how the SCP metric responds to the hyperparameters of the Welch coherence estimator. Second, we analyze whether conclusions drawn from LUR remain consistent under different frequency-band partitions.

#### D.1.1 SENSITIVITY TO WELCH PARAMETERS

We assess the sensitivity of SCP to the Welch coherence hyperparameters. On the Weather dataset, with the forecast horizon fixed to  $N = 96$ , we vary three factors: the window-length fraction  $L_w/N$ , the overlap ratio  $\rho$ , and the tapering window (Hann, Hamming, Blackman). We adopt  $(L_w/N, \rho, \text{window}) = (0.25, 0.5, \text{Hann})$  as the default configuration, and vary one hyperparameter at a time while keeping the others fixed.

As shown in Figure 6 and Table 2, the mean SCP changes only slightly across the entire parameter range, indicating that the estimator remains stable under standard spectral settings. The relatively large standard deviations stem from the data rather than the estimator: different temporal regions exhibit distinct intrinsic predictability, consistent with the non-stationary structure illustrated in Figure 3.

#### D.1.2 SENSITIVITY TO FREQUENCY-BAND PARTITIONING

As illustrated in Figures 7 and 8, changing the band boundaries affects the absolute LUR values within each band, but the qualitative conclusions remain unchanged. Across all configurations, iTransformer consistently achieves higher LUR in the low-frequency region where most signal energy concentrates, whereas DLinear performs better in the high-frequency bands. The frequency

Table 2: SCP and linear MSE lower bound ( $\text{MSE}_{\text{lb}}$ ) under different Welch configurations.

Parameter	Value	SCP (mean $\pm$ std)	$\text{MSE}_{\text{lb}}$ (mean $\pm$ std)
Window-length fraction ( $L_w/N$ )	0.25	$0.344 \pm 0.109$	$0.186 \pm 0.133$
	0.30	$0.345 \pm 0.110$	$0.186 \pm 0.133$
	0.35	$0.367 \pm 0.108$	$0.183 \pm 0.133$
Overlap ( $\rho$ )	0.45	$0.362 \pm 0.109$	$0.184 \pm 0.134$
	0.50	$0.344 \pm 0.109$	$0.185 \pm 0.133$
	0.55	$0.351 \pm 0.109$	$0.185 \pm 0.135$
Window type	Hann	$0.345 \pm 0.110$	$0.186 \pm 0.134$
	Hamming	$0.351 \pm 0.110$	$0.185 \pm 0.133$
	Blackman	$0.335 \pm 0.110$	$0.187 \pm 0.135$

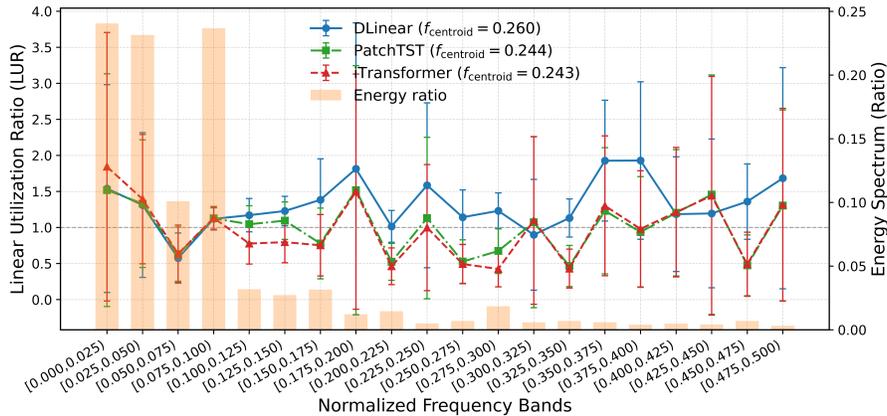


Figure 7: Band-wise normalized energy and LUR on ETTh1 using a 20-band partition.

centroid  $f_{\text{centroid}}$  exhibits the same ordering: DLinear attains the largest centroid, while PatchTST and iTransformer remain close.

## D.2 MULTIVARIATE PREDICTABILITY

To validate that our metric captures multivariate predictability, we construct a controlled synthetic example. The input is a  $d_x$ -dimensional process  $x(n) \in \mathbb{R}^{d_x}$  and the target is scalar ( $d_y = 1$ ). We set  $d_x = 6$ ,  $d_y = 1$ , sequence length  $N = 1024$ , and number of independent sequences  $N_{\text{samples}} = 640$ .

For each input dimension  $i \in \{1, \dots, d_x\}$  we generate a sinusoid  $x_i(n) = \sin(2\pi f_i n + \phi_i)$  for  $n = 0, \dots, N - 1$ , with distinct frequencies  $f_i = k_i/L_w$  for a Welch window length  $L_w = 128$  and  $(k_1, \dots, k_6) = (3, 5, 7, 9, 11, 13)$ , aligned to discrete Fourier bins. The phases are drawn independently as  $\phi_i \sim \text{Unif}[0, 2\pi)$  for each  $i$  and each sequence. The target signal is defined as a noisy sum of all input components,  $y(n) = \sum_{i=1}^{d_x} x_i(n) + \epsilon(n)$ , with  $\epsilon(n) \sim \mathcal{N}(0, 0.05)$ , so that most of the target energy is linearly generated by the  $d_x$  inputs.

For each  $m \in \{1, \dots, d_x\}$  we only reveal the first  $m$  input dimensions  $(x_1, \dots, x_m)$  and compute the resulting multivariate SCP  $\mathcal{P}_{\text{lin}}^{\text{multi}}(m)$  and multivariate MSE lower bound  $\text{MSE}_{\text{lb}}^{\text{multi}}(m)$ . Both quantities are averaged over the  $N_{\text{samples}}$  sequences, and we also report their empirical standard deviations. The numerical results are summarized in Table 3, and the corresponding curves are shown in Fig. 9.

The results exhibit a clear, approximately monotonic trend. As the number of observed input dimensions  $m$  increases, the multivariate SCP  $\mathcal{P}_{\text{lin}}^{\text{multi}}(m)$  rises almost linearly, while the multivariate  $\text{MSE}_{\text{lb}}^{\text{multi}}(m)$  decreases accordingly. As  $m$  approaches  $d_x$ ,  $\mathcal{P}_{\text{lin}}^{\text{multi}}(m)$  approaches the ideal predictability implied by the signal-to-noise ratio (but does not reach 1 due to spectral estimation error and the injected noise), and  $\text{MSE}_{\text{lb}}^{\text{multi}}(m)$  correspondingly approaches zero. Taken together, the

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

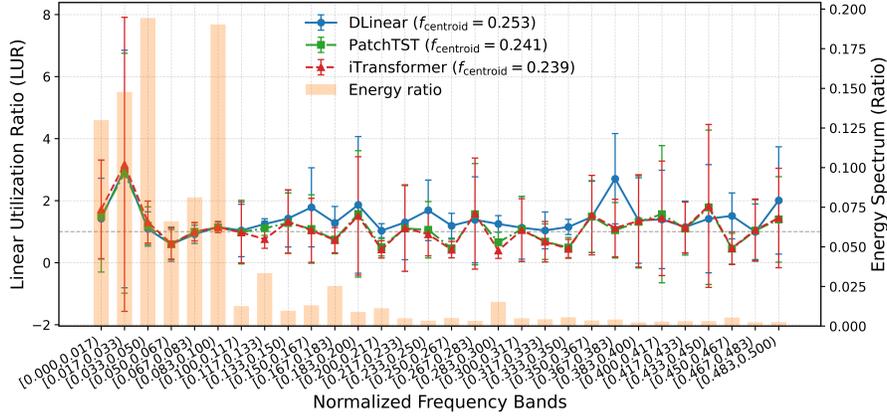


Figure 8: Band-wise normalized energy and LUR on ETTh1 using a 30-band partition.

Table 3: Multivariate SCP and MSE lower bound versus the number of observed input dimensions  $m$  in the synthetic sinusoid mixture experiment.

$m$	$\mathcal{P}_{\text{lin}}^{\text{multi}}$ (mean $\pm$ std)	$\text{MSE}_{\text{lb}}^{\text{multi}}$ (mean $\pm$ std)
1	0.117 $\pm$ 0.121	0.486 $\pm$ 0.067
2	0.208 $\pm$ 0.145	0.436 $\pm$ 0.080
3	0.360 $\pm$ 0.130	0.353 $\pm$ 0.072
4	0.523 $\pm$ 0.169	0.263 $\pm$ 0.094
5	0.662 $\pm$ 0.151	0.186 $\pm$ 0.083
6	0.848 $\pm$ 0.030	0.084 $\pm$ 0.017

table and figure confirm that multivariate SCP faithfully tracks the gain in predictability contributed by additional informative input dimensions.

### D.3 VARIABLE HISTORY WINDOW ( $N_x \neq N_y$ )

Let  $N_x$  and  $N_y$  denote the history and future lengths used for SCP. The construction only requires that a contiguous history–future pair exists around the boundary;  $N_x$  and  $N_y$  need not coincide. Given a Welch segment length  $L_w$  and overlap ratio  $\text{overlap} \in [0, 1)$ , the effective shift between consecutive segments is

$$\Delta = L_w (1 - \text{overlap}),$$

and the approximate number of Welch segments for a sequence of length  $N$  is

$$K(N; L_w, \text{overlap}) \approx \left\lfloor \frac{N - L_w}{L_w(1 - \text{overlap})} \right\rfloor + 1. \quad (34)$$

As a concrete example, consider a history window  $N_x = 192$ , a longer future horizon  $N_y = 336$ , and a Welch window  $L_w = 64$ . With an overlap of  $\text{overlap} = 0.5$ , the hop size is  $\Delta = 64(1 - 0.5) = 32$ , and the corresponding numbers of Welch segments are

$$K_x \approx K(192; 64, 0.5) = \left\lfloor \frac{192 - 64}{32} \right\rfloor + 1 = 5, \quad K_y \approx K(336; 64, 0.5) = \left\lfloor \frac{336 - 64}{32} \right\rfloor + 1 = 9.$$

For each segment we form windowed signals  $\mathbf{x}_k(t)$  and  $\mathbf{y}_k(t)$  of length  $L_w$ , compute their discrete Fourier transforms  $X_k(f)$  and  $Y_k(f)$ , and define the auto-spectra by Welch averaging

$$\hat{S}_{xx}(f) = \frac{1}{K_x} \sum_{k=1}^{K_x} |X_k(f)|^2, \quad \hat{S}_{yy}(f) = \frac{1}{K_y} \sum_{k=1}^{K_y} |Y_k(f)|^2.$$

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

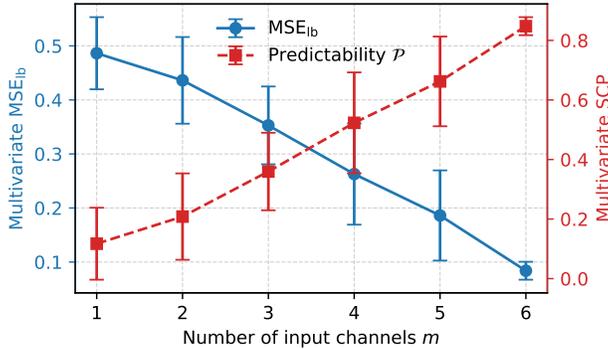


Figure 9: Multivariate SCP and MSE lower bound versus the number of input dimensions  $m$  in the synthetic sinusoid mixture experiment. Error bars indicate standard deviation across  $N_{\text{samples}}$  independent sequences.

Table 4: Effect of history window length  $N_x$  on model error (MSE,  $MSE_{lb}$ ) and correlation  $R$  on ETTh1 (prediction horizon  $N_y = 336$ ).

Model	Metric	History length $N_x$			
		96	192	336	720
iTransformer	MSE	0.491	0.479	0.471	0.480
	$R$	0.808	0.842	0.899	0.832
DLinear	MSE	0.491	0.480	0.447	0.449
	$R$	0.788	0.840	0.872	0.853
Predictability	$MSE_{lb}$	0.431	0.433	0.404	0.424

The cross-spectrum is computed on the aligned history–future portion at the boundary: we use the last  $K_{\text{pair}} = \min(K_x, K_y)$  segments from the history and the first  $K_{\text{pair}}$  segments from the future, denote their transforms by  $X_k^{(\text{hist})}(f)$  and  $Y_k^{(\text{fut})}(f)$ , and set

$$\hat{S}_{xy}(f) = \frac{1}{K_{\text{pair}}} \sum_{k=1}^{K_{\text{pair}}} X_k^{(\text{hist})}(f) \overline{Y_k^{(\text{fut})}(f)}.$$

Thus the shorter side effectively limits  $K_{\text{pair}}$  and hence the stability of  $\hat{S}_{xy}(f)$ , while additional segments on the longer side primarily reduce the variance of the marginal auto-spectra.

We conduct the experiment on ETTh1 with prediction horizon  $N_y = 336$  and history lengths  $N_x \in \{96, 192, 336, 720\}$ , using both iTransformer and DLinear. Table 4 summarizes the results.

On ETTh1, the MSE of both models varies as the history length changes. However, Across all history lengths, the correlation  $R$  between SCP and per-sample errors remains high (typically  $R \geq 0.80$ ). This shows that SCP provides a robust measure of intrinsic predictability that is largely insensitive to the exact history length, even though the models’ absolute accuracy can be sensitive to this choice.

#### D.4 ADDITIONAL DATASET EVALUATION

To complement the main experiments, we further evaluate our framework on the full set of eight long-horizon multivariate forecasting benchmarks: ETTh1, ETTh2, ETTm1, ETTm2, ECL, Weather, Traffic, and ILI. These datasets span electricity, meteorology, transportation, and epidemiology, and cover a wide range of dimensionalities and sampling frequencies. Table 5 summarizes the basic statistics and forecasting horizon settings used in this extended evaluation.

Table 5: Detailed descriptions of the datasets used in our extended evaluation. “Number of variables” gives the dimensionality of each dataset. “Dataset size” denotes the total number of time points in the training, validation, and test splits. “Prediction length” denotes the forecasting horizon; four horizon settings are used for each dataset. “Frequency” is the sampling interval.

Dataset	Dim	Prediction Length	Dataset Size	Frequency	Information
ETTh1, ETTh2	7	{96, 192, 336, 720}	(8545, 2881, 2881)	Hourly	Electricity
ETTm1, ETTm2	7	{96, 192, 336, 720}	(34465, 11521, 11521)	15min	Electricity
ECL	321	{96, 192, 336, 720}	(18317, 2633, 5261)	Hourly	Electricity
Weather	21	{96, 192, 336, 720}	(36792, 5271, 10540)	10min	Weather
Traffic	862	{96, 192, 336, 720}	(12185, 1757, 3509)	Hourly	Transportation
ILI	7	{60, 72}	(617, 74, 170)	Weekly	Epidemiology

Table 6: Long-term multivariate forecasting results on Traffic and ILI datasets. We report MSE, MAE, NMSE, and R. **Bold** marks the best (lowest MSE/MAE) per column across models. *Average* rows give the column-wise mean across models. Predictability reports the per-task linear MSE lower bound ( $MSE_{lb}$ ) and SCP  $\mathcal{P}$  (higher is easier).

Models	Metric	Traffic				ILI	
		96	192	336	720	60	72
iTransformer (Liu et al., 2024)	MSE	<b>0.394</b>	<b>0.385</b>	<b>0.388</b>	<b>0.416</b>	2.001	2.186
	MAE	<b>0.269</b>	0.269	0.274	0.290	0.954	1.033
	NMSE	0.303	0.302	0.265	0.273	0.806	1.032
	R	0.849	0.917	0.959	0.968	0.687	0.847
TimeMixer (Wang et al., 2024)	MSE	0.485	0.423	0.407	0.437	2.272	<b>1.928</b>
	MAE	0.319	0.285	0.275	0.297	0.977	<b>0.938</b>
	NMSE	0.376	0.327	0.275	0.292	1.055	0.714
	R	0.919	0.939	0.971	0.960	0.783	0.781
DLinear (Zeng et al., 2023)	MSE	0.649	0.459	0.436	0.450	2.671	2.661
	MAE	0.396	0.305	0.296	0.306	1.083	1.114
	NMSE	0.541	0.370	0.305	0.302	1.077	1.099
	R	0.899	0.929	0.970	0.968	0.850	0.919
PatchTST (Nie et al., 2023)	MSE	0.451	0.402	0.401	0.434	<b>1.758</b>	2.010
	MAE	0.288	<b>0.263</b>	<b>0.267</b>	<b>0.289</b>	<b>0.863</b>	0.948
	NMSE	0.356	0.321	0.277	0.288	0.753	0.792
	R	0.893	0.920	0.966	0.965	0.675	0.852
TimesNet (Wu et al., 2023)	MSE	0.606	0.608	0.630	0.672	2.160	1.994
	MAE	0.327	0.329	0.347	0.357	0.961	0.974
	NMSE	0.370	0.376	0.331	0.335	0.866	0.697
	R	0.960	0.969	0.946	0.983	0.820	0.742
Average	MSE	0.517	0.455	0.452	0.482	2.172	2.156
	MAE	0.320	0.290	0.292	0.308	0.968	1.001
	NMSE	0.389	0.339	0.291	0.298	0.911	0.867
	R	0.904	0.935	0.962	0.969	0.763	0.828
Predictability	$MSE_{lb}$	<b>0.803</b>	<b>0.616</b>	<b>0.400</b>	<b>0.636</b>	<b>2.151</b>	<b>2.681</b>
	$\mathcal{P}$	<b>0.514</b>	<b>0.619</b>	<b>0.760</b>	<b>0.610</b>	<b>0.560</b>	<b>0.466</b>

Table 6 reports detailed long-horizon multivariate forecasting results on the Traffic and Illness datasets for five representative architectures under a matched-information protocol: the history length equals the prediction horizon ( $N \in \{96, 192, 336, 720\}$  for Traffic and  $N \in \{60, 72\}$  for ILI), with identical preprocessing and no drop-last. We report MSE, MAE, normalized MSE (NMSE), and correlation coefficient  $R$ , together with the linear MSE lower bound  $MSE_{lb}$  and SCP-based predictability  $\mathcal{P}$  for each task.

On Traffic, iTransformer consistently achieves the lowest MSE across all horizons. On Illness, TimeMixer and PatchTST achieve better accuracy than the other baselines. Across both datasets, the SCP and linear MSE lower bound remain well aligned with the empirical results, indicating that our predictability-aware metrics continue to agree with, and help interpret, standard error-based evaluations.

## E LIMITATIONS AND DISCUSSION

Our estimates rely on windowed Welch spectra and an assumption of local second-order stationarity within each window. These conditions are used to derive a tight connection between SCP and the Bayes-optimal MSE lower bound in the linear setting. In practice, however, our real-world benchmarks are visibly nonstationary, with predictability drifting over time (e.g., Fig. 4), yet SCP

1188 remains well aligned with observed forecasting errors, including those of nonlinear architectures.  
1189 This supports the use of SCP/LUR as robust diagnostics of intrinsic difficulty even when the formal  
1190 tightness guarantees no longer strictly hold.

1191 Moreover, the framework is not restricted to univariate series: our matrix-based formulation natu-  
1192 rally extends SCP/LUR to multivariate settings, and nonlinear preprocessing can be applied to map  
1193 raw observations into a feature space before computing SCP, so that the linear prediction lower  
1194 bound is evaluated in that transformed space. A full treatment of conditional and partial coherence  
1195 structures remains an interesting direction for future work.

1196 Finally, we view predictability-aware learning schemes as a promising next step. Extending SCP  
1197 to guide active data selection, curriculum-style scheduling across predictability regimes, and hard-  
1198 sample mining based on  $\mathcal{P}$  and LUR are natural applications of the framework beyond the purely  
1199 diagnostic use studied in this paper.  
1200

1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241