
ProbTS: Benchmarking Point and Distributional Forecasting across Diverse Prediction Horizons

Jiawen Zhang*
HKUST (GZ)
Guangzhou, China
jiawe.zh@gmail.com

Xumeng Wen
Microsoft Research Asia
Beijing, China
xumengwen@microsoft.com

Zhenwei Zhang*
Tsinghua University
Beijing, China
zzw20@mails.tsinghua.edu.cn

Shun Zheng
Microsoft Research Asia
Beijing, China
shun.zheng@microsoft.com

Jia Li
HKUST (GZ)
Guangzhou, China
jialeel@ust.hk

Jiang Bian
Microsoft Research Asia
Beijing, China
jiang.bian@microsoft.com

Abstract

Delivering precise point and distributional forecasts across a spectrum of prediction horizons represents a significant and enduring challenge in the application of time-series forecasting within various industries. Prior research on developing deep learning models for time-series forecasting has often concentrated on isolated aspects, such as long-term point forecasting or short-term probabilistic estimations. This narrow focus may result in skewed methodological choices and hinder the adaptability of these models to uncharted scenarios. While there is a rising trend in developing universal forecasting models, a thorough understanding of their advantages and drawbacks, especially regarding essential forecasting needs like point and distributional forecasts across short and long horizons, is still lacking. In this paper, we present ProbTS, a benchmark tool designed as a unified platform to evaluate these fundamental forecasting needs and to conduct a rigorous comparative analysis of numerous cutting-edge studies from recent years. We dissect the distinctive data characteristics arising from disparate forecasting requirements and elucidate how these characteristics can skew methodological preferences in typical research trajectories, which often fail to fully accommodate essential forecasting needs. Building on this, we examine the latest models for universal time-series forecasting and discover that our analyses of methodological strengths and weaknesses are also applicable to these universal models. Finally, we outline the limitations inherent in current research and underscore several avenues for future exploration.¹

1 Introduction

Time-series forecasting has extensive applications in various industries, including traffic flow forecasting [43], renewable energy forecasting [67], and diverse forecasting demands in retail [8], finance [29], physical system [39], and climate [48]. It is crucial to provide forecasts across different prediction horizons, addressing both short- and long-term planning needs [13, 26, 4, 61]. Moreover, modern decision-making processes typically require not only point forecasts to quantify planning efficiency but also robust distributional estimations to manage uncertainty effectively [24, 30]. The fundamental need to produce accurate point and distributional forecasts across various horizons presents significant challenges to existing forecasting approaches.

*This work was done during the internship at Microsoft Research Asia.

¹Project repository: <https://github.com/microsoft/ProbTS>

Nevertheless, much of the previous research on developing deep learning models for time-series forecasting has often focused on isolated aspects, such as long-term point forecasting or short-term distribution estimations. This narrow focus may result in skewed methodological choices and hinder the adaptability of these models to rarely evaluated scenarios. For example, studies such as [78, 72, 38, 76, 11, 71, 49, 73, 40] have primarily explored neural architecture designs tailored for long-term point forecasting with strong trending and seasonal patterns. However, it remains unclear how these advancements can be effectively extended to capture complicated distributions and whether these designs maintain their effectiveness in short-term scenarios. Conversely, research such as [58, 57, 62, 7, 33] adapts deep generative models [17, 28] for probabilistic forecasting, specializing in characterizing complex data distributions. Yet, these models have mainly been developed and evaluated in short-term scenarios, leaving questions about their effectiveness in long-term forecasting and their ability to preserve point forecasting performance.

Despite the recent surge in building time-series foundation models over the past year [18, 56, 14, 10, 15, 25, 41, 20, 70, 2, 23, 74], our understanding of their advantages and limitations, especially regarding essential forecasting needs like point and distributional forecasts across various horizons, is still limited. Many of these models claim to support arbitrary prediction horizons, employing different mechanisms that come with their own set of advantages and drawbacks. Among them, a select few offer capabilities for distributional forecasting, which, however, are typically confined to predefined closed-form distributions [56, 70] or discrete distributions with value quantization [2]. The emergence of these foundation models has brought about unprecedented zero-shot forecasting capabilities. Consequently, it is both timely and crucial to delve into an evaluation of their strengths and weaknesses, especially in relation to the fundamental forecasting needs mentioned earlier.

In this study, we present ProbTS, a benchmark tool crafted to serve as a comprehensive platform for assessing those key forecasting needs and for performing a detailed comparison of several state-of-the-art models developed in recent years. To address the core forecasting requirements, ProbTS includes a broad array of datasets and spans various forecasting horizons. It also utilizes both point and distributional metrics to facilitate a thorough performance evaluation.

Our research reveals that the specific data characteristics inherent to different forecasting requirements often play a crucial role in guiding the selection of model designs. As a result, it is crucial to have a comprehensive view of the essential forecasting needs. To aid in the analysis and interpretation of performance, we measure three essential data characteristics in ProbTS: the strength of trends and seasonality, and the complexity of the data distribution. Moreover, we have explicitly distinguished three fundamental methodological aspects within ProbTS that differentiate the existing forecasting models, largely influencing their pros and cons. The first aspect involves the approach to distributional forecasting, ranging from models focused on point forecasts [49, 40] to those using pre-defined distribution heads based on specific data assumptions [56, 70]. The second aspect is the decoding scheme used to generate multi-step forecasts, which can be either autoregressive (AR) or non-autoregressive (NAR). The third aspect pertains to the normalization choice, where the long-term point forecasting models typically employ reversible instance normalization (RevIN) [32] while short-term probabilistic ones often use mean scaling strategies [58, 57].

By utilizing ProbTS, we conduct a systematic comparison between studies that focus on long-term point forecasting and those aimed at short-term distributional estimation, employing various forecasting horizons and evaluation metrics. Our overarching finding is that the strengths of these methods tend to diminish in scenarios they are rarely evaluated in, highlighting several important but unresolved research questions. Notably, while recent probabilistic forecasting approaches have shown proficiency in short-term distribution estimation, we find that long-term distributional forecasting remains a significant challenge. This challenge stems from achieving distribution estimation that remains both efficient and effective as the prediction horizon extends—a topic that has not been thoroughly investigated in existing literature. Additionally, our analysis uncovers a clear divide in the choice of decoding schemes: most long-term point forecasting methods opt for NAR, whereas choices in short-term forecasting studies are more evenly split. Further investigation suggests that the preference for NAR methods stems from existing AR models’ difficulty in managing error accumulation, particularly over extended horizons with strong trends. However, we observe that a proper normalization strategy can significantly improve AR models in long-term forecasting, opening new possibilities for AR-based approaches. Moreover, AR decoding performs better in scenarios with pronounced seasonality, indicating potential for refining these strategies, particularly for long-

term forecasts. Given the inefficiency of existing NAR-based probabilistic methods like CSDI, our comparison highlights the need for further exploration of decoding strategies in future research.

Furthermore, we have expanded the analytical framework of ProbTS to include an examination of several very recently developed time-series foundation models, which has allowed us to re-validate some of our earlier findings. Interestingly, there appears to be a relatively even split in their preference for AR and NAR decoding schemes. Our analysis reaffirms the limitation of AR in handling time-series data, as we observe that AR-based foundation models tend to excel at shorter horizons. However, their performance advantages often significantly diminish over longer forecasting periods. This underscores the critical need for future research to focus on addressing the issue of error accumulation in AR-based foundation models. Besides, our exploration reveals that current probabilistic foundation models may face challenges when dealing with complex data distributions. This observation suggests that the integration of more sophisticated distribution estimation techniques could enhance the development of time-series foundation models.

In summary, we have made the following contributions.

- Introduction of ProbTS, a benchmark tool designed for a thorough evaluation of essential forecasting needs, towards precise point and probabilistic forecasting across varied horizons.
- Comprehensive analysis of methodological variations within forecasting models, especially regarding distributional estimation methods and decoding schemes (AR vs. NAR), which illuminates significant yet previously underexplored research challenges.
- Extension of our analytical framework to include the latest time-series foundation models, providing insights into the implications of their methodological choices and underscoring important directions for future research endeavors.

2 Related Work

Classical Time-series Forecasting Models In recent years, classical research in time-series forecasting has bifurcated into two distinct but complementary streams. The first stream has concentrated on refining neural architecture designs for long-term forecasting, primarily employing non-autoregressive decoding schemes to address scenarios with pronounced trend and seasonality. This stream has evolved from enhancing multi-layer perceptrons [53, 76] to developing specialized recurrent or convolutional neural networks [34, 37], and introducing Transformer-based models [66, 49, 40]. Despite achieving advancements in point forecasts, these efforts mainly capture average future changes, with only a few adopting approaches like quantile regression to partially overcome this limitation [69, 36]. On the other hand, the second stream, probabilistic time-series forecasting specializes in capturing the intricate data distribution of future time series. It encompasses a spectrum of techniques, from utilizing predefined likelihood functions [55, 60] and Gaussian copulas [59, 19] to exploring advanced deep generative models [58, 7]. Unlike the first stream, this branch employs both AR [58, 57] and NAR decoding schemes [62, 7, 33], often utilizing standard neural network architectures to represent time series [16, 58, 57, 7, 19], though some studies propose customized designs [62, 35, 5]. Together, these streams highlight the diverse approaches to forecasting, ranging from point predictions focusing on the mean future variations to probabilistic forecasts that capture the full distribution of future values. In Appendix A.1, we summarize a comparison of these models on the coverage of essential forecasting needs and their methodological preferences.

Universal Time-series Foundation Models Over the past year, the development of time-series foundation models has greatly accelerated, driven by the success of language foundation models [9]. This wave has seen models such as Lag-Llama [56], TimesFM [15], Timer [41], and Chronos [2] adopting the decoder-only Transformer architecture with an AR decoding scheme. Conversely, models like ForecastPFN [18], MOIRAI [70], TTM [20], and UniTS [23] employ the NAR decoding, often using variable-length placeholders to indicate prediction positions for different horizons. Probabilistic forecasting is less common, with MOIRAI and Lag-Llama integrating pre-defined distribution heads (Student-t for Lag-Llama and a mixture for MOIRAI) while Chronos uses quantized bins to accommodate time-series values and adopts Softmax outputs for distribution approximation. The strategic choice between AR and NAR decoding and the method for distributional estimation highlight distinct trade-offs. For an extensive comparison, see Appendix A.2.

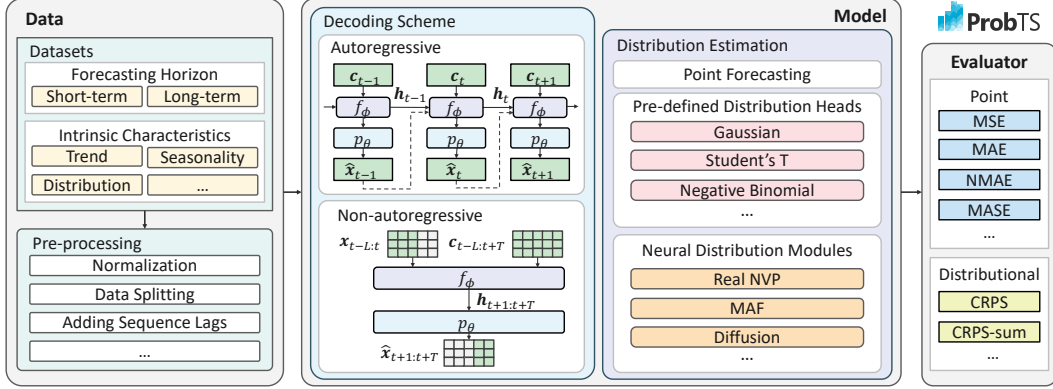


Figure 1: An overview of ProbTS.

Toolkits for Time-series Forecasting. We observe a plethora of toolkits that have been developed for time-series forecasting. These range from those primarily designed for point forecasting, such as Prophet [63], sktime [42], tsai [52], and TSLib [71], to others that incorporate probabilistic forecasting, including GluonTS [1], PyTorchTS [58], PyTorchForecasting², and NeuralForecast³. In creating ProbTS, we built upon the foundations laid by tools like PyTorchTS, GluonTS, and TSLib. Our unique contribution is a detailed approach that supports both precise point and probabilistic forecasting over various horizons, and examines methodological differences in forecasting models, especially regarding distributional estimation and decoding schemes (AR vs. NAR). Additionally, ProbTS integrates cutting-edge time-series foundation models, making it a comprehensive benchmark tool for tackling current and future challenges in time-series forecasting. A comparison of ProbTS with existing toolkits, focusing on functionalities and features, is provided in Appendix A.3.

3 The ProbTS Tool

This section offers a concise overview of the ProbTS tool’s design and implementation. The core modules and the primary pipeline of ProbTS are depicted in Figure 1.

Data We aggregate publicly accessible datasets used for both short-term and long-term forecasting. Initial data visualization analyses reveal that the data domains and forecasting horizons significantly influence specific data characteristics within a given forecasting horizon. For instance, many long-term forecasting scenarios exhibit clear trend and seasonality patterns within a forecasting window, while numerous short-term forecasting cases display irregular variations within a short sliding window. Consequently, we have developed quantified indicators, such as trend and seasonality strengths, along with *non-Gaussianity* to indicate the complexity of data distribution within a forecasting window. Detailed information about dataset statistics, visualization analyses, and quantified measures can be found in Appendix B.1.1, B.1.2, B.1.3, and B.1.4. The quantified measurements for all forecasting scenarios are compiled in Table 1.

Metrics ProbTS incorporates a broad range of evaluation metrics to enable a thorough assessment of both point and distributional forecasts. These metrics are elaborated in detail in Appendix B.2. In this paper, we primarily use the normalized mean absolute error (NMAE) for point forecasts and the continuous ranked probability score (CRPS) for distributional forecasts to succinctly communicate the critical insights discovered. It is noteworthy that some methods reproduced in ProbTS, their original papers reported certain point forecast metrics before de-normalizing forecasts to the initial scale [75, 71, 49] or primarily reveal aggregated distributional metrics over all time-series variates, namely CRPS-sum [59, 57, 58]. We have verified our reproduced results align with their reported results and utilized the unified metrics in this study to offer a comprehensive and fair comparison of these studies from different research threads.

²github.com/jdb78/pytorch-forecasting

³github.com/Nixtla/neuralforecast

Table 1: This table includes a quantitative assessment of the inherent characteristics for all forecasting scenarios, each corresponding to a dataset with a specific forecasting horizon. We use the suffixes "-S" and "-L" to differentiate between short-term and long-term scenarios. Quantified indicators encompass trend and seasonality strengths, as well as non-Gaussianity, where a higher value signifies a greater deviation from a Gaussian distribution.

Dataset-Horizon	Exchange-S	Solar-S	Electricity-S	Traffic-S	Wikipedia-S	ETTm1-L	ETTm2-L
Trend F_T	0.9982	0.1688	0.6443	0.2880	0.5253	0.9462	0.9770
Seasonality F_S	0.1256	0.8592	0.8323	0.6656	0.2234	0.0105	0.0612
Non-Gaussianity	0.2967	0.5004	0.3579	0.2991	0.2751	0.0833	0.1701

Dataset-Horizon	ETTh1-L	ETTh2-L	Electricity-L	Traffic-L	Weather-L	Exchange-L	ILI-L
Trend F_T	0.7728	0.9412	0.6476	0.1632	0.9612	0.9978	0.5438
Seasonality F_S	0.4772	0.3608	0.8344	0.6798	0.2657	0.1349	0.6075
Non-Gaussianity	0.0719	0.1422	0.1533	0.1378	0.1727	0.1082	0.1112

Model The model module in ProbTS explicitly differentiates critical methodological decisions, especially the decoding scheme (AR vs NAR) and the distributional estimation approach. Specifically, we employ the following mathematical formulation. We denote an element of a multivariate time series as $x_t^k \in \mathbb{R}$, where k represents the variate index and t denotes the time index. At time step t , we have a multivariate vector $\mathbf{x}_t \in \mathbb{R}^K$. Each x_t^k is associated with covariates $\mathbf{c}_t^k \in \mathbb{R}^N$, which encapsulates auxiliary information about the observations. Given a length- T forecast horizon, a length- L observation history $\mathbf{x}_{t-L:t}$ and corresponding covariates $\mathbf{c}_{t-L:t}$, the objective in time series forecasting is to generate the vector of future values $\mathbf{x}_{t+1:t+T}$. Based on established conventions, we categorize forecast as short-term if the horizon $T \leq \mathcal{T}$ [57, 62], and long-term if $T \gg \mathcal{T}$ [75, 49, 40], where \mathcal{T} represents the primary periodicity of the data (e.g., 24 for hourly frequency). To represent point and distributional forecasting in a unified way, here we divide a model into an encoder f_ϕ and a forecaster p_θ . An encoder is tasked with generating expressive hidden states $\mathbf{h} \in \mathbb{R}^D$. Under *autoregressive* decoding scheme, encoder forecasts variates using their past values: $\mathbf{h}_t = f_\phi(\mathbf{x}_{t-1}, \mathbf{c}_t, \mathbf{h}_{t-1})$. Under the *non-autoregressive* scheme, the encoder generates all the forecasts in one step: $\mathbf{h}_{t+1:t+T} = f_\phi(\mathbf{x}_{t-L:t}, \mathbf{c}_{t-L:t+T})$. A forecaster p_θ is employed either to directly estimate *point forecasts* as $\hat{\mathbf{x}}_t = p_\theta(\mathbf{h}_t)$, or to perform sampling based on the estimated *probabilistic distributions* as $\hat{\mathbf{x}}_t \sim p_\theta(\mathbf{x}_t | \mathbf{h}_t)$. In addition, the normalization choices utilized by different research branches vary, with a detailed analysis provided in Appendix D.1.

4 Results and Analyses

Utilizing ProbTS, we conducted a comprehensive benchmarking and analysis of a diverse range of state-of-the-art models from different strands of research. We mainly assessed these models using NAME and CRPS metrics across multiple forecasting horizons, repeating each experiment five times with different seeds to ensure result reliability.

Selected Models for Comparison. Our selection criteria for models focused on a balance of performance, reproducibility, and simplicity. For long-term point forecasting, we included models like iTransformer [40], PatchTST [49], TimesNet [71], N-HiTS [11], and LTSF-Linear [75]. Probabilistic forecasting methods selected include GRU NVP, GRU MAF, Trans MAF [58], TimeGrad [57], and CSDI [62]. Additionally, general architectures like Linear, GRU [12], and Transformer [66], along with simple non-parametric baselines, were evaluated as a reference. For foundation models, reproducible methods such as Lag-Llama [56], TimesFM [15], Timer [41], MOIRAI [70], Chronos [2], and UniTS [23] were included. Detailed implementation specifics are in Appendix B.3.

Due to space constraints, comprehensive comparison results are placed in Appendix C, with detailed results for short-term and long-term forecasting in Tables 9 and 10, respectively. Zero-shot evaluations of pre-trained time-series foundation models are detailed in Tables 11 and 12. Our evaluation highlights the critical relationship between forecasting requirements, data properties, and modeling strategies. It aims to shed light on the strengths and limitations of current approaches, paving the way for uncovering novel research avenues.

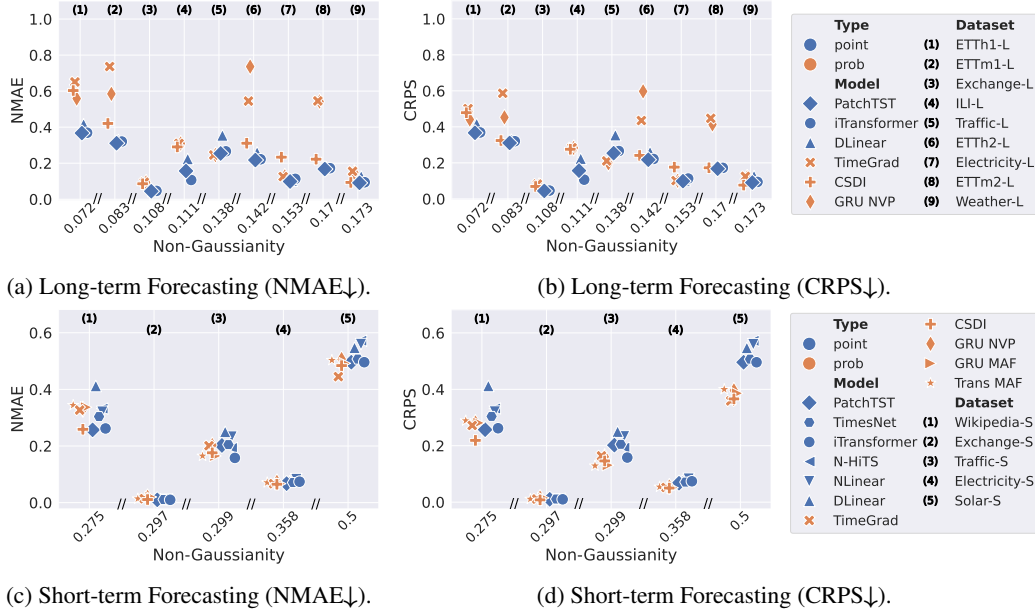


Figure 2: We present a comprehensive comparison between classical models designed for long-term point forecasting and short-term distributional forecasting across various prediction horizons. It utilizes a non-Gaussianity score to highlight the complexity of the data distribution across different datasets. The aggregated performance metrics are derived from Tables 9 and 10.

4.1 Analyzing Classical Models for Time-series Forecasting

We examine traditional non-universal time-series models from distinct research branches: one branch focuses on developing customized neural architectures tailored for long-term point forecasting, while the other branch concentrates on creating advanced probabilistic methods for short-term distributional forecasting. Our investigation confirms the effectiveness of these models for their intended purposes. However, we observe a notable trend: the strengths of these methods tend to diminish in scenarios where they are seldom tested.

Diminishing Advantages of Customized Architectures in Short-term Forecasting Scenarios

The comparative analysis presented in Figures 2a and 2c showcases the performance of point and probabilistic forecasting methods with respect to the NMAE metric. These figures also illustrate how NMAE values correlate with non-Gaussianity, a measure we employ to evaluate the complexity of data distributions. It becomes evident that customized architectures, originally crafted for long-term forecasting, tend to lose their competitive performance in short-term scenarios. This phenomenon could be attributed to the increased importance of accurately characterizing complex data distributions within shorter forecasting windows, where higher non-Gaussianity scores are indicative of this necessity. A closer look at Figures 2c and 2d further reveals that the performance disparity measured with CRPS becomes even more pronounced for datasets characterized by significant non-Gaussianity, such as Solar-S. This observation underscores the critical need for incorporating short-term patterns and distributional estimation capabilities into the design of new forecasting architectures.

Significant Performance Degradation for Existing Probabilistic Methods in Long-term Distributional Forecasting

The performance of current probabilistic forecasting models in long-term scenarios, even when assessed using distributional metrics such as CRPS, reveals notable limitations. This is highlighted by the comparison between Figures 2b and 2d, which shows a significant drop in performance for models like TimeGrad on ETTm1-L, GRU NVP on ETTh2-L, and Weather-L datasets. The decline in performance can be attributed to the fact that these probabilistic models were not specifically designed with the unique challenges of long-term forecasting in mind. This oversight has mixed influences. On the positive side, the design of these methods has led to a more balanced approach in the choice between AR and NAR decoding schemes, providing a versatile

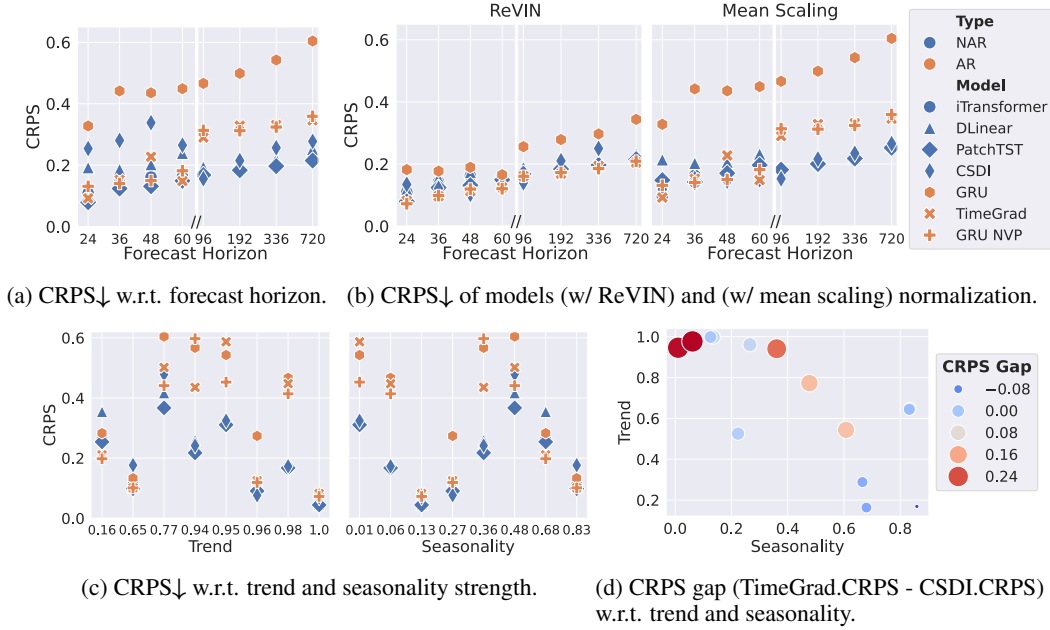


Figure 3: We explore the challenges faced by current models in conducting long-term distributional forecasting, with insights drawn from Table 10 and Table 16. Subplot (a) shows significant error increases in AR-based models, averaged across all datasets except Traffic. Subplot (b) demonstrates how the instance-level normalization impacts performance in long-term forecasting. Subplots (c) examine how trends and seasonality impact performance across all long-term forecasting datasets and horizons.

Subplot (d) further investigates the combined effects of trend and seasonality, using lighter and smaller circles to indicate situations where AR-based models are favored over NAR-based ones.

foundation for probabilistic forecasting. However, the downside is more significant: no matter using which decoding schemes, existing probabilistic models face considerable challenges when applied to long-term distributional forecasting. We will dive deeper into the specific challenges associated with each decoding scheme next. Here the performance gap underscores the need for future research to systematically investigate long-term distributional forecasting.

Different Decoding Schemes & Challenges in Long-term Distributional Forecasting Existing probabilistic forecasting methods exhibit a balanced preference for both AR and NAR decoding schemes. For instance, TimeGrad employs an AR decoding scheme, whereas CSDI utilizes an NAR approach. This contrasts starkly with the aforementioned customized architectures, which solely opt for NAR decoding. These two types of decoding schemes, however, confront distinctive challenges when applied to long-term probabilistic forecasting. With original normalization strategy, AR probabilistic models like TimeGrad struggle with error accumulation, particularly as the forecast horizon extends or trends strengthen, the performance gap widens, as shown in Figures 3a and 3c. On the other hand, NAR models such as CSDI encounter memory constraints in long-term forecasts (detailed in Appendix D.7). Moreover, Table 10 reveals that even on smaller datasets, such as ETTm2 and ETTh1, CSDI’s performance in long-term scenarios is less than optimal, indicating reduced learning efficiency by the extension of the forecasting horizon.

The Unexpected Superiority of AR Decoding in Addressing Strong Seasonality Despite its drawbacks, AR-based models, as used in TimeGrad, excels in capturing strong seasonality, outperforming models like PatchTST in scenarios such as the Traffic dataset (Table 10). This advantage is further analyzed in Figures 3c and 3d, showing AR’s increasing benefit with stronger seasonal patterns. This suggests AR’s potential in long-term forecasting could be revitalized with solutions to its error accumulation challenge in long horizons.

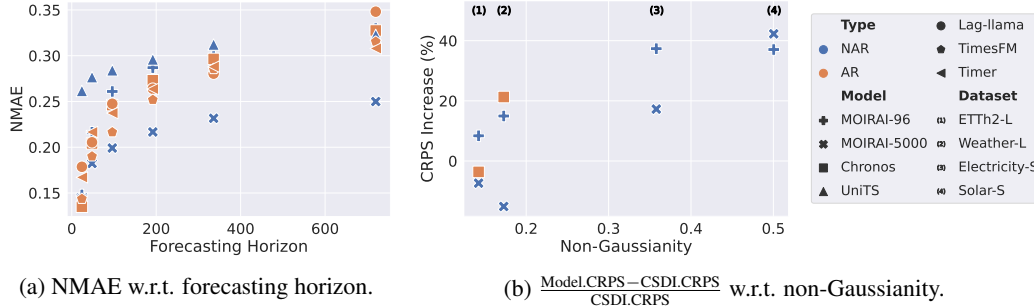


Figure 4: We evaluate the efficacy of time-series foundation models for various forecasting horizons and distributional estimation. Subplot (a), derived from Table 11 and excluding results from the Electricity dataset, demonstrates the short-term forecasting capabilities and long-term error accumulation of AR-based models. Subplot (b), draw from Table 12, investigates short-term distributional estimation, highlighting the performance challenges of foundation models compared to CSDI in handling complex data distributions. Note that we include MOIRAI with two different context lengths, 96 and 5000, as context length significantly affects its transfer performance.

ReVIN’s Effectiveness in Long-term Forecasting with Exceptions. ReVIN significantly enhances AR-based models in long-term forecasting, as shown in Figures 3b and 7. Notably, on the ETTh1 dataset, GRU NVP (w/ ReVIN) even outperforms PatchTST (w/ ReVIN). While ReVIN offers substantial improvements for most models across most datasets, it brings negative impact on the Traffic dataset. The Traffic dataset features strong seasonality but minimal trends, thus we speculate that the major distribution shift addressed by ReVIN is related to normalizing the effect of trending. These findings indicate that normalizing the trending effect could be a direction to alleviate error accumulation of AR-based models in long-term forecasting. However, we also observe that ReVIN does not seem to be an ideal match for the NAR probabilistic model. For instance, CSDI (w/ ReVIN) performs worse than CSDI (w/ Scaling) on the Weather dataset. Further research in developing effective normalization strategies for NAR probabilistic models is necessary.

No Dominating Normalization Strategies in Short-term Forecasting. While ReVIN is effective for long-term scenarios, it does not adequately address the challenges faced by short-term probabilistic models. As shown in Table 14, ReVIN fails to consistently deliver significant improvements for models such as CSDI, TimeGrad, and GRU NVP in short-term forecasting. The mean scaling strategy has proven to be the most reliable option for these models, explaining its widespread use. Although omitting instance-level normalization is occasionally acceptable, it can lead to significant issues, as seen with TimeGrad (without normalization) on the Wikipedia and Solar datasets, and GRU NVP (without normalization) on Electricity. We provide detailed analysis in Appendix D.1.

4.2 Analyzing Foundation Models for Universal Time-series Forecasting

We next explore the capabilities of recent foundation models in universal time-series forecasting, focusing on their performance across different prediction horizons and their ability to estimate distributions, especially regarding their zero-shot transfer capabilities on unseen datasets. Table 11 showcases their significant progress, sometimes outperforming traditional models without re-training. Using the analytic framework of ProBTs, we delve into their methodological pros and cons.

Navigating the AR Decoding Challenge over Extended Forecasting Horizons Figure 4a illustrates how the performance of various time-series foundation models evolves in relation to expanding forecasting horizons. It is evident that for shorter horizons, AR-based foundation models such as TimesFM and Timer exhibit highly competitive performance, on par with NAR-based models like MOIRAI. However, the advantage of NAR-based decoding becomes increasingly apparent as the forecasting horizon lengthens, as demonstrated by the widening performance gap between TimesFM and MOIRAI. This trend is consistent with our earlier observation that without proper normalization strategies, AR-based methods could suffer from significant error accumulation when applied to long-term time-series forecasting. Given the inherent strengths of AR decoding, such as its superiority at capturing strong seasonality and its robust performance in certain short-term forecasting

scenarios, it is clear that further research is warranted to explore ways to overcome its limitations in long-term forecasting contexts. This could potentially unlock new avenues for enhancing the versatility and effectiveness of AR-based time-series foundation models across a broader range of forecasting horizons.

The Critical Role of Addressing Complex Data Distributions Figure 4a illustrates the incremental changes in CRPS among leading probabilistic time-series foundation models, such as MOIRAI and Chronos, compared to the best-performing short-term probabilistic model, CSDI. It becomes apparent that in scenarios characterized by complex data distributions, indicated by higher non-Gaussianity scores, the performance decline of MOIRAI in relation to CSDI becomes notably more pronounced. This phenomenon may be attributed to MOIRAI’s approach to supporting distributional forecasting, which involves utilizing a mixture of predefined distribution heads. While this method is efficient and effective for certain applications, it may lack the expressiveness required to accurately model more complex data distributions. Furthermore, these observations underscore that, in specific contexts, foundation models might not be able to fully replace traditional models that have been specifically tailored and trained for particular domains. Additionally, the prospect of fine-tuning these foundation models as a remedy is less economically viable, primarily due to their significantly larger size. This highlights the importance of not only continuing to refine foundation models to enhance their adaptability and performance across a wide spectrum of data distributions but also recognizing the continued relevance of domain-specific models, especially for handling intricate data distributions where a more nuanced approach may be necessary.

5 Conclusion

In this study, we introduced ProbTS, a benchmark tool tailored for evaluating essential forecasting needs, which facilitates a detailed comparison of various state-of-the-art models in the context of time-series forecasting. Through our comprehensive analysis, we identified significant challenges and opportunities in the realm of time-series forecasting, particularly highlighting the need for models that can effectively address both point and probabilistic forecasting across diverse horizons.

Limitations While our study represents a significant step forward in understanding and evaluating time-series forecasting models, it does come with many limitations. A predominant focus of our work is on empirical analysis, relying heavily on intuitions and experimental observations, which may lack the depth that theoretical foundations could provide. Additionally, our exploration, though extensive, might not encompass all the nuanced factors that influence model performance. By concentrating on major methodological decisions such as AR versus NAR decoding schemes, we may inadvertently overlook other critical aspects that could play a decisive role in forecasting accuracy. Moreover, the datasets employed for evaluation, despite their diversity and relevance to current research threads, may not fully capture the vast spectrum of real-world forecasting challenges. This limitation is particularly pronounced when comparing different foundation models, as their pre-training often involves an even broader array of data, potentially skewing the comparative analysis.

Future Directions The insights derived from our study open the door to numerous promising research directions. Addressing the shortcomings of AR and NAR decoding schemes, especially in their application across varying forecasting horizons, emerges as a critical area for future exploration. Innovating effective architecture designs that can navigate the intricacies of short-term forecasting challenges and devising efficient methods for long-term probabilistic forecasting stand out as urgent needs. For AR-based models, reducing error accumulation remains essential, with ReVIN-style normalization showing potential for improving long-term forecasting. Additionally, exploring effective normalization strategies for NAR-based probabilistic models is an underdeveloped yet promising area. Equally important is the enhancement of models’ abilities to characterize complex data distributions, which could significantly improve the adaptability and effectiveness of foundation models. Beyond these technical endeavors, expanding the scope of datasets used for evaluation to encompass a wider range of real-world scenarios will be crucial for validating the robustness and versatility of future forecasting models. Lastly, integrating theoretical insights with empirical findings could provide a more holistic understanding of model behaviors, contributing to the development of more sophisticated and nuanced forecasting solutions.

References

- [1] Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Danielle C. Maddix, Syama Rangapuram, David Salinas, Jasper Schulz, Lorenzo Stella, Ali Caner Türkmen, and Yuyang Wang. 2020. GluonTS: Probabilistic and Neural Time Series Modeling in Python. *JMLR* 21, 116 (2020), 1–6.
- [2] Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Syndar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, Jasper Zschiegner, Danielle C. Maddix, Hao Wang, Michael W. Mahoney, Kari Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Yuyang Wang. 2024. Chronos: Learning the Language of Time Series. *arXiv preprint arXiv:2403.07815* (2024).
- [3] George Athanasopoulos, Rob J Hyndman, Haiyan Song, and Doris C Wu. 2011. The tourism forecasting competition. *International Journal of Forecasting* 27, 3 (2011), 822–844.
- [4] Joaquim Barros, Miguel Araujo, and Rosaldo JF Rossetti. 2015. Short-term real-time traffic prediction methods: A survey. In *International Conference on Models and Technologies for Intelligent Transportation Systems*.
- [5] Shane Bergsma, Timothy Zeyl, Javad Rahimipour Anaraki, and Lei Guo. 2022. C2FAR: Coarse-to-Fine Autoregressive Networks for Precise Probabilistic Forecasting. In *NeurIPS*.
- [6] Aadyot Bhatnagar, Paul Kassianik, Chenghao Liu, Tian Lan, Wenzhuo Yang, Rowan Cassius, Doyen Sahoo, Devansh Arpit, Sri Subramanian, Gerald Woo, et al. 2021. Merlion: A machine learning library for time series. *arXiv preprint arXiv:2109.09265* (2021).
- [7] Marin Bilos, Kashif Rasul, Anderson Schneider, Yuriy Nevmyvaka, and Stephan Günnemann. 2023. Modeling Temporal Data as Continuous Functions with Stochastic Process Diffusion. In *ICML*. 2452–2470.
- [8] Joos-Hendrik Böse, Valentin Flunkert, Jan Gasthaus, Tim Januschowski, Dustin Lange, David Salinas, Sebastian Schelter, Matthias Seeger, and Yuyang Wang. 2017. Probabilistic demand forecasting at scale. *VLDB* (2017).
- [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *NeurIPS*.
- [10] Defu Cao, Furong Jia, Sercan O Arik, Tomas Pfister, Yixiang Zheng, Wen Ye, and Yan Liu. 2024. TEMPO: Prompt-based Generative Pre-trained Transformer for Time Series Forecasting. In *ICLR*.
- [11] Cristian Challu, Kin G. Olivares, Boris Oreshkin, Federico Ramirez, Max Canseco, and Artur Dubrawski. 2023. NHITS: Neural Hierarchical Interpolation for Time Series Forecasting. In *AAAI*. 6989–6997.
- [12] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014).
- [13] Estee Y Cramer, Evan L Ray, Velma K Lopez, Johannes Bracher, Andrea Brennen, Alvaro J Castro Rivadeneira, Aaron Gerding, Tilmann Gneiting, Katie H House, Yuxin Huang, et al. 2022. Evaluation of individual and ensemble probabilistic forecasts of COVID-19 mortality in the United States. *Proceedings of the National Academy of Sciences* (2022).
- [14] Luke Nicholas Darlow, Qiwen Deng, Ahmed Hassan, Martin Asenov, Rajkarn Singh, Artjom Joosen, Adam Barker, and Amos Storkey. 2024. DAM: Towards a Foundation Model for Forecasting. In *ICLR*.

- [15] Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. 2024. A decoder-only foundation model for time-series forecasting. In *ICML*.
- [16] Emmanuel de Bézenac, Syama Sundar Rangapuram, Konstantinos Benidis, Michael Bohlke-Schneider, Richard Kurlle, Lorenzo Stella, Hilaf Hasson, Patrick Gallinari, and Tim Januschowski. 2020. Normalizing Kalman Filters for Multivariate Time Series Analysis. In *NeurIPS*. 2995–3007.
- [17] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. 2017. Density Estimation using Real NVP. In *ICLR*.
- [18] Samuel Dooley, Gurnoor Singh Khurana, Chirag Mohapatra, Siddartha V Naidu, and Colin White. 2023. ForecastPFN: Synthetically-trained zero-shot forecasting. In *NeurIPS*.
- [19] Alexandre Drouin, Étienne Marcotte, and Nicolas Chapados. 2022. TACTiS: Transformer-Attentional Copulas for Time Series. In *ICML*. 5447–5493.
- [20] Vijay Ekambaram, Arindam Jati, Nam H Nguyen, Pankaj Dayama, Chandra Reddy, Wesley M Gifford, and Jayant Kalagnanam. 2024. Tiny Time Mixers (TTMs): Fast Pre-trained Models for Enhanced Zero/Few-Shot Forecasting of Multivariate Time Series. *arXiv preprint arXiv:2401.03955* (2024).
- [21] Shereen Elsayed, Daniela Thyssens, Ahmed Rashed, Hadi Samer Jomaa, and Lars Schmidt-Thieme. 2021. Do we really need deep learning models for time series forecasting? *arXiv preprint arXiv:2101.02118* (2021).
- [22] William Falcon and The PyTorch Lightning team. 2019. PyTorch Lightning. <https://doi.org/10.5281/zenodo.3828935>
- [23] Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and Marinka Zitnik. 2024. UniTS: Building a Unified Time Series Model. *arXiv preprint arXiv:2403.00131* (2024).
- [24] Tilmann Gneiting and Matthias Katzfuss. 2014. Probabilistic Forecasting. *Annual Review of Statistics and Its Application* (2014).
- [25] Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski. 2024. MOMENT: A Family of Open Time-series Foundation Models. In *ICML*.
- [26] Luis Hernandez, Carlos Baladron, Javier M Aguiar, Belén Carro, Antonio J Sanchez-Esguevillas, Jaime Lloret, and Joaquim Massana. 2014. A survey on electric power demand forecasting: Future trends in smart grids, microgrids and smart buildings. *IEEE Communications Surveys & Tutorials* (2014).
- [27] Julien Herzen, Francesco LÃ¶ssig, Samuele Giuliano Piazzetta, Thomas Neuer, LÃ©o Tafti, Guillaume Raille, Tomas Van Pottelbergh, Marek Pasiaka, Andrzej Skrodzki, Nicolas Huguenin, Maxime Dumonal, Jan KoÅcisz, Dennis Bader, FrÃ©dÃ©rick Gusset, Mounir Benheddi, Camila Williamson, Michal Kosinski, Matej Petrik, and GaÅl Grosch. 2022. Darts: User-Friendly Modern Machine Learning for Time Series. *JMLR* 23, 124 (2022), 1–6.
- [28] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic Models. In *NeurIPS*.
- [29] Min Hou, Chang Xu, Yang Liu, Weiqing Liu, Jiang Bian, Le Wu, Zhi Li, Enhong Chen, and Tie-Yan Liu. 2021. Stock trend prediction with multi-granularity data: A contrastive learning approach with adaptive fusion. In *CIKM*.
- [30] Rob J Hyndman and George Athanasopoulos. 2018. *Forecasting: principles and practice*.
- [31] Xiaodong Jiang, Sudeep Srivastava, Sourav Chatterjee, Yang Yu, Jeffrey Handler, Peiyi Zhang, Rohan Bopardikar, Dawei Li, Yanjun Lin, Uttam Thakore, Michael Brundage, Ginger Holt, Caner Komurlu, Rakshita Nagalla, Zhichao Wang, Hechao Sun, Peng Gao, Wei Cheung, Jun Gao, Qi Wang, Marius Guerard, Morteza Kazemi, Yulin Chen, Chong Zhou, Sean Lee, Nikolay Laptev, Tihamér Levendovszky, Jake Taylor, Huijun Qian, Jian Zhang, Aida Shoydokova, Trisha

- Singh, Chengjun Zhu, Zeynep Baz, Christoph Bergmeir, Di Yu, Ahmet Koylan, Kun Jiang, Ploy Temiyasathit, and Emre Yurtbay. 2022. *Kats*. <https://github.com/facebookresearch/Kats>
- [32] Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. 2022. Reversible Instance Normalization for Accurate Time-Series Forecasting against Distribution Shift. In *ICLR*.
- [33] Marcel Kollovich, Abdul Fatir Ansari, Michael Bohlke-Schneider, Jasper Zschiegner, Hao Wang, and Yuyang Wang. 2023. Predict, Refine, Synthesize: Self-Guiding Diffusion Models for Probabilistic Time Series Forecasting. In *NeurIPS*.
- [34] Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. 2018. Modeling Long- and Short-Term Temporal Patterns with Deep Neural Networks. In *SIGIR*. 95–104.
- [35] Yan Li, Xinjiang Lu, Yaqing Wang, and Dejing Dou. 2022. Generative Time Series Forecasting with Diffusion, Denoise, and Disentanglement. In *NeurIPS*.
- [36] Bryan Lim, Sercan Ö Arık, Nicolas Loeff, and Tomas Pfister. 2021. Temporal Fusion Transformers for Interpretable Multi-horizon Time Series Forecasting. *International Journal of Forecasting* (2021).
- [37] Minhao LIU, Ailing Zeng, Muxi Chen, Zhijian Xu, Qiuxia LAI, Lingna Ma, and Qiang Xu. 2022. SCINet: Time Series Modeling and Forecasting with Sample Convolution and Interaction. In *NeurIPS*.
- [38] Shizhan Liu, Hang Yu, Cong Liao, Jianguo Li, Weiyao Lin, Alex X. Liu, and Schahram Dustdar. 2022. Pyraformer: Low-Complexity Pyramidal Attention for Long-Range Time Series Modeling and Forecasting. In *ICLR*.
- [39] Yang Liu, Jiashun Cheng, Haihong Zhao, Tingyang Xu, Peilin Zhao, Fugee Tsung, Jia Li, and Yu Rong. 2024. SEGNO: Generalizing Equivariant Graph Neural Networks with Physical Inductive Biases. In *ICLR*.
- [40] Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long. 2024. itransformer: Inverted transformers are effective for time series forecasting. In *ICLR*.
- [41] Yong Liu, Haoran Zhang, Chenyu Li, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. 2024. Timer: Transformers for Time Series Analysis at Scale. In *ICML*.
- [42] Markus Löning, Anthony Bagnall, Sajaysurya Ganesh, Viktor Kazakov, Jason Lines, and Franz J Király. 2019. sktime: A unified interface for machine learning with time series. *arXiv preprint arXiv:1909.07872* (2019).
- [43] Yisheng Lv, Yanjie Duan, Wenwen Kang, Zhengxi Li, and Fei-Yue Wang. 2014. Traffic flow prediction with big data: A deep learning approach. *IEEE Transactions on Intelligent Transportation Systems* (2014).
- [44] Spyros Makridakis and Michele Hibon. 1997. ARMA models and the Box–Jenkins methodology. *Journal of forecasting* 16, 3 (1997), 147–163.
- [45] Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. 2020. The M4 Competition: 100,000 time series and 61 forecasting methods. *International Journal of Forecasting* 36, 1 (2020), 54–74.
- [46] Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. 2022. The M5 competition: Background, organization, and implementation. *International Journal of Forecasting* 38, 4 (2022), 1325–1336.
- [47] James E Matheson and Robert L Winkler. 1976. Scoring Rules for Continuous Probability Distributions. *Management Science* 22, 10 (1976), 1087–1096.
- [48] Manfred Mudelsee. 2019. Trend analysis of climate time series: A review of methods. *Earth-science reviews* (2019).

- [49] Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. 2023. A Time Series is Worth 64 Words: Long-term Forecasting with Transformers. In *ICLR*.
- [50] Frank Nielsen. 2019. On the Jensen–Shannon symmetrization of distances relying on abstract means. *Entropy* 21, 5 (2019), 485.
- [51] Nixtla. 2024. NeuralForecast. <https://github.com/Nixtla/neuralforecast>. GitHub repository.
- [52] Ignacio Oguiza. [n. d.]. tsai - A State-of-the-art Deep Learning Library for Time Series and Sequential Data. Github. <https://github.com/timeseriesAI/tsai>
- [53] Boris N. Oreshkin, Dmitri Carпов, Nicolas Chapados, and Yoshua Bengio. 2020. N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. In *ICLR*.
- [54] Xiangfei Qiu, Jilin Hu, Lekui Zhou, Xingjian Wu, Junyang Du, Buang Zhang, Chenjuan Guo, Aoying Zhou, Christian S. Jensen, Zhenli Sheng, and Bin Yang. 2024. TFB: Towards Comprehensive and Fair Benchmarking of Time Series Forecasting Methods. *VLDB* 17, 9 (2024), 2363–2377.
- [55] Syama Sundar Rangapuram, Matthias W. Seeger, Jan Gasthaus, Lorenzo Stella, Yuyang Wang, and Tim Januschowski. 2018. Deep State Space Models for Time Series Forecasting. In *NeurIPS*. 7796–7805.
- [56] Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Hena Ghonia, Rishika Bhagwatkar, Arian Khorasani, Mohammad Javad Darvishi Bayazi, George Adamopoulos, Roland Riachi, Nadhir Hassen, et al. 2023. Lag-Llama: Towards Foundation Models for Probabilistic Time Series Forecasting. *arXiv preprint arXiv:2310.08278* (2023).
- [57] Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. 2021. Autoregressive Denoising Diffusion Models for Multivariate Probabilistic Time Series Forecasting. In *ICML*. 8857–8868.
- [58] Kashif Rasul, Abdul-Saboor Sheikh, Ingmar Schuster, Urs M. Bergmann, and Roland Vollgraf. 2021. Multivariate Probabilistic Time Series Forecasting via Conditioned Normalizing Flows. In *ICLR*.
- [59] David Salinas, Michael Bohlke-Schneider, Laurent Callot, Roberto Medico, and Jan Gasthaus. 2019. High-dimensional Multivariate Forecasting with Low-rank Gaussian Copula Processes. In *NeurIPS*. 6824–6834.
- [60] David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. 2020. DeepAR: Probabilistic Forecasting with Autoregressive Recurrent Networks. *International Journal of Forecasting* 36, 3 (2020), 1181–1191.
- [61] Fei Su, Honghui Dong, Limin Jia, Yong Qin, and Zhao Tian. 2016. Long-term forecasting oriented to urban expressway traffic situation. *Advances in mechanical engineering* (2016).
- [62] Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. 2021. CSDI: Conditional Score-based Diffusion Models for Probabilistic Time Series Imputation. In *NeurIPS*. 24804–24816.
- [63] Sean J Taylor and Benjamin Letham. 2018. Forecasting at Scale. *The American Statistician* (2018).
- [64] Pytorch Forecasting Team. 2024. PyTorch Forecasting. <http://github.com/jdb78/pytorch-forecasting>. GitHub repository.
- [65] Pytorch Transformer TS Team. 2024. Pytorch Transformer based Time Series Models. <https://github.com/kashif/pytorch-transformer-ts>. GitHub repository.
- [66] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *NeurIPS*. 5998–6008.
- [67] Huaizhi Wang, Zhenxing Lei, Xian Zhang, Bin Zhou, and Jianchun Peng. 2019. A review of deep learning for renewable energy forecasting. *Energy Conversion and Management* (2019).

- [68] Xiaozhe Wang, Kate Smith, and Rob Hyndman. 2006. Characteristic-based Clustering for Time Series Data. *Data Mining and Knowledge Discovery* 13 (2006), 335–364.
- [69] Ruofeng Wen, Kari Torkkola, Balakrishnan Narayanaswamy, and Dhruv Madeka. 2017. A Multi-horizon Quantile Recurrent Forecaster. *arXiv preprint arXiv:1711.11053* (2017).
- [70] Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo. 2024. Unified Training of Universal Time Series Forecasting Transformers. In *ICML*.
- [71] Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. 2023. TimesNet: Temporal 2D-Variation Modeling for General Time Series Analysis. In *ICLR*.
- [72] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. 2021. Autoformer: Decomposition Transformers with Auto-Correlation for Long-Term Series Forecasting. In *NeurIPS*. 22419–22430.
- [73] Zhijian Xu, Ailing Zeng, and Qiang Xu. 2024. FITS: Modeling Time Series with 10^4 Parameters. In *ICLR*.
- [74] Jiexia Ye, Weiqi Zhang, Ke Yi, Yongzi Yu, Ziyue Li, Jia Li, and Fugee Tsung. 2024. A Survey of Time Series Foundation Models: Generalizing Time Series Representation with Large Language Mode. *arXiv preprint arXiv:2405.02358* (2024).
- [75] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. 2023. Are Transformers Effective for Time Series Forecasting?. In *AAAI*. 11121–11128.
- [76] Tianping Zhang, Yizhuo Zhang, Wei Cao, Jiang Bian, Xiaohan Yi, Shun Zheng, and Jian Li. 2022. Less is More: Fast Multivariate Time Series Forecasting with Light Sampling-Oriented MLP Structures. *arXiv preprint arXiv:2207.01186* (2022).
- [77] Weiqi Zhang, Jianfeng Zhang, Jia Li, and Fugee Tsung. 2023. A co-training approach for noisy time series learning. In *CIKM*. 3308–3318.
- [78] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. 2021. Informer: Beyond Efficient Transformer for Long Sequence Time-series Forecasting. In *AAAI*. 11106–11115.

A Additional Related Work

A.1 A Comparison with Traditional Models on Covering Essential Forecasting Needs and methodological decisions

Table 2 presents a comparative summary of our approach, which adopts an integrated perspective, and representative studies from the existing literature.

Table 2: We provide a concise comparison between the methodologies presented in this paper and those from two distinct research branches. The comparison is based on data scenarios (short-term versus long-term forecasting), primary evaluation metrics (point versus distributional forecasts), and key methodological choices (general or customized neural architecture designs, and autoregressive or non-autoregressive decoding schemes).

Method	Pred. Horizon		Paradigm		Arch. Design		Dec. Scheme	
	Short	Long	Point	Distr.	General	Customized	AR	Non-AR
N-BEATS [53]	✗	✓	✓	✗	✗	✓	✗	✓
Autoformer [72]	✗	✓	✓	✗	✗	✓	✗	✓
Informer [78]	✗	✓	✓	✗	✗	✓	✗	✓
Pyraformer [38]	✗	✓	✓	✗	✗	✓	✗	✓
N-HiTS [11]	✗	✓	✓	✗	✗	✓	✗	✓
LTSF-Linear [75]	✗	✓	✓	✗	✗	✓	✗	✓
PatchTST [49]	✗	✓	✓	✗	✗	✓	✗	✓
TimesNet [71]	✓	✓	✓	✗	✗	✓	✗	✓
iTransformer [40]	✗	✓	✓	✗	✗	✓	✗	✓
DeepAR [60]	✓	✗	✗	✓	✓	✗	✓	✗
GP-copula [59]	✓	✗	✗	✓	✓	✗	✓	✗
LSTM NVP [58]	✓	✗	✗	✓	✓	✗	✓	✗
LSTM MAF [58]	✓	✗	✗	✓	✓	✗	✓	✗
Trans MAF [58]	✓	✗	✗	✓	✓	✗	✓	✗
TimeGrad [57]	✓	✗	✗	✓	✓	✗	✓	✗
CSDI [62]	✓	✗	✗	✓	✗	✓	✗	✓
SPD [7]	✓	✗	✗	✓	✓	✗	✗	✓
TSDiff [33]	✓	✗	✗	✓	✗	✓	✗	✓
This Study	✓	✓	✓	✓	✓	✓	✓	✓

A.2 A Comparison of Pre-trained Time-series Foundation Models

We have incorporated eight recently emerged time series foundation models, namely Lag-Llama [56], Chronos [2], TimesFM [15], Timer [41], MOIRAI [70], UniTS [23], ForecastPFN [18], and TTM [20], into our framework. These foundation models are categorized based on their capabilities, such as zero-shot forecasting, adaptability to varying prediction lengths, and support for probabilistic predictions, as well as their architectural designs, including whether they are auto-regressive and the nature of their backbone networks. Additionally, we have detailed their training processes, including the lengths of prediction horizons used during pre-training and the sizes of look-back windows. These details are summarized in Table 3.

Furthermore, we have compiled a summary of these foundation models’ pre-training and evaluation on several classical time series forecasting datasets. This compilation is presented in Table 4.

A.3 A Comparison with Existing Libraries on the Coverage of Data, Model, and Metrics

ProbTS is a research toolkit designed to advance forecasting research across varied horizons, focusing on both point and distributional forecasting. To achieve these objectives, ProbTS includes state-of-the-art models, comprehensive evaluation protocols (point vs. distributional), and explores different methodological aspects of forecasting models, particularly in terms of distributional estimation

Table 3: Foundation Models for Time Series. **Zero-shot** indicates whether the original work tests zero-shot capabilities. **Any-horizon** indicates if the same pre-trained model can adapt to prediction tasks of varying lengths. **AR** denotes if the model performs auto-regressive forecasting. **Prob.** indicates if the model natively supports probabilistic forecasting. **Arch.** denotes the model’s backbone architecture: D-O for decoder-only transformer, E-O for encoder-only transformer, E-D for encoder-decoder transformer, and unique for specially designed backbones. **Multi-variate** indicates if the model explicitly handles multivariate relationships. **Pre-train Horizons** specifies the forecasting task horizons during pre-training. **Look-back Window** specifies the context history length settings used in the original experiments.

Model	Zero-shot	Any-horizon	AR	Prob.	Arch.	Multi-variate	Pre-train Horizons	Look-back Window
Lag-Llama	✓	✓	✓	✓	D-O	✗	24~60	32~1024
Chronos	✓	✓	✓	✓	E-D	✗	64	512
TimesFM	✓	✓	✓	✗	D-O	✗	–	512
Timer	✓	✓	✓	✗	D-O	✗	up to 1440	672
MOIRAI	✓	✓	✗	✓	E-O	○	varying	100~5000
UniTS	✓	✓	✗	✗	E-O	✗	–	60~720
ForecastPFN	✓	✓	✗	✗	E-O	✗	0 50	50 500
TTM	✓	✗	✗	✗	Unique	✓	96	512

Table 4: Evaluation Datasets for Time-series Foundation Models. We selected several popular datasets to evaluate time-series foundation models. ✓ indicates pre-training on the dataset, ○ indicates zero-shot evaluation on the dataset, few indicates few-shot evaluation on the dataset, and ✗ indicates the dataset is not mentioned in the paper or documentation. ‘*’ indicates that the data comes from the same source but may be processed differently.

Model	Solar	Wikipedia	ETTh1	ETTm2	ETTh1	ETTh2	Electricity	Traffic	Weather	Exchange	ILI
MOIRAI	○	✓	○	○	○	○	○	✓	○	✗	✗
Lag-Llama	✓*	✗	✓	○	✓	✓	✓	✓	○	○	✗
TimesFM	○*	✓*	○	○	○	○	✓	✓	✓	✗	✗
Chronos	✓	✓*	○	○	○	○	✓	○	○	○	✗
TTM	✗	✗	○	○	○	○	○	○	○	✗	✗
UniTS	○	✗	<i>few</i>	✗	✓	<i>few</i>	✓	✓	✓	✓	✓
Timer	✗	✗	<i>few</i>	<i>few</i>	<i>few</i>	<i>few</i>	<i>few</i>	<i>few</i>	<i>few</i>	✗	✗

methods and decoding schemes (AR vs. NAR). In Table 5, we provide a comprehensive comparison of ProbTS with existing libraries in terms of toolkit functionalities and the benchmarking aspects we aim to investigate.

B More Details on ProbTS

B.1 Data

The data module unifies varied data scenarios to facilitate thorough evaluation and implements standardized pre-processing techniques to ensure fair comparison.

Moreover, we utilize a quantitative approach to visually delineate datasets’ intrinsic characteristics, which employs decomposition to assess trends and seasonality in a time series and evaluate the similarity between data distribution and a Gaussian to depict the complexity of data distribution.

B.1.1 Time-series Forecasting Datasets

Table 6 provides a summary of the public datasets employed in our study. These datasets have been sourced from recent research studies in the field of deep time-series forecasting.

Table 5: Comparison of Various Time Series Tools.

Tools	Features		Benchmarking			
	SOTA Model	Dist. Evaluation	Short-term	Long-term	AR vs. NAR	Point vs Prob.
Merlion [6]	✗	✗	✓	✗	✗	✗
Kats [31]	✗	✗	✗	✗	✗	✗
pytorch-transformer-ts [65]	✗	✗	✗	✗	✗	✗
Prophet [63]	✗	✗	✓	✗	✗	✗
Darts [27]	✗	✓	✗	✗	✗	✗
sktime [42]	✗	✓	✗	✗	✗	✗
pytorch-forecasting [64]	✗	✓	✗	✗	✗	✗
NeuralForecast [51]	✓	✓	✗	✗	✗	✗
tsai [52]	✓	✗	✗	✗	✗	✗
TFB [54]	✓	✗	✗	✓	✗	✗
TSlib [71]	✓	✗	✓	✓	✓	✗
GluonTS [1]	✓	✓	✓	✗	✗	✓
ProbTS	✓	✓	✓	✓	✓	✓

Table 6: Dataset Summary.

Horizon	Dataset	#var.	range	freq.	timesteps	Description
Long-term	ETTh1/h2	7	\mathbb{R}^+	H	17,420	Electricity transformer temperature per hour
	ETTm1/m2	7	\mathbb{R}^+	15min	69,680	Electricity transformer temperature every 15 min
	Electricity	321	\mathbb{R}^+	H	26,304	Electricity consumption (Kwh)
	Traffic	862	(0,1)	H	17,544	Road occupancy rates
	Exchange	8	\mathbb{R}^+	Busi. Day	7,588	Daily exchange rates of 8 countries
	ILI	7	(0,1)	W	966	Ratio of patients seen with influenza-like illness
	Weather	21	\mathbb{R}^+	10min	52,696	Local climatological data
Short-term	Exchange	8	\mathbb{R}^+	Busi. Day	6,071	Daily exchange rates of 8 countries
	Solar	137	\mathbb{R}^+	H	7,009	Solar power production records
	Electricity	370	\mathbb{R}^+	H	5,833	Electricity consumption
	Traffic	963	(0,1)	H	4,001	Road occupancy rates
	Wikipedia	2,000	N	D	792	Page views of 2000 Wikipedia pages

B.1.2 Data Visualization

To provide a more tangible understanding of the different forecasting scenarios, we visualize time-series segments from both short-term and long-term forecasting datasets. The segments’ window size is determined by the specific forecasting setup.

In Figure 5, we present samples extracted from short-term forecasting scenarios. At this scale, the series primarily exhibit local variations, and the compact window size often obscures pronounced seasonal or trending patterns. However, these short-term scenarios may reveal irregularly varied patterns, suggesting a more complex underlying data distribution.

On the contrary, Figure 6 illustrates long-term forecasting scenarios. With extended forecasting horizons, as showcased in datasets like Traffic, Electricity, and ETT, the series display more pronounced seasonality and trends. These characteristics render the series more regular patterns in the long-term scenarios.

It’s important to note that these visualizations are not selectively chosen or "cherry-picked". We have depicted multiple time-series segments from various time steps, and the observed patterns remain consistent across these different instances.

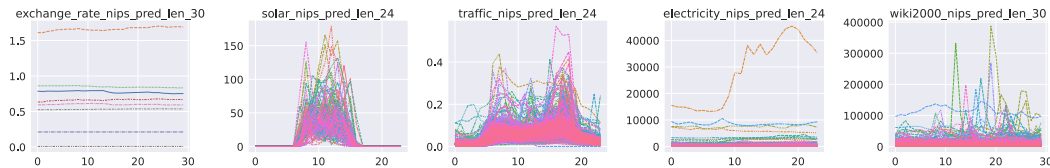


Figure 5: We have sampled and visualized multiple time-series segments from the short-term forecasting datasets. The size of the segment window is set equal to the prediction horizon.

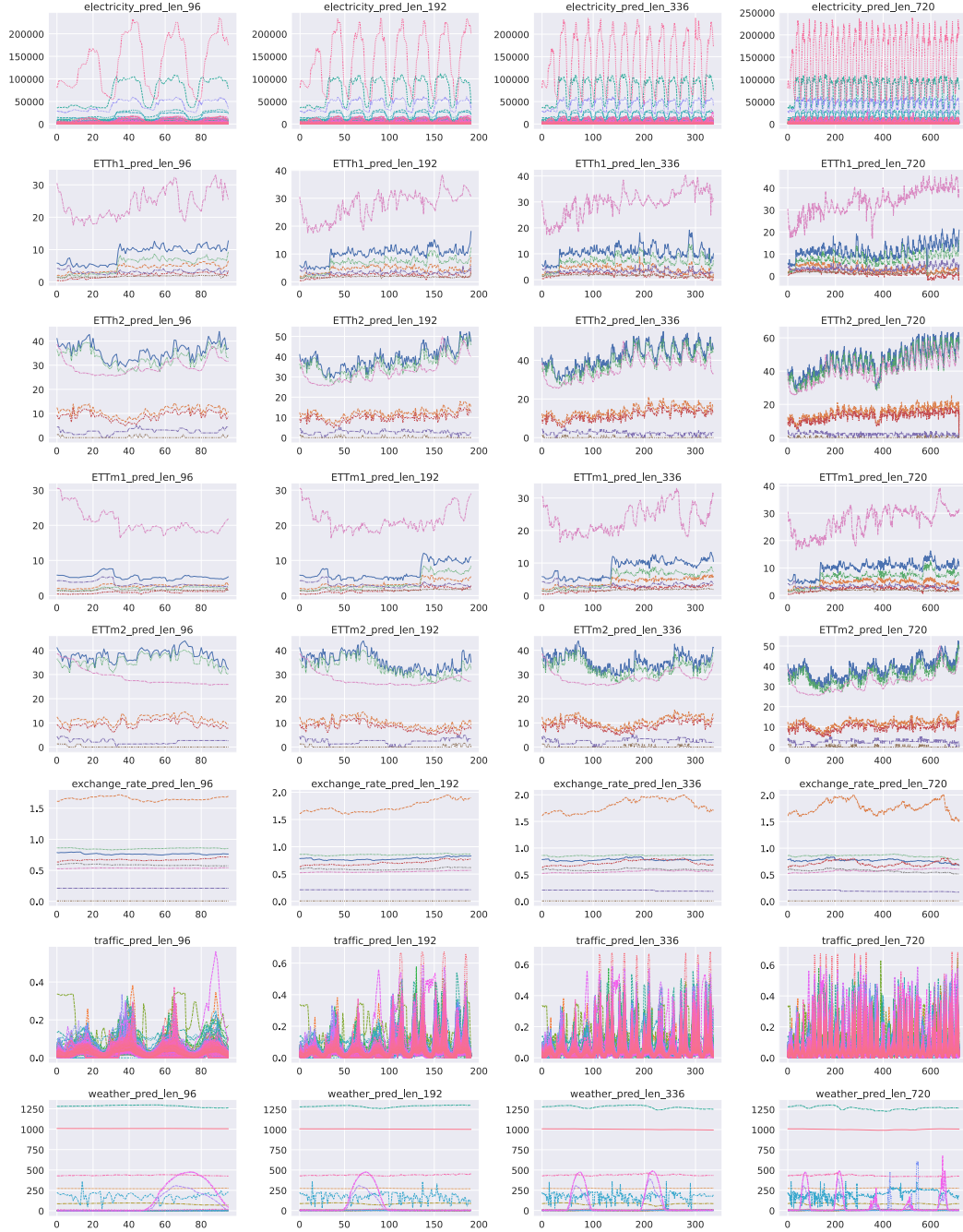


Figure 6: We have also sampled and visualized multiple time-series segments from the long-term forecasting datasets, where the size of the segment window matches the prediction horizon.

B.1.3 Quantifying Trend and Seasonality Strengths

Based on the intuition obtained from data visualization, we would like to quantify the strengths of trend and seasonality for a time-series segment with a predefined window size (corresponding to the prediction horizon). Then we can quantify the trend and seasonality strengths at the dataset level by averaging over all time-series segments of a dataset.

To quantify the strengths of trend and seasonality for a fixed-length time-series segment, we draw upon methodologies outlined in the work of [68]. In particular, we employed a time series decomposition

model expressed as:

$$y_t = T_t + S_t + R_t,$$

where T_t represents the smoothed trend component, S_t signifies the seasonal component, and R_t denotes the remainder component. In order to obtain each component, we followed the STL decomposition approach ⁴.

In the case of strongly trended data, the variation within the seasonally adjusted data should considerably exceed that of the remainder component. Consequently, the ratio $\text{Var}(R_t)/\text{Var}(T_t + R_t)$ is expected to be relatively small. As such, the measure of trend strength can be formulated as:

$$F_T = \max\left(0, 1 - \frac{\text{Var}(R_t)}{\text{Var}(T_t + R_t)}\right).$$

The quantified trend strength, ranging from 0 to 1, characterizes the degree of trend presence. Similarly, the evaluation of seasonal intensity employs the detrended data:

$$F_S = \max\left(0, 1 - \frac{\text{Var}(R_t)}{\text{Var}(S_t + R_t)}\right).$$

A series with F_S near 0 indicates minimal seasonality, while strong seasonality is indicated by F_S approaching 1 due to the considerably smaller variance of $\text{Var}(R_t)$ in comparison to $\text{Var}(S_t + R_t)$.

Tables 1 depict the results for each dataset. Notably, the ETT datasets and the Exchange dataset manifest conspicuous trends, whereas the Electricity, Solar, and Traffic datasets showcase marked seasonality. Additionally, the Exchange dataset stands out with distinctive features. Figure 6 also illustrates that with shorter prediction windows, the Exchange dataset sustains comparatively minor fluctuations, almost forming a linear trajectory. This enables effective forecasting through a straightforward batch mean approach. As the forecasting horizon extends, the dataset appears a more pronounced trend while retaining minimal seasonality.

B.1.4 Quantifying Data Distribution Complexity

To differentiate between methods optimized for point or distributional forecasts, we aim to quantify the complexity of data distribution within a time-series segment whose window size equals the prediction horizon length. Such complexities may arise from the unpredictability of the data itself or from noises accidentally introduced during the data collection process [77].

We propose that assessing non-Gaussianity, i.e., how closely the distribution of time-series values within that window resembles a Gaussian distribution, could serve as a meaningful measure. This is because point forecasting methods, optimized with mean squared loss, are essentially equivalent to probabilistic counterparts that include a Gaussian output head and employ maximum a posteriori estimation. This suggests that point forecasting methods inherently assume that time-series values adhere to a Gaussian distribution. In contrast, advanced probabilistic methods, which do not make prior assumptions about data distribution, can adapt to complex data distributions in a data-driven manner.

Hence, we use the Jensen–Shannon divergence [50] to measure the similarity between the actual value distribution of a time-series segment and a Gaussian distribution fitted to the observed values. Short-term datasets used a window size of 30, while long-term datasets used a size of 336. By averaging the calculated divergence values across all time-series segments of a dataset, we obtain a dataset-level measure of non-Gaussianity. A larger divergence value indicates a larger deviation from a Gaussian distribution in the data.

B.2 Metrics

In ProbTS, we integrate an extensive variety of metrics that take into account both point and distributional forecasts, thereby providing a comprehensive and multifaceted assessment of forecasting models.

⁴<https://otexts.com/fpp2/stl.html>

B.2.1 Metrics for Point Forecasts

Mean Absolute Error (MAE) The Mean Absolute Error (MAE) quantifies the average absolute deviation between the forecasts and the true values. Since it averages the absolute errors, MAE is robust to outliers. Its mathematical formula is given by:

$$\text{MAE} = \frac{1}{K \times T} \sum_{k=1}^K \sum_{t=1}^T |x_t^k - \hat{x}_t^k|,$$

where K is the number of variates, T is the length of series, x_t^k and \hat{x}_t^k denotes the ground-truth value and the predicted value, respectively. For multivariate time series, we also provide the aggregated version:

$$\text{MAE}_{\text{sum}} = \frac{1}{T} \sum_{t=1}^T |x_t^{\text{sum}} - \hat{x}_t^{\text{sum}}|,$$

where x_t^{sum} and \hat{x}_t^{sum} are the summation across the dimension K of x_t^k and \hat{x}_t^k , respectively.

Normalized Mean Absolute Error (NMAE) The Normalized Mean Absolute Error (NMAE) is a normalized version of the MAE, which is dimensionless and facilitates the comparability of the error magnitude across different datasets or scales. The mathematical representation of NMAE is given by:

$$\text{NMAE} = \frac{\sum_{k=1}^K \sum_{t=1}^T |x_t^k - \hat{x}_t^k|}{\sum_{k=1}^K \sum_{t=1}^T |x_t^k|}.$$

Its aggregated version is:

$$\text{NMAE}_{\text{sum}} = \frac{\sum_{t=1}^T |x_t^{\text{sum}} - \hat{x}_t^{\text{sum}}|}{\sum_{t=1}^T |x_t^{\text{sum}}|}.$$

Mean Squared Error (MSE) The Mean Squared Error (MSE) is a quantitative metric used to measure the average squared difference between the observed actual value and forecasts. It is defined mathematically as follows:

$$\text{MSE} = \frac{1}{K \times T} \sum_{k=1}^K \sum_{t=1}^T (x_t^k - \hat{x}_t^k)^2.$$

For multivariate time series, we also provide the aggregated version:

$$\text{MSE}_{\text{sum}} = \frac{1}{T} \sum_{t=1}^T (x_t^{\text{sum}} - \hat{x}_t^{\text{sum}})^2.$$

Normalized Root Mean Squared Error (NRMSE) The Normalized Root Mean Squared Error (NRMSE) is a normalized version of the Root Mean Squared Error (RMSE), which quantifies the average squared magnitude of the error between forecasts and observations, normalized by the expectation of the observed values. It can be formally written as:

$$\text{NRMSE} = \frac{\sqrt{\frac{1}{K \times T} \sum_{k=1}^K \sum_{t=1}^T (x_t^k - \hat{x}_t^k)^2}}{\frac{1}{K \times T} \sum_{k=1}^K \sum_{t=1}^T |x_t^k|}.$$

For multivariate time series, we also provide the aggregated version:

$$\text{NRMSE}_{\text{sum}} = \frac{\sqrt{\frac{1}{T} \sum_{t=1}^T (x_t^{\text{sum}} - \hat{x}_t^{\text{sum}})^2}}{\frac{1}{T} \sum_{t=1}^T |x_t^{\text{sum}}|}.$$

Mean Absolute Scaled Error (MASE) The Mean Absolute Scaled Error (MASE) divides the MAE of forecasted values by MAE of the in-sample one-step naive forecast, which is a scale-invariant metrics:

$$\text{MASE} = \frac{\frac{1}{K \times T} \sum_{k=1}^K \sum_{t=1}^T |x_t^k - \hat{x}_t^k|}{\frac{1}{K \times T} \sum_{k=1}^K \sum_{t=1}^T |x_t^k - x_{t-1}^k|}.$$

B.2.2 Metrics for Distributional Forecasts

Continuous Ranked Probability Score (CRPS) The Continuous Ranked Probability Score (CRPS) [47] quantifies the agreement between a cumulative distribution function (CDF) F and an observation x , represented as:

$$\text{CRPS} = \int_{\mathbb{R}} (F(z) - \mathbb{I}\{x \leq z\})^2 dz,$$

where $\mathbb{I}\{x \leq z\}$ denotes the indicator function, equating to one if $x \leq z$ and zero otherwise.

Being a proper scoring function, CRPS reaches its minimum when the predictive distribution F coincides with the data distribution. When using the empirical CDF of F , denoted as $\hat{F}(z) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{X_i \leq z\}$, where n represents the number of samples $X_i \sim F$, CRPS can be precisely calculated from the simulated samples of the conditional distribution $p_{\theta}(\mathbf{x}_t | \mathbf{h}_t)$. In our practice, 100 samples are employed to estimate the empirical CDF.

For multivariate time series, the aggregate CRPS, denoted as CRPS_{sum} , is derived by summing across the K time series, both for the ground-truth data and sampled data, and subsequently averaging over the forecasting horizon. Formally, it is represented as:

$$\text{CRPS}_{\text{sum}} = \mathbb{E}_t \left[\text{CRPS} \left(\hat{F}_{\text{sum}}(t), \sum_{i=1}^K x_{i,l}^0 \right) \right].$$

B.3 Baselines

To ensure the integrity of the results, ProbTS adheres to a standard implementation process, employing unified data splitting, standardization techniques, and adopting fair settings for hyperparameter tuning across all methods.

Implementation Details ProbTS was developed using PyTorch Lightning [22]. During training, we sampled 100 batches per epoch and limited training to 50 epochs, using the CRPS metric for checkpointing. All experiments employed the Adam optimizer and were run on single NVIDIA Tesla V100 GPUs with CUDA 11.3. To enable evaluation of distribution-level metrics, we conducted 100 samplings to calculate metrics on the test set.

Following the most commonly adopted settings [75, 49, 71], in the long-term forecasting context, all of the models are following the same experimental setup with prediction length $T \in \{24, 36, 48, 60\}$ for ILI-L dataset and $T \in \{96, 192, 336, 720\}$ for other datasets. Note that the lookback window here is 96 for all the models, to ensure a fair comparison. In the short-term forecasting context, the length of the lookback window is the same as the forecasting horizons, which are 30 for Exchange-S dataset and Wikipedia-S dataset, and 24 for the rest, the same as [59].

Hyper-parameter Tuning For a fair comparison, we conducted a comprehensive grid search for critical hyperparameters across all models in this study. Table 7 details the shared hyperparameters tuned within the ProbTS pipeline, along with those kept constant. Due to the vast array of model-specific hyperparameters, we present an example configuration in Table 8. Complete hyperparameter configurations for each model, identified through this process, will be made available in a public GitHub repository for transparency and reproducibility.

Implementation Details on Foundation Models We used reference implementations of eight time series foundation models into ProbTS.

For Lag-Llama [56], we use its official code⁵ and integrate the `LagLlamaEstimator` with its pre-trained checkpoint. The look-back window is uniformly set to 512, irrespective of the forecast horizon. The other hyper-parameters are aligned with the recommended settings.

For Chronos [2], we use its official code⁶ and integrate the `ChronosPipeline` into ProbTS using `amazon/chronos-t5` checkpoints (three models were tested: small, base and large). Two look-back

⁵<https://github.com/time-series-foundation-models/lag-llama>

⁶<https://github.com/amazon-science/chronos-forecasting>

Table 7: Hyper-parameters values fixed or range searched in hyper-parameter tuning.

Hyper-parameter	Value or Range Searched
learning rate	[1e-4, 1e-3, 1e-2]
dropout	[0, 0.1, 0.2]
batch_size	[8, 16, 32, 64]
use_lags	[True, False]
use_feat_idx_emb	[True, False]
use_time_feat	[True, False]
autoregressive	[True, False]
scaler	[Standard, Scaling, None]
limit_train_batches	100
num_samples	100
quantiles_num	20

Table 8: Hyperparameter settings for Electricity-S dataset.

Model	Hyperparameter
DLinear	learning_rate=0.01, kernel_size=3, f_hidden_size=40
PatchTST	learning_rate=0.0001, stride=3, patch_len=6, n_layers=3, n_heads=8, dropout=0.1, kernel_size=3, f_hidden_size=32
TimesNet	learning_rate=0.001, n_layers=2, num_kernels=6, top_k=5, f_hidden_size=64, d_ff=64
GRU NVP	learning_rate=0.001, f_hidden_size=40, num_layers=2, n_blocks=3, hidden_size=100, conditional_length=200
GRU MAF	learning_rate=0.001, f_hidden_size=40, num_layers=2, n_blocks=4, hidden_size=100, conditional_length=200
Trans MAF	learning_rate=0.001, f_hidden_size=32, num_heads=8, n_blocks=4, hidden_size=100, conditional_length=200
TimeGrad	learning_rate=0.001, f_hidden_size=128, num_layers=4, conditional_length=100, beta_end=0.1, diff_steps=100
CSDI	learning_rate=0.001, channels=64, emb_time_dim=128, emb_feature_dim=16, num_steps=50, num_heads=8, n_layers=4

windows are used during evaluation: 96 and 512. We set `limit_prediction_length=False` to enable it to predict horizons longer than 64. However, this may potentially lead to a decrease in predictive performance since the model was only trained to consider prediction lengths of 64 or less during pre-training.

For TimesFM [15], we modify its official code⁷ into ProbTS and load checkpoints from `google/timesfm-1.0-200m`. Look-back window is set to 96 for a fair comparison.

For Timer [41], we modify its official code⁸ into ProbTS and `Timer_67M_UTSD_4G` checkpoint downloaded from its repo. Look-back window is also set to 96 for a fair comparison.

For MOIRAI [70], we employ its official code⁹ and load checkpoints from `Salesforce/moirai-1.0-R-base`. Two look-back windows are utilized during evaluation: 96 and 5000. The original experiments suggest that MOIRAI’s forecasting capability can be consistently enhanced by increasing the look-back window. Consequently, we have included a 5000 look-back window to test the model’s performance.

For UniTS [23], we have adapted its official code¹⁰ into our ProbTS framework and loaded the `saved_weights` from its repo. The look-back window is set to 96 to ensure a fair comparison. It is important to note that this checkpoint was originally used for the Zero-Shot New-length Forecasting experiment in the original work (where models are challenged to predict new lengths by adjusting from the trained length, with offsets ranging from 0 to 384), which differs from the objectives of our experiments.

For ForecastPFN [18], we have integrated its official code¹¹ into our ProbTS framework and have utilized the `units_x128_pretrain_checkpoint`. However, we have encountered some challenges

⁷<https://github.com/google-research/timesfm>

⁸<https://github.com/thuml/Large-Time-Series-Model>

⁹<https://github.com/SalesforceAIRResearch/uni2ts>

¹⁰<https://github.com/abacusai/ForecastPFN>

¹¹<https://github.com/SalesforceAIRResearch/uni2ts>

Table 9: Results (mean_{std}) on short-term forecasting scenarios, each containing five independent runs with different seeds.

Model	Exchange-S		Solar-S		Electricity-S		Traffic-S		Wikipedia-S	
	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE
Glob. mean	0.188	0.188	1.403	1.403	0.412	0.412	0.540	0.540	0.577	0.577
Batch mean	0.012	0.012	1.244	1.244	0.365	0.365	0.503	0.503	0.336	0.336
Linear	0.012 _{.001}	0.012 _{.001}	0.704 _{.036}	0.704 _{.036}	0.138 _{.009}	0.138 _{.009}	0.327 _{.032}	0.327 _{.032}	0.874 _{.151}	0.874 _{.151}
GRU	0.013 _{.002}	0.013 _{.002}	0.594 _{.144}	0.594 _{.144}	0.134 _{.009}	0.134 _{.009}	0.193 _{.002}	0.193 _{.002}	0.394 _{.013}	0.394 _{.013}
Transformer	0.016 _{.001}	0.016 _{.001}	0.538 _{.066}	0.538 _{.066}	0.115 _{.005}	0.115 _{.005}	0.204 _{.006}	0.204 _{.006}	0.408 _{.011}	0.408 _{.011}
N-HiTS	0.012 _{.000}	0.012 _{.000}	0.572 _{.020}	0.572 _{.020}	0.074 _{.003}	0.074 _{.003}	0.193 _{.002}	0.193 _{.002}	0.332 _{.011}	0.332 _{.011}
NLinear	0.010 _{.000}	0.010_{.000}	0.560 _{.002}	0.560 _{.002}	0.083 _{.002}	0.083 _{.002}	0.233 _{.001}	0.233 _{.001}	0.321 _{.001}	0.321 _{.001}
DLinear	0.012 _{.001}	0.012 _{.001}	0.547 _{.009}	0.547 _{.009}	0.076 _{.003}	0.076 _{.003}	0.250 _{.002}	0.250 _{.002}	0.412 _{.001}	0.412 _{.001}
PatchTST	0.010 _{.000}	0.010_{.000}	0.496 _{.002}	0.496 _{.002}	0.067 _{.001}	0.067 _{.001}	0.202 _{.001}	0.202 _{.001}	0.257 _{.001}	0.257_{.001}
TimesNet	0.011 _{.001}	0.011 _{.001}	0.507 _{.019}	0.507 _{.019}	0.071 _{.002}	0.071 _{.002}	0.205 _{.002}	0.205 _{.002}	0.304 _{.002}	0.304 _{.002}
iTransformer	0.010 _{.000}	0.010 _{.000}	0.496 _{.000}	0.496 _{.000}	0.074 _{.000}	0.074 _{.000}	0.158 _{.000}	0.158_{.000}	0.262 _{.000}	0.262 _{.000}
GRU NVP	0.016 _{.003}	0.020 _{.003}	0.396 _{.021}	0.507 _{.022}	0.055 _{.002}	0.073 _{.003}	0.161 _{.006}	0.203 _{.009}	0.282 _{.003}	0.330 _{.003}
GRU MAF	0.015 _{.001}	0.020 _{.001}	0.386 _{.026}	0.492 _{.027}	0.051 _{.001}	0.067 _{.001}	0.131 _{.006}	0.165 _{.009}	0.281 _{.004}	0.337 _{.005}
Trans MAF	0.011 _{.001}	0.014 _{.001}	0.400 _{.022}	0.503 _{.022}	0.054 _{.004}	0.071 _{.005}	0.129_{.004}	0.165 _{.006}	0.289 _{.008}	0.344 _{.008}
TimeGrad	0.011 _{.001}	0.014 _{.002}	0.359_{.011}	0.445_{.023}	0.052 _{.001}	0.067 _{.001}	0.164 _{.091}	0.201 _{.115}	0.272 _{.008}	0.327 _{.011}
CSDI	0.008_{.000}	0.011 _{.000}	0.366 _{.005}	0.484 _{.008}	0.050_{.001}	0.065_{.001}	0.146 _{.012}	0.176 _{.013}	0.219_{.006}	0.259 _{.009}

in replicating the performance levels reported in the original paper across various datasets. It appears that our experience aligns with the observations made in the Chronos paper [2], specifically as depicted in Figure 5, where ForecastPFN’s performance was not as robust as initially anticipated.

For TTM [20], we have adapted its official code¹² into our ProBTs framework and have loaded the `ibm-granite/granite-timeseries-ttm-v1` checkpoint. However, this method does not support arbitrary lengths for forecasting. The publicly available model currently supports a forecast length of 96 only, and thus, we have not conducted evaluations at other forecast lengths.

It is worth noting that, due to time constraints, we did not adjust the context window for all models. Instead, we chose 96 as a balanced and fair window size, which might result in suboptimal performance for some models.

B.4 Data and Code Availability

We release the ProBTs toolkit, documentation, and running scripts at <https://github.com/microsoft/ProBTs> under the MIT license. The repository includes parameter configurations for benchmarking experiments, ensuring reproducibility of all results presented in the paper. Most datasets used in this paper licensed under Creative Commons Attribution 4.0 International (CC BY 4.0), accessible via instructions in the repository.

C Overall Comparison Results

C.1 An Overall Comparison of Traditional Time-series Models on Short-term Forecasting

Table 9 presents a comprehensive comparison of various time-series models in ProBTs on short-term forecasting scenarios. The results, reported as mean \pm standard deviation, are derived from five independent runs with different seeds for each scenario.

C.2 An Overall Comparison of Traditional Time-series Models on Long-term Forecasting

Table 10 presents a comprehensive comparison of various time-series models in ProBTs on long-term forecasting scenarios. The results, reported as mean \pm standard deviation, are derived from five independent runs with different seeds for each scenario. The input sequence length is set to 36 for the ILI dataset and 96 for the others. Due to the excessive time and memory consumption of CSDI in producing long-term forecasts, its results are unavailable in some datasets.

¹²https://github.com/ibm-granite/granite-tsfm/blob/main/notebooks/hfdemo/ttm_getting_started.ipynb

Table 10: Results (mean_{std}) on long-term forecasting scenarios, each containing five independent runs with different seeds. The input sequence length is set to 36 for the ILL dataset and 96 for the others. Due to the excessive time and memory consumption of CSDI in producing long-term forecasts, its results are unavailable in some datasets.

Dataset	Pred len	iTransformer		PatchTST		DLinear		Autoformer		CSDI		TimeGrad		GRU NVP	
		CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE
ETTm1-L	96	0.271.000	0.271.000	0.272.001	0.272.001	0.282.002	0.282.002	0.388.001	0.388.001	0.236.006	0.308.005	0.522.105	0.645.129	0.383.053	0.488.058
	192	0.301.000	0.301.000	0.295.001	0.295.001	0.309.004	0.309.004	0.442.001	0.442.001	0.291.025	0.377.026	0.603.092	0.748.084	0.396.030	0.514.042
	336	0.333.000	0.333.000	0.323.001	0.323.001	0.338.008	0.338.008	0.429.000	0.429.000	0.322.083	0.419.042	0.601.028	0.759.015	0.486.032	0.630.029
	720	0.376.000	0.376.000	0.353.001	0.353.001	0.387.006	0.387.006	0.440.000	0.440.000	0.448.038	0.578.051	0.621.037	0.793.034	0.546.036	0.707.050
ETTm2-L	96	0.137.000	0.137.000	0.132.001	0.132.001	0.138.000	0.138.000	0.158.000	0.158.000	0.115.009	0.146.012	0.427.042	0.525.047	0.319.044	0.413.059
	192	0.161.000	0.161.000	0.157.001	0.157.001	0.163.003	0.163.003	0.175.000	0.175.000	0.147.008	0.189.012	0.424.061	0.530.060	0.326.025	0.427.033
	336	0.180.000	0.180.000	0.176.000	0.176.000	0.188.001	0.188.001	0.191.000	0.191.000	0.190.018	0.248.024	0.469.049	0.566.047	0.449.145	0.580.169
	720	0.211.000	0.211.000	0.205.001	0.205.001	0.219.003	0.219.003	0.217.000	0.217.000	0.239.035	0.306.040	0.470.054	0.561.044	0.561.273	0.749.385
ETTh1-L	96	0.321.000	0.321.000	0.328.003	0.328.003	0.352.011	0.352.011	0.367.000	0.367.000	0.437.018	0.557.022	0.455.046	0.585.058	0.379.030	0.481.037
	192	0.359.000	0.359.000	0.359.002	0.359.002	0.393.001	0.393.001	0.392.000	0.392.000	0.496.051	0.625.065	0.516.038	0.680.058	0.425.019	0.531.018
	336	0.388.000	0.388.000	0.384.002	0.384.002	0.419.007	0.419.007	0.398.000	0.398.000	0.454.025	0.574.026	0.512.026	0.666.047	0.458.054	0.580.064
	720	0.408.000	0.408.000	0.397.002	0.397.002	0.502.029	0.502.029	0.433.000	0.433.000	0.528.012	0.657.014	0.523.027	0.672.015	0.502.039	0.643.046
ETTh2-L	96	0.177.000	0.177.000	0.177.000	0.177.000	0.211.027	0.211.027	0.203.000	0.203.000	0.164.013	0.214.018	0.358.026	0.448.031	0.432.141	0.548.158
	192	0.203.000	0.203.000	0.201.001	0.201.001	0.238.028	0.238.028	0.226.000	0.226.000	0.266.018	0.294.027	0.457.081	0.575.089	0.625.170	0.766.223
	336	0.243.000	0.243.000	0.240.001	0.240.001	0.284.008	0.284.008	0.264.000	0.264.000	0.274.022	0.353.028	0.481.076	0.606.095	0.793.319	0.942.408
	720	0.264.000	0.264.000	0.252.000	0.252.000	0.307.000	0.307.000	0.287.000	0.287.000	0.302.040	0.382.030	0.445.018	0.550.018	0.539.090	0.688.161
Electricity-L	96	0.098.000	0.098.000	0.086.001	0.086.001	0.090.001	0.090.001	0.140.000	0.140.000	0.153.137	0.203.189	0.096.002	0.119.003	0.094.003	0.118.003
	192	0.106.000	0.106.000	0.092.001	0.092.001	0.095.001	0.095.001	0.136.000	0.136.000	0.200.094	0.264.129	0.100.004	0.124.005	0.097.002	0.121.003
	336	0.115.000	0.115.000	0.100.000	0.100.000	0.104.000	0.104.000	0.147.000	0.147.000	—	—	0.102.007	0.126.008	0.099.001	0.123.001
	720	0.133.000	0.133.000	0.116.000	0.116.000	0.122.001	0.122.001	0.159.000	0.159.000	—	—	0.108.003	0.134.004	0.114.013	0.144.017
Traffic-L	96	0.246.000	0.246.000	0.248.001	0.248.001	0.356.009	0.356.009	0.293.000	0.293.000	—	—	0.202.004	0.234.006	0.187.002	0.231.003
	192	0.259.000	0.259.000	0.245.001	0.245.001	0.346.009	0.346.009	0.318.000	0.318.000	—	—	0.208.003	0.239.004	0.192.001	0.236.002
	336	0.283.000	0.283.000	0.257.002	0.257.002	0.350.008	0.350.008	0.332.000	0.332.000	—	—	0.213.003	0.246.003	0.201.004	0.248.006
	720	0.275.000	0.275.000	0.266.001	0.266.001	0.365.009	0.365.009	0.341.003	0.341.003	—	—	0.220.002	0.263.001	0.211.004	0.264.006
Weather-L	96	0.089.000	0.089.000	0.087.002	0.087.002	0.112.001	0.112.001	0.239.004	0.239.004	0.068.008	0.087.012	0.130.017	0.164.023	0.116.013	0.145.017
	192	0.093.000	0.093.000	0.090.001	0.090.001	0.122.001	0.122.001	0.213.000	0.213.000	0.068.006	0.086.007	0.127.019	0.158.024	0.122.021	0.147.025
	336	0.096.000	0.096.000	0.092.002	0.092.002	0.130.002	0.130.002	0.176.000	0.176.000	0.083.002	0.098.002	0.130.006	0.162.006	0.128.011	0.160.012
	720	0.099.000	0.099.000	0.094.001	0.094.001	0.144.001	0.144.001	0.170.001	0.170.001	0.087.003	0.102.005	0.113.011	0.136.020	0.110.004	0.135.008
Exchange-L	96	0.025.000	0.025.000	0.023.000	0.023.000	0.024.000	0.024.000	0.032.000	0.032.000	0.028.003	0.036.005	0.068.003	0.079.002	0.071.006	0.091.009
	192	0.036.000	0.036.000	0.034.000	0.034.000	0.035.000	0.035.000	0.041.000	0.041.000	0.045.003	0.058.005	0.087.013	0.100.019	0.068.004	0.087.005
	336	0.048.000	0.048.000	0.048.000	0.048.000	0.048.001	0.048.001	0.056.000	0.056.000	0.060.004	0.076.006	0.074.009	0.086.008	0.072.002	0.091.002
	720	0.076.000	0.076.000	0.072.000	0.072.000	0.075.002	0.075.002	0.112.002	0.112.002	0.143.020	0.173.020	0.099.015	0.113.016	0.079.009	0.103.009
ILL	24	0.094.000	0.094.000	0.169.005	0.169.005	0.213.038	0.213.038	0.122.000	0.122.000	0.250.013	0.263.012	0.275.047	0.296.044	0.257.003	0.283.001
	36	0.102.000	0.102.000	0.156.005	0.156.005	0.230.015	0.230.015	0.111.000	0.111.000	0.285.010	0.298.011	0.272.057	0.298.048	0.281.004	0.307.007
	48	0.103.000	0.103.000	0.156.008	0.156.008	0.221.009	0.221.009	0.134.000	0.134.000	0.285.036	0.301.034	0.295.033	0.320.025	0.288.008	0.314.009
	60	0.128.000	0.128.000	0.147.003	0.147.003	0.230.013	0.230.013	0.144.000	0.144.000	0.283.012	0.299.013	0.295.083	0.325.068	0.307.005	0.333.005

C.3 An Overall Comparison of Time-series Foundation Models on Diverse Prediction Horizons

Based on the results in Table 4, we selected several datasets that most foundation models have not been pre-trained on for zero-shot evaluation. We compared these foundation models with fully-supervised traditional non-universal time-series models. The results, presented in Table 11, are reported as normalized MAE (NMAE). To comprehensively evaluate short-term and long-term forecasting performance, we selected prediction horizons of {24, 48, 96, 192, 336, 720}. Due to time constraints, the context window for most models was set to 96 unless otherwise specified. For Lag-Llama, we chose a context window of 512 based on explicit recommendations from the original paper and model specifications.

Additionally, on the MOIRAI model, we explored the impact of longer context windows. In some datasets (e.g., ETTh2-L, Weather-L), there were significant improvements, which we retained for reference. Other models also utilized longer context windows, but without consistent performance gains. It is worth noting that the longer context window of MOIRAI had an adverse effect on the Electricity dataset. To achieve optimal performance on unseen data, these models may require a hyperparameter search using validation data, which we leave for future work.

Table 11: Results of time-series foundation models on diverse prediction horizons. Mean NMAE value of five independent runs with different seeds is reported. The input sequence length is set to 96 if not specified. For every model, we exclude the evaluation results on its pre-trained datasets

Dataset	Pred	MOIRAI-5000	MOIRAI	Lag-Llama-512	Chronos	TimesFM	Timer	UniTS	ForecastPFN	CSDI	DLinear	PatchTST	iTransformer
ETTh1-L	24	0.144	0.128	0.231	0.113	0.133	0.182	0.353	0.878	<u>0.115</u>	0.121	0.199	0.116
	48	0.269	0.343	0.244	0.307	0.269	0.339	0.432	0.997	0.266	<u>0.238</u>	0.199	0.243
	96	0.296	0.449	0.399	0.393	0.324	0.384	0.457	0.961	0.303	0.283	<u>0.271</u>	0.269
	192	0.311	0.479	0.416	0.422	0.378	0.423	0.466	1.031	0.389	0.309	0.295	<u>0.304</u>
	336	0.321	0.471	0.429	0.439	0.435	0.446	0.475	1.012	0.449	0.337	<u>0.324</u>	0.337
720	0.351	0.515	0.472	0.467	0.511	0.480	0.499	1.053	0.530	0.385	<u>0.353</u>	0.380	
ETTh2-L	24	<u>0.100</u>	0.113	0.128	0.110	0.117	0.130	0.173	1.506	0.096	0.103	0.162	0.106
	48	0.116	0.132	0.162	0.132	0.132	0.141	0.165	1.471	0.121	0.121	0.162	<u>0.118</u>
	96	0.141	0.168	0.188	0.158	0.156	0.153	0.169	1.386	<u>0.133</u>	0.138	0.133	0.141
	192	0.159	0.192	0.205	0.183	0.186	0.175	0.184	1.397	0.201	0.167	0.158	<u>0.158</u>
	336	0.173	0.209	0.226	0.206	0.209	0.193	0.199	1.422	0.226	0.188	<u>0.176</u>	0.184
720	0.200	0.242	0.249	0.240	0.246	0.220	0.223	1.485	0.264	0.222	<u>0.206</u>	0.213	
ETTh1-L	24	0.304	0.297	0.313	0.265	0.277	0.315	0.453	1.141	0.292	0.277	0.356	<u>0.275</u>
	48	0.312	0.326	0.341	0.294	0.307	0.339	0.461	1.161	0.392	0.311	0.356	<u>0.304</u>
	96	0.324	0.354	0.353	<u>0.321</u>	0.332	0.361	0.469	1.157	0.534	0.340	0.326	0.319
	192	0.350	0.393	0.376	0.375	0.383	0.399	0.482	1.226	0.698	0.394	<u>0.357</u>	0.364
	336	0.366	0.417	0.393	0.410	0.411	0.433	0.500	1.183	0.603	0.423	<u>0.383</u>	0.389
720	0.380	0.448	0.424	0.432	0.418	0.460	0.499	1.037	0.664	0.501	<u>0.399</u>	0.402	
ETTh2-L	24	0.126	0.142	0.160	<u>0.127</u>	0.134	0.146	0.197	1.650	0.133	0.146	0.213	0.135
	48	0.147	0.170	0.195	<u>0.161</u>	0.166	0.166	0.199	1.667	0.186	0.176	0.213	0.161
	96	0.162	0.200	0.203	0.194	0.190	0.179	0.204	1.685	0.204	0.209	<u>0.176</u>	0.177
	192	0.189	0.244	0.220	0.225	0.220	0.205	0.223	1.768	0.272	0.206	<u>0.200</u>	0.205
	336	0.224	0.269	0.245	0.252	0.266	0.247	0.262	1.719	0.321	0.293	<u>0.240</u>	0.244
720	0.243	0.293	<u>0.248</u>	0.290	0.277	0.254	0.264	1.542	0.417	0.307	0.251	0.263	
Weather-L	24	0.043	0.060	0.062	0.060	—	0.063	—	1.916	<u>0.050</u>	0.095	0.083	0.074
	48	0.066	0.110	<u>0.086</u>	0.128	—	0.098	—	1.924	0.093	0.125	0.156	0.090
	96	0.074	0.134	0.096	0.163	—	0.109	—	1.924	0.092	0.113	<u>0.086</u>	0.089
	192	0.074	0.126	0.103	0.161	—	0.116	—	1.926	<u>0.079</u>	0.121	0.089	0.099
	336	0.075	0.129	0.109	0.177	—	0.121	—	1.930	0.100	0.131	<u>0.092</u>	0.100
720	0.076	0.151	0.120	0.208	—	0.124	—	1.970	0.108	0.145	<u>0.094</u>	0.111	
Electricity-L	24	0.227	<u>0.091</u>	—	—	—	0.114	—	—	—	0.093	—	0.085
	48	0.210	0.100	—	—	—	0.128	—	—	—	<u>0.097</u>	—	0.089
	96	0.194	0.101	—	—	—	0.130	—	—	0.153	0.090	0.086	<u>0.098</u>
	192	0.200	0.107	—	—	—	0.147	—	—	0.200	<u>0.095</u>	0.092	0.106
	336	0.202	0.119	—	—	—	0.168	—	—	—	<u>0.104</u>	0.100	0.115
720	0.217	0.190	—	—	—	0.205	—	—	—	<u>0.121</u>	0.116	0.133	

C.4 An Overall Comparison of Time-series Foundation Models on Short-term Probabilistic Forecasting

Table 12 presents a comparison of two time-series probabilistic foundation models in short-term forecasting scenarios. We excluded the results of datasets that has been used in pre-training based on table 4. We also explored both short and long context windows for MOIRAI.

Table 12: Results of probabilistic foundation models on short-term distributional forecasting. For every model, we exclude the evaluation results on its pre-trained datasets.

Model	Exchange-S		Solar-S		Electricity-S		Traffic-S	
	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE
CSDI	0.008 _{.000}	0.011 _{.000}	0.366 _{.005}	0.484 _{.008}	0.050 _{.001}	0.065 _{.001}	0.146 _{.012}	0.176 _{.013}
Chronos	0.007	<u>0.010</u>	–	–	–	–	<u>0.178</u>	<u>0.211</u>
MOIRAI-96	0.007	0.010	<u>0.502</u>	<u>0.681</u>	0.069	0.084	–	–
MOIRAI-5000	<u>0.007</u>	0.010	0.521	0.702	<u>0.059</u>	<u>0.079</u>	–	–

D Additional Results and Experiments

D.1 The Impact of Normalization

Normalization is a crucial aspect of time-series models, with different research branches adopting distinct strategies that can affect model performance across various data scenarios. Motivated by this, we provide a detailed analysis of different normalization methods in this section. Unless stated otherwise, our benchmarking follows the default normalization methods used by each model.

D.1.1 Different Normalization Choice Between Distinct Research Branches

In time series forecasting, normalization typically occurs in two stages. First, during preprocessing, *dataset-level normalization* is applied, where global statistics (e.g., mean and standard deviation) from the training set are used to normalize all time-series values. Then, a local normalization module might be used to perform *instance-level normalization* when feeding a batch of time-series segments into the model.

Different research branches prefer distinct instance-level normalization strategies. Long-term point forecasting models [49, 40] typically adopt the RevIN [32]. Given a batch of time-series segments within a lookback window, RevIN applies a per-series z-score normalization, augmented with learnable affine parameters. Its main advantage is its effectiveness in addressing distribution shifts, particularly in long-term forecasting.

In contrast, most short-term probabilistic forecasting models [57, 58] employ an ad-hoc but still effective normalization strategy. For example, given a batch of time-series segments $X \in \mathbb{R}^{K \times L}$ (where K is the number of variables and L is the length of the lookback window), a per-series scaling is applied as $X_i^{\text{norm}} = \frac{X_i}{\sum_{t=1}^L |X_{i,t}|/L}$, $i = 1, \dots, K$ to stabilize value ranges. For simplicity, we refer to this type of normalization as *Mean Scaling*.

We summarize the instance-level normalization choices originally used by each model in the Table 13.

Table 13: Original instance-level normalization choices of each model.

Normalization Choice	Model
ReVIN	iTransformer, PatchTST
Mean Scaling	TimeGrad, GRU NVP
w/o Norm	GRU, CSDI, DLinear

Existing probabilistic models rarely use RevIN and are seldom combined with AR-based models that employ RevIN-style normalization. Similarly, the mean scaling strategy, commonly used in probabilistic forecasting models, is rarely applied to models designed for long-term forecasting. To better understand the effects of different normalization strategies, we selected representative models from both categories and combined them with three normalization methods: *RevIN*, *Scaling* (i.e., mean scaling), and *w/o Norm* (no instance-level normalization, using time-series values as provided by the dataset-level preprocessing). The results of these experiments are presented in Table 14, 15, 16, and 17.

Table 14: The impact of different normalization methods in short-term forecasting scenarios (CRPS).

Model	PatchTST			CSDI			TimeGrad			GRU NVP		
	ReVIN	Scaling	w/o Norm	ReVIN	Scaling	w/o Norm	ReVIN	Scaling	w/o Norm	ReVIN	Scaling	w/o Norm
Electricity-S	0.0659	0.0645	0.0660	0.0524	-	0.0502	0.0673	0.0563	0.9681	0.0659	0.0706	0.1607
Exchange Rate-S	0.0102	0.0108	0.0111	0.0070	0.0083	0.0110	0.0100	0.0093	0.0170	0.0090	0.0147	0.0133
Solar-S	0.6275	0.7105	0.7169	0.4903	0.4347	0.4603	0.4945	0.5455	0.8356	0.9293	0.5926	0.4393
Traffic-S	0.2001	0.2036	0.2168	0.1505	0.1552	0.1389	0.1806	0.1280	0.1400	0.1827	0.1770	0.2277
Wikipedia-S	0.2529	0.3245	0.3695	0.2164	0.2060	0.2276	0.2757	0.2773	0.9969	0.3317	0.3187	0.4561

Table 15: The impact of different normalization methods in short-term forecasting scenarios (NMAE).

Model	PatchTST			CSDI			TimeGrad			GRU NVP		
	ReVIN	Scaling	w/o Norm	ReVIN	Scaling	w/o Norm	ReVIN	Scaling	w/o Norm	ReVIN	Scaling	w/o Norm
Electricity-S	0.0659	0.0645	0.0660	0.0666	-	0.0648	0.0852	0.0710	0.9742	0.0861	0.0929	0.2246
Exchange Rate-S	0.0102	0.0108	0.0111	0.0096	0.0111	0.0151	0.0115	0.0118	0.0220	0.0114	0.0189	0.0170
Solar-S	0.6275	0.7105	0.7169	0.5988	0.5616	0.5680	0.6041	0.7011	0.9162	1.1931	0.7424	0.5893
Traffic-S	0.2001	0.2036	0.2168	0.1752	0.1877	0.1669	0.2167	0.1516	0.1693	0.2257	0.2216	0.2837
Wikipedia-S	0.2529	0.3245	0.3695	0.2585	0.2437	0.2698	0.3278	0.3257	0.9998	0.4041	0.3559	0.5131

D.1.2 Analysis of Instance-level Normalization

ReVIN Significantly Improves Most Models in Long-term Forecasting Scenarios, with Some Exceptions. ReVIN’s ability to mitigate the effects of data distribution shifts leads to significant performance improvements in most models for long-term forecasting, as shown in Table 16. This benefit extends beyond models like PatchTST and iTransformer, which originally employed ReVIN, to others such as DLinear that do not inherently use this approach. Notably, ReVIN has greatly enhanced AR-based models in long-term scenarios. For instance, on the ETT datasets, GRU NVP (w/ ReVIN) outperforms even PatchTST (w/ ReVIN), suggesting that normalizing trend effects can help reduce error accumulation in AR-based models.

However, ReVIN can have a negative impact in certain cases. On the Traffic dataset, GRU (w/ ReVIN) and GRU NVP (w/ ReVIN) perform worse than without normalization, as shown in Table 16. Interestingly, this aligns with our analysis of data characteristics: the Traffic dataset displays strong seasonality but less trending. We speculate that ReVIN’s effectiveness in other datasets stems from its ability to normalize trend-related distribution shifts, which is less relevant for the Traffic dataset. Additionally, ReVIN appears less suited for NAR probabilistic models. For instance, CSDI (w/ ReVIN) performs worse than CSDI (w/ Scaling) on the Weather, Electricity, Exchange, and ILI dataset. Further research is needed to develop more effective normalization strategies for NAR probabilistic models.

No Dominating Normalization Strategies in Short-term Forecasting. As shown in Table 14, ReVIN does not consistently provide robust or significant improvements for models such as CSDI, TimeGrad, and GRU NVP in short-term forecasting. The Mean Scaling strategy, though empirical, proves to be the most reliable choice for these probabilistic models, likely explaining its widespread use. In some cases, instance-level normalization can be omitted, but this approach can lead to serious issues, as seen with TimeGrad (w/o normalization) on the Wikipedia and Solar datasets, and GRU NVP (w/o normalization) on the Electricity dataset. Developing effective instance-level normalization methods for complex data distributions in short-term forecasting remains an important yet often overlooked research direction.

D.2 The Impact of Data Scale

To further explore critical characteristics of time-series forecasting, we have examined the correlation between model performance gains, relative to the baseline model (GRU), and dataset dimensions, length, and volume (see Table 18). However, our analysis does not identify a significant correlation between these factors and model performance.

D.3 Statistical and Gradient Boosting Decision Tree Baselines

To enhance the empirical robustness of our study, we integrate classical statistical models, including ARIMA [44] and ETS [30], along with the Gradient Boosting Decision Tree (GBDT) model, XGBoost, into the ProBTS framework. The results in Table 19 clearly demonstrate the superior performance of deep learning methods over simple statistical baselines, emphasizing the importance

Table 16: The impact of different normalization methods in long-term forecasting scenarios (CRPS).

Dataset	Pred len	iTransformer		DLinear		PatchTST		CSDI		TimeGrad		GRU NVP		GRU		
		w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	Scaling
ETTh1-L	96	0.3514	0.3148	0.3613	0.3210	0.3447	0.3212	0.3437	0.3212	0.3630	0.2764	0.5434	0.2958	0.6963	0.4444	0.5946
	192	0.4427	0.3479	0.3843	0.3574	0.3754	0.3562	0.3943	0.3553	0.4352	0.3119	0.5698	0.3119	0.6410	0.5056	0.3076
	336	0.4192	0.3654	0.4001	0.3750	0.3751	0.3806	0.4099	0.3737	0.4213	0.4814	0.4814	0.3614	0.4192	0.6192	0.3081
720	0.5172	0.3902	0.5576	0.4446	0.3878	0.4738	0.4648	0.3909	0.4881	0.5609	0.3960	0.4283	0.5837	0.8670	0.5862	0.3641
ETTh2-L	96	0.2790	0.1796	0.2099	0.2034	0.2179	0.2031	0.1892	0.1763	0.1900	0.1681	0.1446	0.1590	0.6859	0.1738	0.4399
	192	0.3322	0.2044	0.2346	0.2700	0.2449	0.2572	0.2189	0.2012	0.2167	0.2078	0.1734	0.1983	0.6811	0.1905	0.5601
	336	0.4475	0.2209	0.2726	0.2617	0.2339	0.2604	0.2446	0.2244	0.2446	0.2728	0.2094	0.2124	0.6393	0.2213	0.4178
720	0.4315	0.2329	0.3148	0.3009	0.2768	0.2942	0.2819	0.2501	0.2829	0.3062	0.2054	0.2512	0.4990	0.2198	0.3736	
ETTh1-L	96	0.3222	0.2845	0.2890	0.2662	0.2674	0.2699	0.2989	0.2649	0.2858	0.2569	0.2255	0.2597	0.6828	0.2886	0.4473
	192	0.3446	0.3040	0.3122	0.2914	0.2911	0.3007	0.3349	0.2925	0.3077	0.3280	0.3313	0.3495	0.7930	0.2975	0.5341
	336	0.3738	0.3276	0.3303	0.3180	0.3118	0.3215	0.3618	0.3107	0.3338	0.3663	0.3243	0.4120	0.6675	0.3201	0.5657
720	0.4248	0.3660	0.3831	0.3508	0.3468	0.3683	0.3790	0.3468	0.3780	0.3933	0.3789	0.4367	0.8248	0.3357	0.5826	
ETTh2-L	96	0.1505	0.1381	0.1528	0.1438	0.1390	0.1436	0.1522	0.1335	0.1587	0.1314	0.1134	0.1340	0.4493	0.1206	0.2376
	192	0.1833	0.1674	0.1750	0.1645	0.1625	0.1642	0.1884	0.1601	0.1940	0.1391	0.1386	0.1583	0.3959	0.1399	0.3469
	336	0.2628	0.1929	0.2122	0.1933	0.1828	0.1939	0.1994	0.1798	0.2249	0.2139	0.1592	0.1782	0.4661	0.1561	0.3187
720	0.3098	0.2142	0.2278	0.2184	0.2094	0.2183	0.2940	0.2063	0.2905	0.2218	0.1892	0.2002	0.4260	0.1822	0.2896	
Electricity-L	96	0.0887	0.0827	0.0845	0.0872	0.0862	0.0871	0.0848	0.0857	0.0867	0.0735	0.0761	0.0735	0.0961	0.0771	0.0904
	192	0.0911	0.0892	0.0906	0.0977	0.0931	0.0934	0.0910	0.0912	0.0918	0.2475	0.2749	0.2584	0.0980	0.0806	0.0936
	336	0.1026	0.0984	0.1024	0.1020	0.1018	0.1019	0.1006	0.1001	0.1005	0.3136	0.2153	0.2835	0.1125	0.0905	0.1002
720	0.1148	0.1110	0.1123	0.1193	0.1170	0.1192	0.1178	0.1160	0.1187	0.2667	25.7527	0.2945	0.1096	0.1164	0.0957	
Exchange-L	96	0.0461	0.0251	0.0318	0.0244	0.0240	0.0243	0.0299	0.0235	0.0254	0.0343	0.0216	0.0210	0.0637	0.0279	0.0478
	192	0.0694	0.0345	0.0425	0.0343	0.0341	0.0338	0.0404	0.0336	0.0365	0.0418	0.0383	0.0388	0.0638	0.0364	0.0743
	336	0.0843	0.0460	0.0572	0.0450	0.0472	0.0451	0.0597	0.0462	0.0532	0.0510	0.0508	0.0485	0.1030	0.0510	0.1142
720	0.1317	0.0769	0.1145	0.0704	0.0777	0.0830	0.0824	0.0777	0.0812	0.0964	0.0881	0.1042	0.1177	0.0800	0.1428	
ILL	24	0.3792	0.1093	0.0971	0.1920	0.1183	0.2138	0.2665	0.0794	0.1478	0.2541	0.1344	0.1182	0.2776	0.0815	0.0920
	36	0.3453	0.1422	0.1632	0.1864	0.1551	0.2024	0.2955	0.1239	0.1422	0.2807	0.1214	0.1526	0.3164	0.0962	0.1524
	48	0.3349	0.1617	0.1653	0.2021	0.1732	0.2059	0.2957	0.1312	0.1704	0.3383	0.1026	0.1453	0.3229	0.1124	0.2275
60	0.3365	0.1631	0.2148	0.2389	0.1678	0.2321	0.3034	0.1493	0.1964	0.2651	0.1333	0.1496	0.3162	0.1250	0.1484	
Traffic-L	96	0.3495	0.2379	0.2377	0.3528	0.3581	0.3546	0.2526	0.2553	0.2546	-	-	0.2094	0.2309	0.2033	
	192	0.3821	0.2434	0.2394	0.3345	0.3307	0.3366	0.2483	0.2479	0.2447	-	-	0.2079	0.2110	0.2091	
	336	0.3749	0.2503	0.2467	0.3319	0.3294	0.3316	0.2489	0.2492	0.2470	-	-	0.2271	0.2306	0.2091	
720	0.3940	0.2597	0.2580	0.3672	0.3337	0.3671	0.2577	0.2572	0.2568	-	-	0.2302	0.2659	0.2184		
Weather-L	96	0.1083	0.0930	0.1084	0.1107	0.0959	0.1092	0.1051	0.0837	0.1082	0.0591	0.1202	0.0715	0.2939	0.0737	0.2308
	192	0.1160	0.0940	0.1093	0.1227	0.0984	0.1209	0.1136	0.0858	0.1171	0.0790	0.1599	0.0571	0.3138	0.0788	0.2495
	336	0.1131	0.0932	0.1208	0.1317	0.1009	0.1296	0.1118	0.0903	0.1167	0.0977	0.4296	0.0856	0.2249	0.0860	0.2646
720	0.1277	0.1009	0.1212	0.1442	0.1050	0.1425	0.1117	0.0953	0.1234	0.0975	0.1785	0.1419	0.2069	0.0941	0.1973	
													0.1979	0.0925	0.1146	
													0.1005	0.0837	0.1546	
													0.1177	0.0866	0.1342	
													0.1145	0.0868	0.1715	
													0.1979	0.0925	0.1146	
													0.2823	0.0720	0.1312	
													0.2806	0.0983	0.1403	
													0.2908	0.1193	0.1499	
													0.3254	0.1204	0.1817	
													0.2094	0.1952	0.2021	
													0.2005	0.2110	0.2034	
													0.2271	0.1983	0.2220	
													0.2302	0.2198	0.2524	
													0.2939	0.0737	0.2308	
													0.3138	0.0788	0.2495	
													0.1145	0.0868	0.1715	
													0.1979	0.0925	0.1146	

Table 17: The impact of different normalization methods in long-term forecasting scenarios (NMAE).

Dataset	Pred len	iTransformer		DLinear		PatchTST		CSDI		TimeGrad		GRU/NVP		GRU		
		w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	w/o Norm	ReVIN	Scaling
ETTh1-L	96	0.3514	0.3148	0.3613	0.3210	0.3240	0.3212	0.3437	0.3643	0.4605	0.3783	0.8983	0.3569	0.5452	0.4457	0.9734
	192	0.4427	0.3479	0.3843	0.3574	0.3754	0.3562	0.3943	0.4445	0.5683	0.3926	0.8034	0.3946	0.6226	0.4763	1.0658
	336	0.4192	0.3654	0.4001	0.3750	0.3751	0.3737	0.4213	0.4524	0.5321	0.4942	0.8381	0.3983	0.5773	0.4998	1.1416
720	0.5172	0.3902	0.5576	0.3878	0.4446	0.3909	0.4881	0.4648	0.5118	0.5351	0.4332	1.0233	0.4409	0.5901	1.2715	
ETTh2-L	96	0.2790	0.1796	0.2099	0.2034	0.2179	0.1763	0.1900	0.1879	0.2030	0.2080	0.4498	0.2224	0.4846	0.2913	0.5277
	192	0.3322	0.2044	0.2346	0.2700	0.2449	0.2012	0.2167	0.2683	0.2550	0.2323	0.5783	0.2282	0.5304	0.3334	0.5329
	336	0.4475	0.2209	0.2726	0.2339	0.2604	0.2244	0.2446	0.2671	0.2725	0.2657	0.5096	0.2789	0.8343	0.3658	0.8879
720	0.4315	0.2329	0.3148	0.2768	0.2942	0.2501	0.2829	0.2688	0.3303	0.2589	0.4772	0.2886	1.0349	0.3594	0.6519	
ETTh1-L	96	0.3222	0.2845	0.2890	0.2662	0.2674	0.2649	0.2858	0.2879	0.3228	0.3647	0.5763	0.3848	0.4806	0.4411	0.9213
	192	0.3446	0.3040	0.3122	0.2914	0.2914	0.2925	0.3077	0.4235	0.4310	0.3680	0.7200	0.4158	0.6149	0.4824	0.9978
	336	0.3738	0.3276	0.3303	0.3180	0.3118	0.3107	0.3338	0.3978	0.5297	0.3904	0.7476	0.4171	0.5751	0.6186	1.0089
720	0.4248	0.3660	0.3831	0.3468	0.3508	0.3468	0.3780	0.4748	0.5148	0.4063	0.7444	0.4452	0.5375	0.7808	1.0310	
ETTh2-L	96	0.1505	0.1381	0.1528	0.1438	0.1438	0.1390	0.1587	0.1422	0.1713	0.1513	0.2911	0.1454	0.5324	0.1802	0.4339
	192	0.1833	0.1674	0.1750	0.1625	0.1645	0.1601	0.1940	0.1787	0.2011	0.1742	0.4467	0.1692	0.5476	0.2074	0.4159
	336	0.2628	0.1929	0.2122	0.1933	0.1939	0.1798	0.2249	0.2010	0.2262	0.1879	0.3843	0.1939	0.5435	0.2216	0.4764
720	0.3098	0.2142	0.2278	0.2184	0.2184	0.2094	0.2183	0.2879	0.2546	0.2340	0.3588	0.2200	0.4969	0.3022	0.5101	
Electricity-L	96	0.0887	0.0827	0.0845	0.0872	0.0872	0.0862	0.0871	0.0848	0.0857	0.0975	0.1148	0.1007	0.1141	0.1218	0.2261
	192	0.0911	0.0892	0.0906	0.0977	0.0934	0.0931	0.0918	0.3218	0.3788	0.1013	0.1201	0.1044	0.1170	0.1299	0.2381
	336	0.1026	0.0984	0.1024	0.1020	0.1019	0.1006	0.1005	0.4213	0.4689	0.1134	0.1279	0.1174	0.1225	0.1364	0.2807
720	0.1148	0.1110	0.1123	0.1193	0.1193	0.1170	0.1187	0.3726	0.3873	0.1424	0.1220	0.1339	0.1413	0.1337	0.3228	
Exchange-L	96	0.0461	0.0251	0.0318	0.0244	0.0244	0.0240	0.0243	0.0235	0.0254	0.0272	0.0617	0.0780	0.0737	0.0385	0.1539
	192	0.0694	0.0345	0.0425	0.0343	0.0341	0.0338	0.0404	0.0336	0.0365	0.0465	0.0909	0.1097	0.1097	0.0382	0.1704
	336	0.0843	0.0460	0.0572	0.0450	0.0472	0.0451	0.0597	0.0462	0.0532	0.0543	0.1392	0.0834	0.0834	0.0526	0.1252
720	0.1317	0.0769	0.1145	0.0704	0.0777	0.0830	0.0824	0.0777	0.0810	0.1112	0.1198	0.1328	0.0891	0.0722	0.1648	
ILL	24	0.3792	0.1093	0.0971	0.1920	0.1183	0.2138	0.2665	0.0794	0.1478	0.2718	0.1550	0.0955	0.3129	0.0866	0.1787
	36	0.3453	0.1422	0.1632	0.1864	0.1551	0.2024	0.2955	0.1239	0.1422	0.2966	0.1082	0.1082	0.3115	0.1162	0.1937
	48	0.3349	0.1617	0.1653	0.2021	0.1732	0.2059	0.2957	0.1312	0.1704	0.3601	0.1252	0.2650	0.3165	0.1364	0.1989
60	0.3365	0.1631	0.2148	0.2389	0.1678	0.2321	0.3034	0.1493	0.1964	0.2840	0.1607	0.1738	0.3300	0.1409	0.2397	
Traffic-L	96	0.3495	0.2379	0.2377	0.3528	0.3581	0.3546	0.2526	0.2553	0.2546	-	-	0.2436	0.2791	0.2414	
	192	0.3821	0.2434	0.2394	0.3345	0.3307	0.3366	0.2483	0.2479	0.2447	-	-	0.2432	0.2521	0.2461	
	336	0.3749	0.2503	0.2467	0.3319	0.3294	0.3316	0.2489	0.2492	0.2470	-	-	0.2672	0.2739	0.2471	
720	0.3940	0.2597	0.2580	0.3672	0.3337	0.3671	0.2577	0.2572	0.2568	-	-	0.2711	0.3125	0.2620		
Weather-L	96	0.1083	0.0940	0.1084	0.1107	0.0959	0.1092	0.1051	0.0837	0.1082	0.0746	0.1280	0.0710	0.1225	0.1045	0.2055
	192	0.1160	0.0930	0.1093	0.1227	0.0984	0.1209	0.1136	0.0858	0.1171	0.0891	0.2102	0.0857	0.3731	0.0938	0.3222
	336	0.1131	0.0932	0.1208	0.1317	0.1009	0.1296	0.1118	0.0903	0.1167	0.1143	0.5821	0.1023	0.1367	0.1030	0.2285
720	0.1277	0.1009	0.1212	0.1442	0.1050	0.1425	0.1117	0.0953	0.1234	0.1167	0.2316	0.1162	0.2540	0.1075	0.1458	

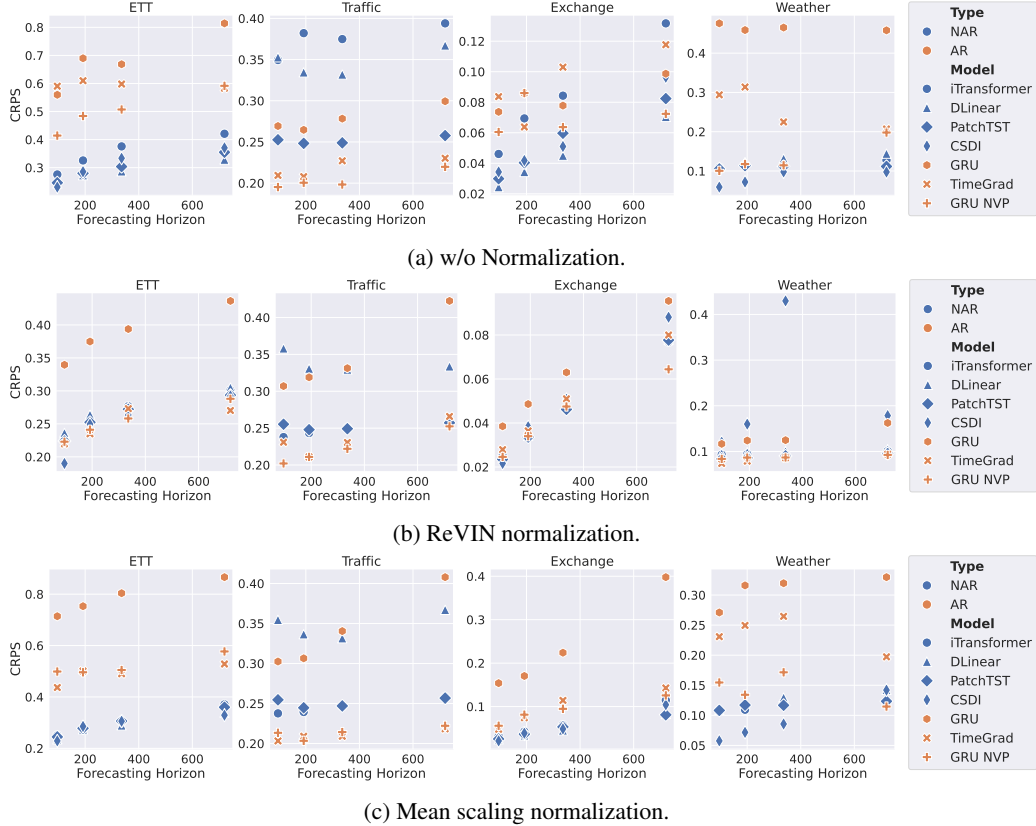


Figure 7: Impact of different instance-level normalization methods on model performance.

Table 18: The correlation coefficient between the data volume and the relative performance improvement compared to the baseline model (GRU).

Model	DLinear		PatchTST		GRU NVP		TimeGrad		CSDI	
	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE
# Var.	0.2422	0.2422	-0.2676	-0.2676	-0.1856	-0.2136	-0.1665	-0.1793	-0.2315	-0.2592
# Total timestep	-0.1422	-0.1422	0.3821	0.3821	0.3072	0.3329	0.2860	0.2971	0.3542	0.3826
# Var. \times Timestep	0.0162	0.0162	0.0166	0.0166	-0.0068	-0.0011	0.0082	0.0117	-0.0053	-0.0133

of capturing non-linear dependencies for accurate forecasts. Notably, ARIMA and ETS exhibit varied performance across different data characteristics. ARIMA struggles with datasets like Solar, characterized by weak trending and strong seasonality, while ETS shows better adaptability. Conversely, in cases of strong trending and weak seasonality, as observed in the 'Wikipedia' dataset, ARIMA significantly outperforms ETS.

Utilizing the implementation from [21], we find that XGBoost competes well, even surpassing neural network models in certain scenarios. However, for datasets with more complex distributions like 'Solar' and 'Electricity,' advanced probabilistic estimation methods demonstrate a substantial advantage over traditional learning methods and point estimation techniques. This highlights the adaptability and strength of advanced probabilistic methods in handling intricate forecasting scenarios.

D.4 Experiments on Univariate Datasets

In pursuit of a comprehensive analysis spanning univariate and multivariate scenarios, we examined a subset of M4 [45], M5 [46], and TOURISM datasets [3]—crucial datasets for univariate time-series forecasting. Table 20 provides a quantitative assessment of the intrinsic characteristics of these new datasets, focusing on trending strength, seasonality, and data distribution complexity, as detailed in

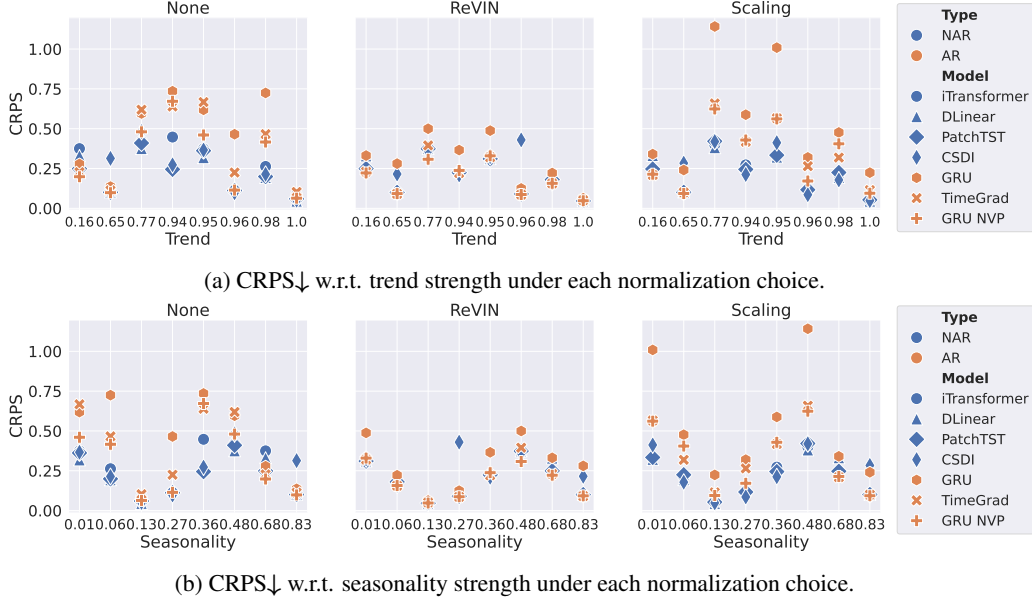


Figure 8: Impact of data characteristics on the effectiveness of different instance-level normalization strategies.

Table 19: Results of statistical models and GBDT baseline on short-term forecasting datasets.

Model	Exchange-S		Solar-S		Electricity-S		Traffic-S		Wikipedia-S	
	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE
ARIMA	0.009	0.009	1.000	1.000	0.164	0.164	0.461	0.461	0.348	0.348
ETS	0.011	0.011	0.580	0.580	0.121	0.121	0.413	0.413	0.685	0.685
ETS-prob	0.008	0.011	0.795	0.695	0.123	0.129	0.380	0.433	0.625	0.697
XGBoost	0.011	0.011	0.599	0.599	0.074	0.074	0.196	0.196	-	-
DLinear	0.012 _{.001}	0.012 _{.001}	0.547 _{.009}	0.547 _{.009}	0.095 _{.006}	0.095 _{.006}	0.273 _{.012}	0.273 _{.012}	1.046 _{.037}	1.046 _{.037}
PatchTST	0.010 _{.000}	0.010_{.000}	0.496 _{.002}	0.496 _{.002}	0.076 _{.001}	0.076 _{.001}	0.202 _{.001}	0.202 _{.001}	0.257 _{.001}	0.257_{.001}
TimesNet	0.011 _{.001}	0.011 _{.001}	0.507 _{.019}	0.507 _{.019}	0.071 _{.002}	0.071 _{.002}	0.205 _{.002}	0.205 _{.002}	0.304 _{.002}	0.304 _{.002}
GRU NVP	0.016 _{.003}	0.020 _{.003}	0.396 _{.021}	0.507 _{.022}	0.055 _{.002}	0.073 _{.003}	0.161 _{.006}	0.203 _{.009}	0.282 _{.003}	0.330 _{.003}
GRU MAF	0.015 _{.001}	0.020 _{.001}	0.386 _{.026}	0.492 _{.027}	0.051 _{.001}	0.067 _{.001}	0.131 _{.006}	0.165 _{.009}	0.281 _{.004}	0.337 _{.005}
Trans MAF	0.011 _{.001}	0.014 _{.001}	0.400 _{.022}	0.503 _{.022}	0.054 _{.004}	0.071 _{.005}	0.129_{.004}	0.165_{.006}	0.289 _{.008}	0.344 _{.008}
TimeGrad	0.011 _{.001}	0.014 _{.002}	0.359_{.011}	0.445_{.023}	0.052 _{.001}	0.067 _{.001}	0.164 _{.091}	0.201 _{.115}	0.272 _{.008}	0.327 _{.011}
CSDI	0.008_{.000}	0.011 _{.000}	0.366 _{.005}	0.484 _{.008}	0.050_{.001}	0.065_{.001}	0.146 _{.012}	0.176 _{.013}	0.219_{.006}	0.259 _{.009}

our paper. Notably, these datasets, except for M4-Daily may exhibit fewer seasonal patterns, do not introduce particularly unique characteristics.

Table 20: Quantitative assessment of the intrinsic characteristics of the univariate datasets. The JS Div. denotes Jensen–Shannon divergence, where a lower score indicates closer approximations to a Gaussian distribution.

Dataset	M4-Weekly	M4-Daily	M5	TOURISM-Monthly
Trend F_T	0.7677	0.9808	0.3443	0.7979
Seasonality F_S	0.3401	0.0467	0.2480	0.6826
JS Div.	0.5106	0.4916	0.6011	0.3291

Table 21 presents experimental results for representative methods, consistent with our initial observations. Probabilistic estimation methods like GRU NVP and TimeGrad excel on datasets with complex distributions (e.g., M4-Weekly and M5), while simpler point forecasting methods such as DLinear and PatchTST perform well on datasets with relatively simple data distribution, like TOURISM-Monthly. Both autoregressive and non-autoregressive decoding schemes show comparable performance in short-term forecasting, as discussed in the main paper.

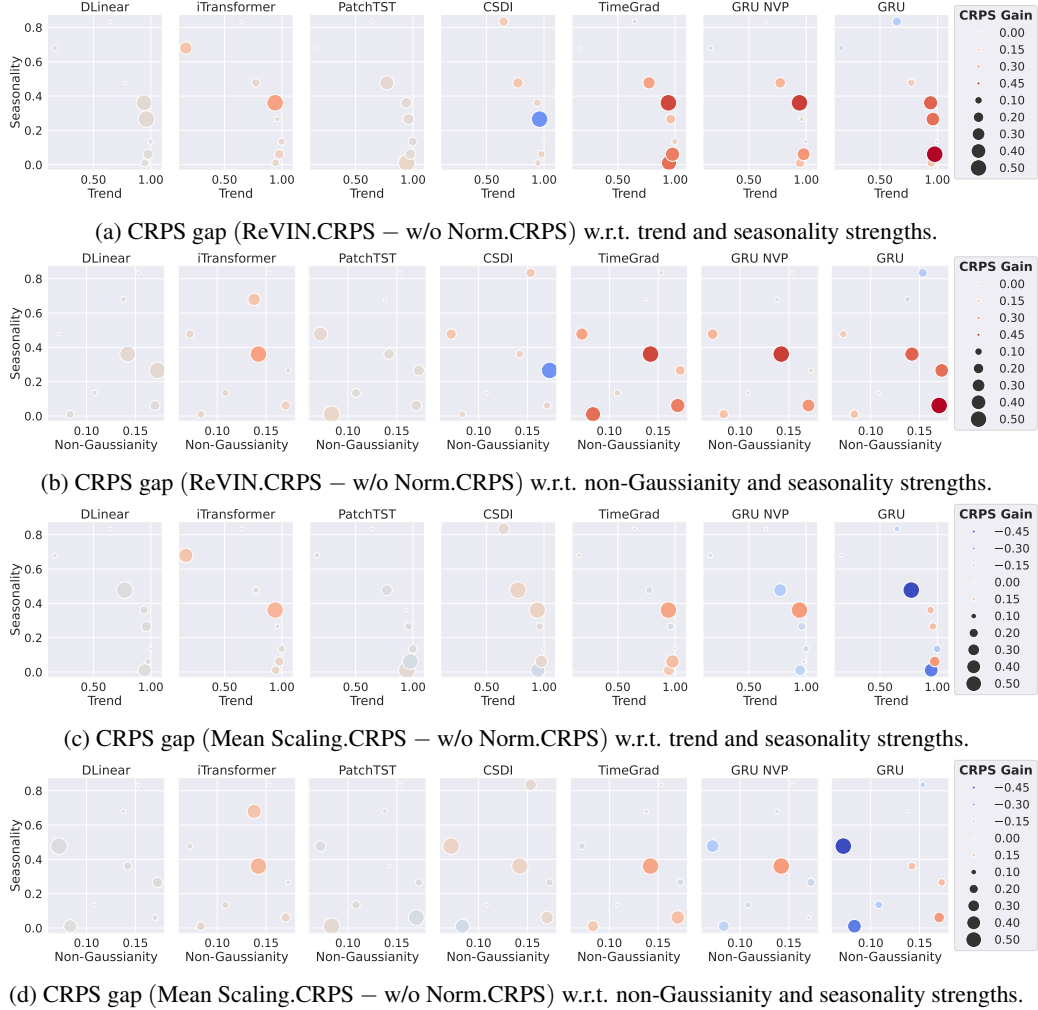


Figure 9: Impact of different instance-level normalization methods on model performance.

Table 21: Results on M4, M5, and TOURISM datasets. We utilize a lookback window of 3H, with 'H' denoting the forecasting horizon.

Dataset	DLinear		PatchTST		GRU NVP		TimeGrad	
	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE
M4-Weekly	0.081	0.081	0.089	0.089	0.066	0.077	0.055	0.065
M4-Daily	0.034	0.034	0.035	0.035	0.030	0.038	0.026	0.032
M5	0.891	0.891	0.898	0.898	0.679	0.864	-	-
TOURISM-Monthly	0.168	0.168	0.136	0.136	0.171	0.223	0.152	0.191

D.5 Experiments on Synthetic Datasets

To enhance the rigor of the insights presented, we employ synthetic datasets created with the GluonTS library¹³, encompassing a baseline dataset and variants with pronounced trends, strong seasonality, and complex data distribution (see Table 22). Specifically, we generate these datasets by superimposing four components - trend, seasonality, noise, and anomaly - each with adjustable intensity parameters. The seasonality component is defined by period hyper-parameters and intensity coefficients; the trend by slope intensity; the noise by Gaussian distribution sampling with adjustable intensity; and the anomaly by occurrence probability and maximum intensity.

¹³https://ts.gluon.ai/stable/tutorials/data_manipulation/index.html

Subsequent experiments on these synthetic datasets (refer to Table 23), using representative models, validate the empirical findings established on other datasets with ProBTs. Key observations include the declining performance of autoregressive decoding models, such as TimeGrad, in the presence of increasing trends, improved performance for models using autoregressive decoding with intensifying seasonality, and the competitive performance of probabilistic methods like CSDI in handling more complex data distributions.

Table 22: Quantitative assessment of intrinsic characteristics for synthetic datasets. The JS Div denotes Jensen–Shannon divergence, where a lower score indicates closer approximations to a Gaussian distribution.

Dataset	Normal	Strong Trend	Strong Seasonality	Complex Distribution
Trend F_T	0.105	0.554	0.105	0.064
Seasonality F_S	0.302	0.302	0.791	0.190
JS Div.	0.261	0.248	0.272	0.469

Table 23: Results on synthetic datasets. The look-back window and forecasting horizon are 30.

Model	Normal		Strong Trend		Strong Seasonality		Complex Distribution	
	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE	CRPS	NMAE
DLinear	0.013	0.013	0.001	0.001	0.014	0.014	0.301	0.301
PatchTST	0.012	0.012	0.001	0.001	0.012	0.012	0.275	0.275
TimeGrad	0.024	0.032	0.042	0.048	0.022	0.028	0.283	0.338
CSDI	0.013	0.014	0.010	0.007	0.020	0.027	0.269	0.301

D.6 Case Study

To intuitively demonstrate the distinct characteristics of point and probabilistic estimations, a case study was conducted on short-term datasets. Figure 10 illustrates that point estimation yields single-valued, deterministic estimates, in contrast to probabilistic methods, which model continuous data distributions as depicted in Figure 11. This modeling of data distributions captures the uncertainty in forecasts, aiding decision-makers in fields such as weather and finance to make more informed choices. It is also observed that while both methods align well with ground truth values in short-term forecasting datasets, they struggle to accurately capture outliers, particularly noted in the Wikipedia dataset.

D.7 Model Efficiency

For reference, detailed results regarding memory usage and time efficiency for five representative models on long-term forecasting datasets are provided here. Table 24 displays the computation memory of various models with a forecasting horizon set to 96. Additionally, Table 25 compares the inference time of these models on long-term forecasting datasets, illustrating the impact of changes in the forecasting horizon.

Table 24: Computation memory. The batch size is 1 and the prediction horizon is set to 96.

Metric	Dataset	DLinear	PatchTST	LSTM NVP	TimeGrad	CSDI
NPARAMS (MB)	ETTm1	0.075	2.145	1.079	1.233	1.720
	Electricity-L	0.076	2.146	3.680	3.472	1.370
	Traffic-L	0.078	2.149	15.926	8.298	1.390
	Weather-L	0.075	2.145	3.085	0.574	1.721
	Exchange-L	0.075	0.135	1.979	0.488	1.720
Max GPU Mem. (GB)	ETTm1	0.002	0.009	0.010	0.012	0.027
	Electricity-L	0.060	0.068	0.129	0.128	1.411
	Traffic-L	0.161	0.168	0.361	0.333	9.102
	Weather-L	0.004	0.012	0.021	0.012	0.070
	Exchange-L	0.002	0.002	0.013	0.008	0.030

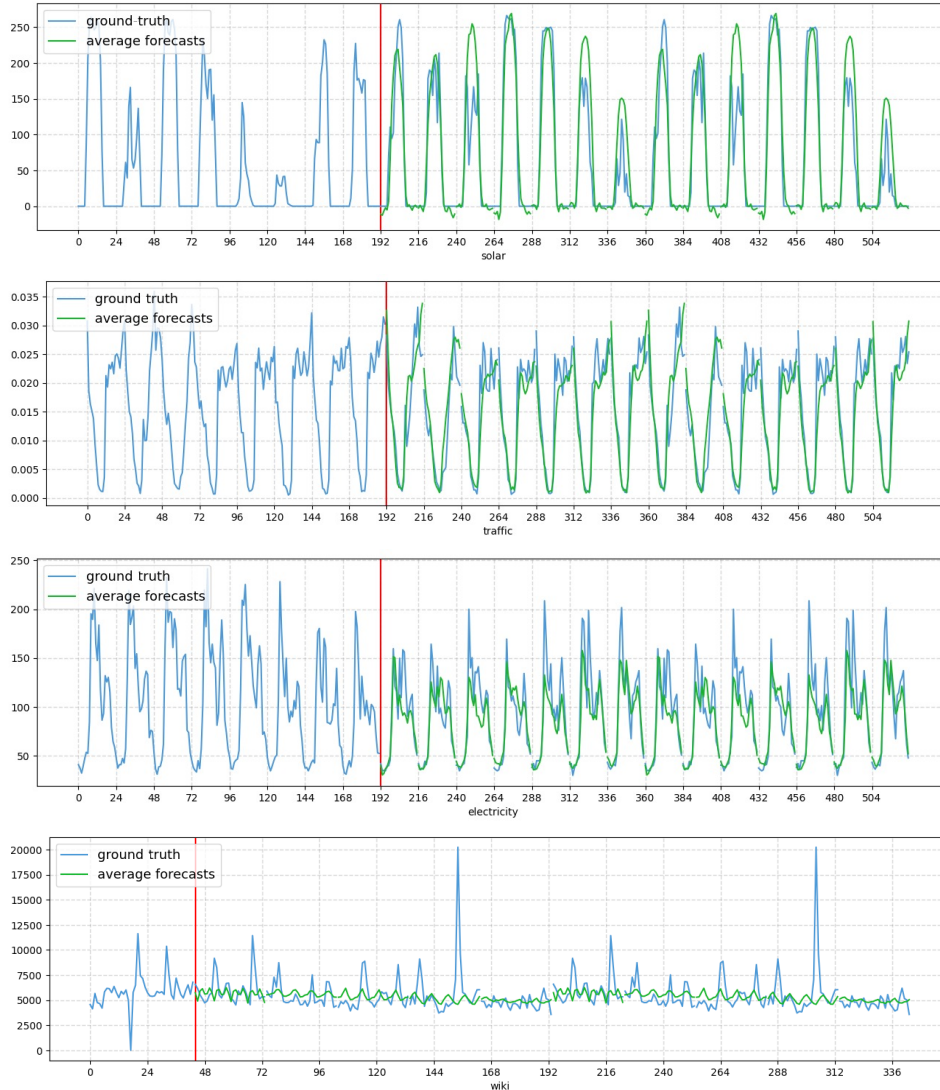


Figure 10: Point forecasts from the PatchTST model and the ground-truth value on short-term forecasting datasets.

D.8 Further discussion on uni- / multi-variate modeling

Our experiments indicate that the choice between univariate and multivariate modeling is not a primary factor in fulfilling the essential forecasting needs considered in this paper.

D.8.1 Discussion

We discuss the differences between univariate and multivariate modeling from two perspectives:

- **Dataset Perspective:** Whether the dataset is prepared for univariate or multivariate benchmarking.
- **Model Perspective:** How the model handles multivariate data, treating each variable channel independently or not.

Dataset Perspective All datasets listed in Table 5 are typically referred to as multivariate datasets, indicating that there may be strong connections across different variables. Despite this implication, when developing forecasting models, we can treat each variable channel independently, essentially turning a multivariate dataset into a univariate setup. In contrast, some datasets, like M4, M5, and

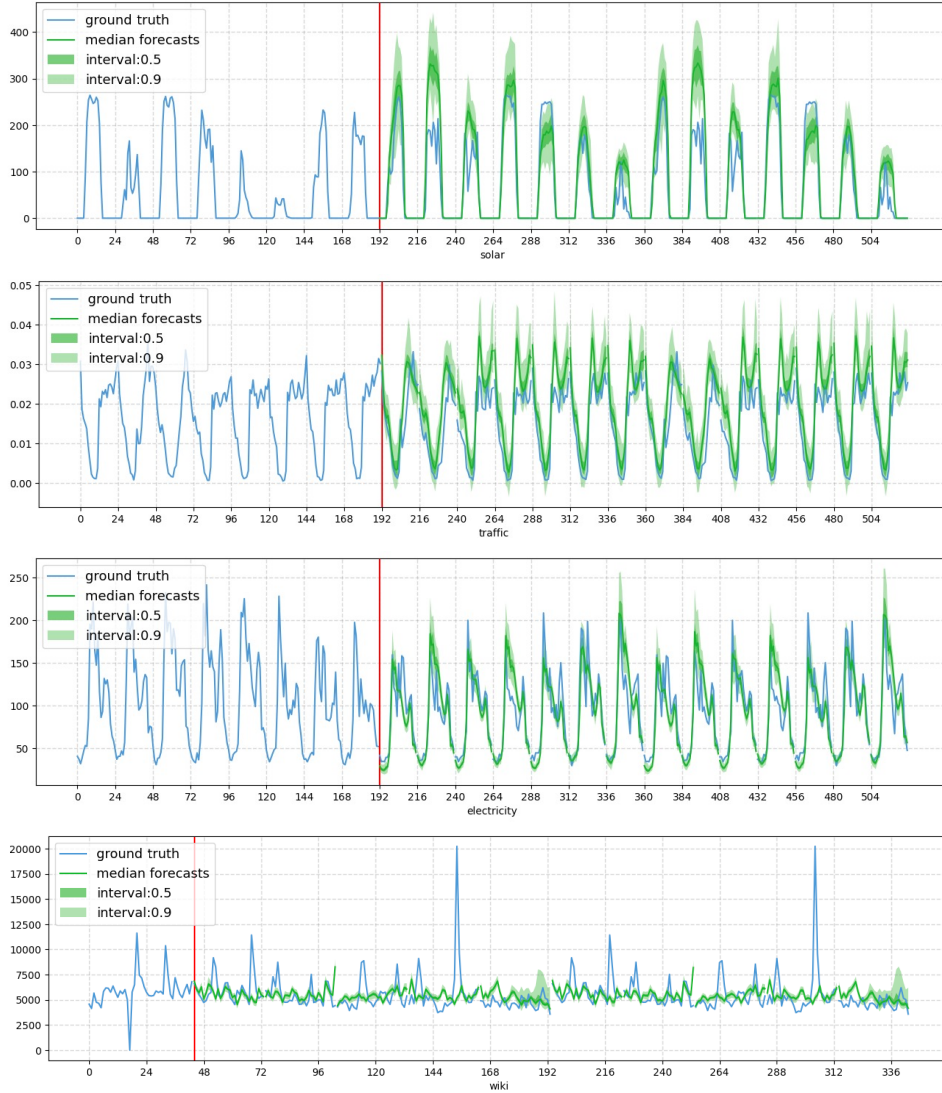


Figure 11: Forecasting intervals from the TimeGrad model and the ground-truth value on short-term forecasting datasets.

TOURISM listed in Table 20, explicitly serve univariate modeling. We rarely see multivariate models being developed for these univariate cases.

Model Perspective We have observed different preferences for univariate and multivariate modeling. Existing models can be categorized into three groups:

- **Native Univariate Models**
 - Classical models like N-BEATS and N-HiTS.
 - Most time-series foundation models, such as TimesFM and Chronos.
- **Native Multivariate Models**
 - Most probabilistic models, such as CSDI and TimeGrad.
 - Some point forecasting models, such as Informer and Autoformer.
- **Hybrid Models of Univariate and Multivariate Modeling**
 - Some classical models, such as PatchTST.
 - Some time-series foundation models, such as MOIRAI.

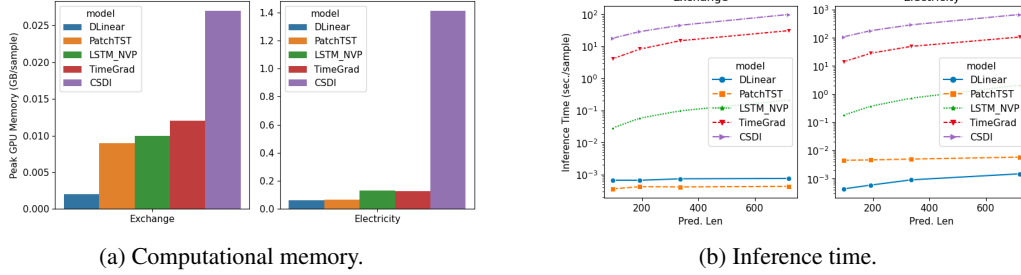


Figure 12: Comparison of computational efficiency. The forecasting horizon is set to 96 for calculating memory usage.

Table 25: Comparison of inference time (sec./sample).

Dataset	Pred len	DLinear	PatchTST	LSTM NVP	TimeGrad	CSDI
ETTm1-L	96	0.0003 ± 0.0000	0.0003 ± 0.0000	0.0352 ± 0.0007	4.1067 ± 0.0504	16.3280 ± 0.0747
	192	0.0003 ± 0.0000	0.0003 ± 0.0000	0.0697 ± 0.0020	7.8979 ± 0.0403	25.8378 ± 0.3124
	336	0.0003 ± 0.0000	0.0003 ± 0.0000	0.1221 ± 0.0044	13.6197 ± 0.1023	39.8832 ± 0.2157
	720	0.0004 ± 0.0000	0.0003 ± 0.0000	0.2603 ± 0.0020	28.6074 ± 1.1346	86.1862 ± 0.1863
Electricity-L	96	0.0004 ± 0.0000	0.0045 ± 0.0001	0.1783 ± 0.0006	13.8439 ± 0.0054	388.3150 ± 0.2155
	192	0.0006 ± 0.0000	0.0046 ± 0.0000	0.3700 ± 0.0010	27.6683 ± 0.0368	659.4284 ± 0.2003
	336	0.0008 ± 0.0000	0.0049 ± 0.0000	0.7157 ± 0.0028	48.4456 ± 0.0279	-
	720	0.0015 ± 0.0000	0.0057 ± 0.0000	2.0785 ± 0.0186	104.1473 ± 0.1465	-
Traffic-L	96	0.0010 ± 0.0001	0.0102 ± 0.0000	0.3695 ± 0.0022	31.7644 ± 0.0101	-
	192	0.0013 ± 0.0000	0.0106 ± 0.0000	0.8287 ± 0.0094	63.5832 ± 0.0060	-
	336	0.0020 ± 0.0000	0.0114 ± 0.0001	1.6945 ± 0.0026	111.4147 ± 0.0169	-
	720	0.0039 ± 0.0000	0.0137 ± 0.0000	5.0963 ± 0.0018	258.1274 ± 0.6088	-
Weather-L	96	0.0002 ± 0.0000	0.0004 ± 0.0000	0.0800 ± 0.0016	4.1261 ± 0.0812	37.8984 ± 0.0782
	192	0.0003 ± 0.0000	0.0004 ± 0.0000	0.1568 ± 0.0008	8.2913 ± 0.5544	62.0223 ± 0.2329
	336	0.0003 ± 0.0000	0.0004 ± 0.0000	0.2482 ± 0.0297	14.2391 ± 0.4891	96.8704 ± 0.2258
	720	0.0003 ± 0.0000	0.0005 ± 0.0000	0.5447 ± 0.0249	29.4407 ± 0.3519	216.6044 ± 0.4253
Exchange-L	96	0.0006 ± 0.0000	0.0004 ± 0.0000	0.0284 ± 0.0001	4.1069 ± 0.0981	17.8655 ± 0.1282
	192	0.0007 ± 0.0000	0.0004 ± 0.0000	0.0563 ± 0.0008	8.1576 ± 0.0911	28.5456 ± 0.0873
	336	0.0007 ± 0.0000	0.0004 ± 0.0000	0.0966 ± 0.0007	14.4593 ± 0.4466	44.9733 ± 0.3820
	720	0.0007 ± 0.0000	0.0004 ± 0.0000	0.2085 ± 0.0046	30.1443 ± 0.5378	97.7417 ± 0.2606
ILI-L	24	0.0002 ± 0.0000	0.0008 ± 0.0001	0.0080 ± 0.0001	1.0427 ± 0.0190	12.4038 ± 0.1681
	192	0.0002 ± 0.0000	0.0008 ± 0.0000	0.0121 ± 0.0003	1.5762 ± 0.0282	12.7187 ± 0.1344
	336	0.0002 ± 0.0000	0.0008 ± 0.0000	0.0155 ± 0.0002	2.1344 ± 0.0660	12.7386 ± 0.1868
	720	0.0002 ± 0.0000	0.0008 ± 0.0000	0.0196 ± 0.0004	2.5787 ± 0.0594	12.5407 ± 0.0481

Native univariate models can also be applied to multivariate datasets by treating them as univariate cases. Similarly, native multivariate models can be applied to univariate datasets by setting the variable dimension to 1. Hybrid models typically include specific modes to activate univariate and multivariate functionalities. For example, in PatchTST, we can use a shared forecasting head for univariate modeling or assign a specific forecasting head for each variable channel to differentiate different variables.

We compile a summary table (Table 26) delineating how models from each branch address the multivariate aspect. Despite a thorough investigation, we have not identified a clear pattern linking the modeling of cross-channel interactions to overall model performance. A notable trend is the prevalent use of a channel-mixing approach in most studies. However, findings are diverse; models like DLinear and PatchTST suggest that processing channels independently can yield superior results, while others like CSDI indicate that explicit modeling of cross-channel interactions offers significant advantages. This diversity underscores the ongoing exploration of the impact of cross-channel interactions on forecasting performance.

Table 26: Summary of how existing models handle multivariate time series.

Model	Research branch	Process channels independently
Customized neural architectures	N-BEATS [53]	✓
	N-HiTS [11]	✓
	Autoformer [72]	✗
	Informer [78]	✗
	LTSF-Linear [75]	✗/✓
	PatchTST [49]	✗/✓
	TimesNet [71]	✗
Probabilistic estimation	DeepAR [60]	✓
	GP-copula [59]	✗
	LSTM NVP [58]	✗
	LSTM MAF [58]	✗
	Trans MAF [58]	✗
	TimeGrad [57]	✗
	CSDI [62]	✗
	SPD [7]	✗

D.8.2 Additional Experiments on uni- / multi-variate modeling

In Table 27, we include additional experiments comparing univariate and multivariate modeling of PatchTST across different datasets. Our observation is that there is no definitive answer as to which approach is superior; it depends on the nature of the dataset. Some datasets benefit from modeling variable correlations, while others perform better with independent modeling. The performance gaps are not significant.

Similar observations have been reported in MOIRAI, which allows either univariate or multivariate modes by controlling its cross-variate attention masks. When applied to a downstream forecasting scenario, it can search over the validation set to determine which configurations to activate. We believe such a design could serve as a good example of unifying univariate and multivariate.

Table 27: The comparison of PatchTST on univariate and multivariate modeling.

Dataset	Pred. Horizon	PatchTST (Multivariate)	PatchTST (Univariate)
ETTh1	96	0.3239	0.3212
	192	0.3609	0.3562
	336	0.3763	0.3737
	720	0.3882	0.3909
ETTm1	96	0.2652	0.2739
	192	0.2926	0.2961
	336	0.3101	0.3188
	720	0.345	0.3463
Electricity	96	0.0832	0.0857
	192	0.0899	0.0912
	336	0.0995	0.1001
	720	0.1183	0.116
Exchange	96	0.0243	0.0235
	192	0.0348	0.0336
	336	0.0471	0.0462
	720	0.0787	0.0777
Weather	96	0.0872	0.0837
	192	0.0924	0.0858
	336	0.0934	0.0903
	720	0.0993	0.0953

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Sections 5.
 - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments (e.g. for benchmarks)...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We release the ProBTs toolkit, documentation, and running scripts at <https://github.com/microsoft/ProBTs>. The repository includes parameter configurations for benchmarking experiments, ensuring reproducibility of all results presented in the paper.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] All datasets in this paper are split according to standard practices in previous works. We have specified the hyperparameters in Appendix B.3 The detailed configuration file for each baseline is available on our GitHub page.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We specify a fixed set of random seeds for each experiment to indicate the error bars. We show the mean and standard deviation of evaluation scores collected from multiple runs.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix B.3.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes] See Appendix B.4.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [Yes] The data utilized in this study is open-source and does not necessitate consent. Acquisition details are available on our GitHub page.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] The data we are using contains no personally identifiable information or offensive content.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]