

# GIFT-EVAL: A BENCHMARK FOR GENERAL TIME SERIES FORECASTING MODEL EVALUATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Time series foundation models excel in zero-shot forecasting, handling diverse tasks without explicit training. However, the advancement of these models has been hindered by the lack of comprehensive benchmarks. To address this gap, we introduce the **General Time Series Forecasting Model Evaluation**, GIFT-Eval, a pioneering benchmark aimed at promoting evaluation across diverse datasets. GIFT-Eval encompasses 23 datasets over 144,000 time series and 177 million data points, spanning seven domains, 10 frequencies, multivariate inputs, and prediction lengths ranging from short to long-term forecasts. To facilitate the effective pretraining and evaluation of foundation models, we also provide a non-leaking pretraining dataset containing approximately 230 billion data points. Additionally, we provide a comprehensive analysis of **R: [20]** baselines, which includes statistical models, deep learning models, and foundation models. We discuss each model in the context of various benchmark characteristics and offer a qualitative analysis that spans both deep learning and foundation models. We believe the insights from this analysis, along with access to this new standard zero-shot time series forecasting benchmark, will guide future developments in time series foundation models.

## 1 INTRODUCTION

The success of foundation model pretraining in language and vision modalities has catalyzed similar progress in time series forecasting. By pretraining on extensive time series datasets, a universal forecasting model can be developed, equipped to address varied downstream forecasting tasks across multiple domains, frequencies, prediction lengths, and number of variates in a zero-shot manner (Woo et al., 2024; Rasul et al., 2023b; Ansari et al., 2024).

A critical aspect of foundation model research is creating a high-quality benchmark that includes large, diverse evaluation data, and preferably non-leaking pretraining data to fairly evaluate models and identify their weaknesses. Research in Natural Language Processing (NLP) has produced key benchmarks such as GLUE, MMLU, *etc.* (Wang et al., 2018; Hendrycks et al., 2020; Srivastava et al., 2022; Chen et al., 2021), which are crucial for developing high-quality models.

Unlike NLP, time series foundation models lack a unified, diverse benchmark for fair comparison. For instance, Woo et al. (2024) introduces LOTSA, which remains the largest collection of time series forecasting pre-training data to date. However, the proposed architecture, *Moirai*, is evaluated on existing benchmarks that are tailored to specific forecasting tasks, such as the LSF (Zhou et al., 2020) dataset for long-term forecast, and the Monash (Godaheva et al., 2021) dataset for univariate forecasts. Both datasets lack sufficient diversity in time series characteristics and forecasting tasks, making it challenging to evaluate the zero-shot capabilities of foundation models in handling broad and generalized forecasting scenarios. This limitation remains in the recent empirical evaluations of other foundation models, including those featured in benchmarks such as TimesFM, Chronos, and Lag-Llama (Das et al., 2023b; Ansari et al., 2024; Rasul et al., 2023b). Furthermore, the inconsistency in pretraining, training, and test splits across various foundation models complicates comparisons and poses a risk of data leakage during in-domain and out-of-domain evaluations. To accelerate the advancement for research on time series model, it is essential to establish a high-quality and diverse benchmark that supports universal forecasting evaluation.

To fill identified gaps, we introduce the **General Time Series Forecasting Model Evaluation** (GIFT-Eval), consisting of distinct pretraining and train/test components. The pretraining component

Table 1: Property comparisons of various forecasting benchmarks.

Benchmark	Property	Data			Forecasting Task		Evaluation	
		Freq. Range	Num. of Domain	Pretraining data	Num. of var.	Pred. Len.	Benchmark Methods	Prob. Forecasting
Monash (Godaheewa et al., 2021)	Secondly ~ Yearly	7	No	Uni	Short	Stat./DL	No	
TFB (Qiu et al., 2024)	Minutely ~ Yearly	6	No	Uni/Multi	Short	Stat./DL	No	
LTSF (Zeng et al., 2022)	Minutely ~ Weekly	5	No	Multi	Long	Stat./DL	No	
BasicTS+ (Shao et al., 2023)	Minutely ~ Daily	3	No	Multi	Short/Long	Stat./DL	No	
R: [ProbTS (Zhang et al., 2023)]	Minutely ~ Weekly	5	No	Multi	Short/Long	Stat./DL/FM	Yes	
GIFT-Eval (our work)	Secondly ~ Yearly	7	Yes	Uni/Multi	Short/Long	Stat./DL/FM	Yes	

features 88 datasets including 240 billion data points (Appendix E lists more details on pretraining data). The train/test component features 23 datasets encompassing 144,000 time series and 177 million data points across seven domains and 10 frequencies, with prediction lengths ranging from short to long-term, as well as univariate and multivariate forecasting settings. Prior to our work, Qiu et al. (2024) introduced TFB, a comprehensive dataset for time series forecasting. While it offered diversity in the number of variates and domains, it lacks the evaluation of foundation models and accompanying pretraining data without leakage. Our benchmark fills these gaps and it also includes a broader range of frequencies, a more diverse taxonomy, and a wider span of prediction lengths. We compare GIFT-Eval with other similar benchmarks in Table 1. Our contributions are three-fold:

- **GIFT-Eval:** We introduce a general time series forecasting benchmark that evaluates the zero-shot and universal forecasting capabilities of foundation models. We provide pretraining and train-test components that ensure diversity across multiple characteristics and time series features.
- **Comprehensive Benchmarking:** We design diverse forecasting tasks and evaluate R: [20 baselines] that encompass statistical, deep learning, and foundational models on GIFT-Eval.
- **Detailed Analysis:** We provide insights into the strengths of different models on all aspects of GIFT-Eval including domains, frequencies, prediction lengths, and the number of variates R: [and also among 6 time series features]. We further provide a qualitative analysis showing failure cases of both deep learning and foundation models. We believe these insights will contribute to the future development of foundation models.

## 2 RELATED WORK

**Forecasting Methods** Time series forecasters can be broadly categorized into statistical models, deep learning models, and, more recently, foundation models. Statistical models rely solely on historical data statistics to predict future values. Among these, ARIMA (Box & Pierce, 1970), ETS (Hyndman et al., 2008), Theta (Garza et al., 2022), and VAR (Godaheewa et al., 2021) are some of the most widely used ones. With the advent of deep learning technologies, models that apply these techniques to time series forecasting have emerged. Examples include DeepAR (Flunkert et al., 2017), N-BEATS (Oreshkin et al., 2019), and DLinear (Zeng et al., 2022), which utilize pre-transformer architectures. Additionally, transformer-based models such as PatchTST (Nie et al., 2022), Autoformer (Wu et al., 2021), and Crossformer (Zhang & Yan, 2023) have been developed. R: [Another line of important work are probabilistic forecasting models. TimeGrad (Rasul et al., 2021) is an autoregressive probabilistic forecasting model utilizing diffusion probabilistic methods, CSDI (Tashiro et al., 2021) is a time series imputation approach that leverages score-based diffusion models conditioned on observed data. Their predecessor GRU NVP (Rasul et al., 2020) on the other hand, models multivariate temporal dynamics in time series forecasting using an autoregressive deep learning model combined with conditioned normalizing flows.] In the last few years, foundation models have been proposed, inspired by their success in other modalities like language and vision. The multivariate Moirai (Woo et al., 2024) forecaster, for instance, is based on an encoder-decoder architecture pretrained on a large dataset. Conversely, Chronos (Ansari et al., 2024) and TimesFM (Das et al., 2023b), R: [and Lag-Llama (Rasul et al., 2023a) are univariate forecasters trained using a decoder-only model. Following these other foundation models have also been proposed Timer (Liu et al., 2024), UniTS (Gao et al., 2024), TTM (Ekambaram et al., 2024), Moment (Goswami et al., 2024), and VisionTS (Chen et al., 2024).] However, the main bottleneck in building and evaluating these foundation models is the lack of a diverse and large benchmark dataset.

**Forecasting Benchmarks** To address this challenge, several efforts have been made to develop extensive time series benchmarks. Woo et al. (2024) introduced LOTSA, which holds the title for

the largest collection of open time series datasets, encompassing 231 billion data points across nine domains. Despite its vast size, the evaluation datasets reuse existing benchmarks from the time series forecasting community and still lack sufficient variety in terms of time series data characteristics and forecasting tasks, which our benchmark aims to augment. Ansari et al. (2024) developed a dataset specifically structured for pretraining, in-domain evaluation, and zero-shot evaluation splits. However, their work is constrained by a limited range of prediction lengths (from 6 to 56), which excludes long-term forecasts, and it restricts the data to univariate forecasting. In contrast, our benchmark encompasses extensive multivariate scenarios and evaluates diverse data across various domains and frequencies. The corpus by Rasul et al. (2023b) presents a diverse array of domains, yet it comprises only univariate datasets totaling 8,000 time series. In contrast, GIFT-Eval dramatically expands this scope with 144,000 time series, enhancing the breadth and depth of the dataset. The benchmark by Qiu et al. (2024) is closely aligned with our work in its aim to curate a diverse and comprehensive set of data. However, it lacks pretraining data, does not evaluate foundation models, and limits the taxonomy to time series features only. Our benchmark not only includes pretraining data (with zero-shot evaluation support) but also provides evaluations for foundation models and offers a taxonomy over both characteristics and time series properties. In summary, our benchmark, GIFT-Eval, builds upon and seeks to address the gaps identified in existing time series forecasting benchmarks. We provide a wider comparison with more benchmarks in Table 1. By providing a more diverse and extensive dataset, we aim to facilitate the development and evaluation of foundation models in time series forecasting.

**Forecasting Tools** R: [Apart from raw benchmarks, there are also frameworks that provide access to a range of time series datasets. Prophet (Taylor & Letham, 2018), sktime (Löning et al., 2019), TSLib (Wu et al., 2022) are primarily implemented for point forecasting. On the other hand PyTorchTS (Rasul, 2021), GluonTS (Alexandrov et al., 2020b), NeuralForecast<sup>1</sup> are Python packages for probabilistic time series forecasting, with PyTorchTS including more advanced probabilistic models based on deep generative models. ProbTS (Zhang et al., 2023), on the other hand, differs from these by providing insights into point vs. probabilistic forecasts and short- vs. long-term forecasting. It is the most relevant to our work, as it includes both classical and foundation models. However, our benchmark is larger in scale, as ProbTS incorporates only 12 multivariate datasets, with no univariate datasets included.]

### 3 GIFT-EVAL

In this section, we first provide a background on time series forecasting tasks and define key characteristics and features of time series data. We then outline the design decisions behind the development of GIFT-Eval, concluding with an analysis that highlights the key features of its final distribution.

#### 3.1 BACKGROUND

We start by defining univariate and multivariate forecasting tasks. After that, we outline the fundamental characteristics of time series datasets which also influenced our data collection process, including domain, frequency, number of variates, and prediction length. We also introduce time series features as part of our data analysis.

##### 3.1.1 TIME SERIES FORECASTING

Time series forecasting is a task of predicting future values over one (univariate) or more (multivariate) variates given historical (most commonly real-valued) data which is sampled at regular time intervals. Suppose  $D = (Y^i, Z^i)_{i=1}^N$  is a dataset of  $N$  time series where  $Y^i = (y_1^i, y_2^i, \dots, y_{T_i}^i) \in \mathbb{R}^{d_{y_i} \times T_i}$  is the target time series with  $d_{y_i}$  variates and  $T_i$  time steps and  $Z^i = (z_1^i, z_2^i, \dots, z_{T_i}^i) \in \mathbb{R}^{d_{z_i} \times T_i}$  are the set of covariates with  $d_{z_i}$  variates. Then the forecasting task can be modeled as the predictive distribution:  $p(Y_{t:t+h} | Y_{t-l:t}, Z_{t-l:t+h})$  where  $l$  is the context length, and  $h$  is the forecast horizon. Univariate forecasting is a special case where the target series is univariate (*i.e.*,  $d_{y_i} = 1$ ), no

<sup>1</sup><https://github.com/Nixtla/neuralforecast>

covariates are used (*i.e.*,  $Z = \emptyset$ ), and only the historical values of the target time series are utilized for prediction.

### 3.1.2 TIME SERIES CHARACTERISTICS AND FEATURES

**Characteristics** Time series datasets possess inherent characteristics that define their structure, and common patterns observed in the data and even choices of modelling techniques. We believe a universal forecasting model should be able to perform irrespective of the domain from which the data is sourced, the granularity at which it was sampled, the length of the forecast horizon and whether it is univariate or multivariate. Thus in our study, we focus on these four characteristics: *(i) Domain*, *i.e.*, the field or industry from which the time series data originates, such as finance, healthcare or meteorology. The domain often has a direct effect on the nature of patterns. Another crucial aspect is *(ii) the frequency* of observations, indicating the time intervals at which the data points are recorded – such as hourly, daily, monthly or annually. *(iii) Prediction length*, or forecast horizon, is the number of future time steps for which predictions are expected. Lastly, *(iv) the number of variates* pertains to the dimensionality of the time series data. A *univariate* time series consists of observations of a single variable over time, whereas a *multivariate* time series involves multiple interrelated variables. The number of variates adds complexity to the modelling process, as models need to account for dependencies among multiple time series. By ensuring diversity across these specific characteristics in our benchmark, we aim to encompass a wide array of real-life scenarios.

**Features** Time series features<sup>2</sup> are statistical properties that capture essential characteristics of the data. We have selected six such properties to analyze our benchmark, grouped into three categories based on the aspects they assess, *c.f.* Appendix B for a detailed explanation and formula of each feature. First, we chose two metrics for assessing the temporal attributes of each time series: *(i) Trend* refers to the progression of the time series, indicating whether the data shows an overall increase, decrease or stability over time, where higher values indicate stronger trends. *(ii) Seasonal strength* measures the extent to which regular, repeating patterns occur at specific intervals, such as daily cycles in energy consumption, or annual peaks in finance. The higher the value the more repeating patterns the data exhibits. Second, to assess the forecastability of the time series, we included two metrics: *(iii) Entropy* measures the “forecastability” of a time series, where low values indicate a high signal-to-noise ratio and high values occur when a series is difficult to forecast. *(iv) Hurst exponent* quantifies the long-term memory or persistence in a time series. It indicates whether future values are likely to be influenced by past trends, revert to the mean, or behave randomly, where higher values indicate more persistence. Lastly, to understand the regularity and variability within the time series, we selected two metrics: *(v) Stability* assesses the inconsistency of the mean of the time series. In simpler terms, it can be defined as the variance of the means. Note that, unlike what the name suggests, lower values indicate more stable data. Finally, *(vi) Lumpiness* quantifies the variability of the variance across different segments of the time series. A high value of lumpiness indicates significant fluctuations in variability, which can be challenging to model due to the inconsistent behavior of the data.

## 3.2 DATASETS

To evaluate and advance universal time series forecasting methods, we have curated a comprehensive collection of datasets. Our compilation spans a wide array of domains with varying frequencies, numbers of variates, and prediction lengths. This diversity is crucial for assessing the generalization capabilities of forecasting models across different types of time series data. In the following sections, we provide detailed descriptions of GIFT-Eval and its unique splits, outlining their sources, and key properties. We also conduct a detailed analysis on the test data to gain a better understanding of the datasets’ characteristics and the distribution of time series features.

**Train/Test Data** We curated the train/test portion of GIFT-Eval with 15 univariate and eight multivariate datasets, spanning seven domains and 10 frequencies, totaling 144,000 time series and 177 million data points. We adhere to established prediction lengths for well-known datasets like M4 (Makridakis et al., 2018). For other datasets, we establish three prediction settings—short, medium, and long—based on frequency and domain, with medium and long settings extending the

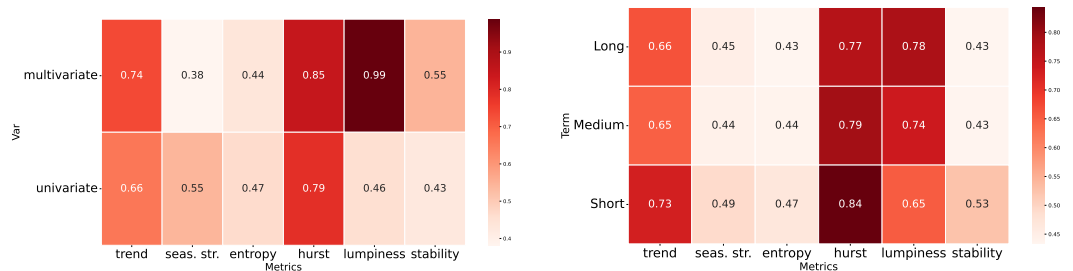
<sup>2</sup>We use the python implementation of tsfeatures library (Garza et al., 2024) to calculate each feature.

short-term length by factors of 10 and 15, respectively. To support models without multivariate forecasting, our framework flattens multivariate datasets for broader compatibility. Data is stored in the Arrow format (Richardson et al., 2023), ensuring efficient integration into deep learning pipelines. Our benchmark features 97 unique triplets of dataset, frequency, and length, with aggregated results for each model reported across these configurations. The sources of each dataset used in train/test split can be found in Appendix D.

We structure the evaluation component of our benchmark by dedicating the final 10% of each dataset in train/test portion to testing, with the rest allocated for training. A non-overlapping rolling evaluation method is employed, setting a predetermined number of windows in the test split, each equal to the dataset’s prediction length. The final window of the training data serves as validation for tuning deep learning model hyperparameters.

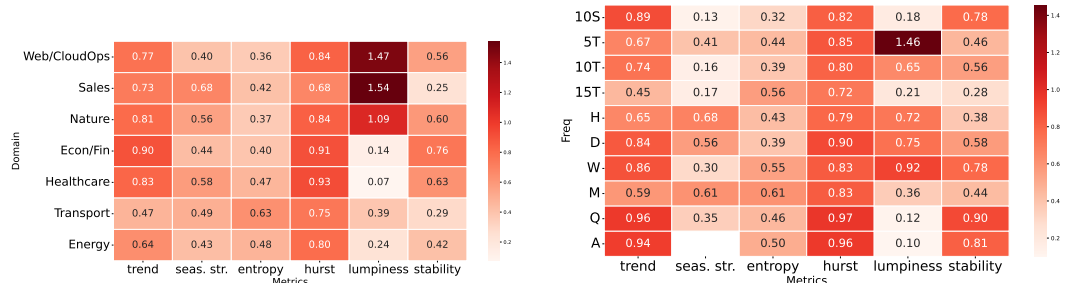
**Analysis over test data** We analyze GIFT-Eval to understand the distribution and characteristics of the time series features across various datasets, for more granular information see Appendix D. Figure 1 illustrates the mean values of each time series features across different dataset characteristics. These heatmaps provide valuable insights into how metrics such as trend, seasonal strength, Hurst exponent, stability, and lumpiness vary across datasets with different domains, frequencies, prediction lengths, and numbers of variates. This visualization aids in identifying patterns and potential biases within the data, ensuring that the benchmark captures a diverse range of time series behaviors. It also facilitates fine-grained analysis of model performance across varying dataset characteristics, offering a comprehensive comparison.

**Number of variates:** Figure 1(a) depicts that multivariate data exhibit higher stability and lumpiness values, suggesting more fluctuation in variance across different segments, indicating multivariate time series are more complex and potentially more challenging to model. Conversely, univariate series show stronger seasonal strength, reflecting more pronounced and regular repeating patterns, making them more predictable over certain periods. Note that the metrics on multivariate time series are calculated individually for each variate and aggregated for each dataset.



(a) Mean values of TS features across univariate and multivariate datasets.

(b) Mean values of TS features across different prediction lengths.



(c) Mean values of TS features across different domains.

(d) Mean values of TS features across different frequencies.

Figure 1: Heatmaps depicting mean values of six time series features across different characteristics.

**Prediction Length:** Figure 1(b) shows that shorter prediction lengths have higher values for both trend and Hurst metrics. This suggests that time series with shorter forecast horizons exhibit stronger directional movements and greater persistence in their trends, making them potentially easier to

270 predict. As the prediction length increases, the trend and hurst values tend to decrease significantly  
 271 which makes forecasting harder. Notably, the stability values decrease from short to long indicating  
 272 higher steadiness in long term while lumpiness increases suggesting higher fluctuations in different  
 273 sections of the data.

274 **Domain:** Figure 1(c) reveals distinct patterns in the metrics. The Web/CloudOps, Sales, and Nature  
 275 domains exhibit notably high lumpiness, indicating significant fluctuations in variance. This may  
 276 reflect the volatile nature of online operations, sales dynamics and weather predictions. On the other  
 277 hand, Transport shows the highest entropy and lowest trend values, indicating less predictability,  
 278 likely due to the variable nature of transportation data influenced by numerous external factors. The  
 279 Econ/Fin domain shows the highest trend values, indicating strong directional movements that may  
 280 imply clear market trends or economic cycles. Finally, healthcare exhibits the highest Hurst and  
 281 lowest lumpiness values, suggesting persistence in the data, possibly due to consistent patient trends  
 282 or medical outcomes over time.

283 **Frequency:** Figure 1(d) lists frequencies from highest to lowest. Data with very short intervals,  
 284 such as secondly (S) and minutely (T) exhibit the lowest seasonal strengths and poor steadiness,  
 285 indicative of the erratic and volatile nature typical at these granular levels. There is a noticeable  
 286 increase in seasonal strength progressing from secondly and minutely data to hourly (H) and daily  
 287 (D). Finally, yearly (A) and quarterly (Q) data demonstrate the strongest trends and Hurst values,  
 288 with notably low lumpiness, suggesting increased persistence and high predictability. Notably, the  
 289 yearly data lack seasonal strength measurements due to the tsfeatures library not providing seasonal  
 290 strength for excessively long time series, a limitation commonly observed in low-frequency datasets.

291 **Summary:** These observations confirm that our benchmark is rich and diverse, representing a broad  
 292 range of real-life time series scenarios. Our dataset encompasses various characteristics—such as  
 293 differing levels of trend, seasonality, persistence, volatility, and complexity—across multiple domains  
 294 and frequencies. For instance, we include data from domains with high volatility and significant  
 295 fluctuations, as well as data exhibiting strong persistence and stability. We also cover a wide spectrum  
 296 of frequencies, from high-frequency data with erratic patterns to low-frequency data with strong  
 297 trends and greater predictability. In a similar manner, our metrics show diversity across variate  
 298 types and prediction lengths. This diversity ensures that models are tested across various temporal  
 299 behaviors, making our benchmark a robust platform for evaluating the general capabilities of unified  
 300 models, particularly foundation models in time series forecasting.

301 **Pretraining Dataset** We have also curated a pretraining dataset aligned with GIFT-Eval that has  
 302 71 univariate and 17 multivariate datasets, spanning seven domains and 13 frequencies, totaling 4.5  
 303 million time series and 230 billion data points. Notably this collection of data has no leakage issue  
 304 with the train/test split and can be used to pretrain foundation models that can be fairly evaluated on  
 305 GIFT-Eval. Further details on pretraining dataset can be found in Appendix E.

## 307 4 EXPERIMENTS

308 In this section, we present the experimental evaluation of GIFT-Eval across various models.

309 **Models** Time series forecasting training and inference may take different forms for different fami-  
 310 lies of models. Statistical models make predictions by directly analyzing patterns in the historical  
 311 data without a separate training phase. We incorporate five statistical models in our benchmark:  
 312 Naive, Seasonal Naive (Hyndman & Athanasopoulos, 2018), Auto\_Arima, Auto\_ETs, and Auto\_Theta (Garza et al., 2022) methods. Deep learning models require training a specific  
 313 model instance for each dataset. Representing deep learning, we select 8 models: DeepAR (Flunkert  
 314 et al., 2017), TFT (Lim et al., 2019), TiDE (Das et al., 2023a), N-BEATS (Oreshkin et al., 2019),  
 315 PatchTST (Nie et al., 2022), DLinear (Zeng et al., 2022), Crossformer (Zhang & Yan, 2023)  
 316 and iTransformer (Liu et al., 2023b). To obtain both point and probabilistic forecasts, we  
 317 either adapt models using gluonts (Alexandrov et al., 2020b) with a small probabilistic head or  
 318 implement our own modifications. We conduct an extensive hyperparameter search for each deep  
 319 learning model, see Appendix A for details. We evaluate four foundation models on our bench-  
 320 mark: TimesFM (Das et al., 2023b), Chronos (Ansari et al., 2024) available in tiny, small, and base  
 321 sizes, Moirai (Woo et al., 2024) available in small, base, and large sizes and, R: [Lag-Llama (Ra-  
 322  
 323

Table 2: Results on GIFT-Eval aggregated by domain. The best results across each row are **bolded**, while the second best results are underlined.

Domain	Metric	Nv.	S.Nv.	A.Ar.	A.Th.	D.AR	TFT	TiDE	N-B.	P.TST	iTr.	T.FM	V.TS	Chr.s	Chr.B	Chr.L	Moi.s	Moi.B	Moi.L	Best
Econ/Fin	MASE	1.43	1.00	8.66e <sup>-1</sup>	9.83e <sup>-1</sup>	1.54	1.03	1.51	8.61e <sup>-1</sup>	9.08e <sup>-1</sup>	9.89e <sup>-1</sup>	8.24e <sup>-1</sup>	9.31e <sup>-1</sup>	7.97e <sup>-1</sup>	<u>7.83e<sup>-1</sup></u>	<b>7.83e<sup>-1</sup></b>	1.04	9.27e <sup>-1</sup>	9.63e <sup>-1</sup>	Chr.L
	CRPS	1.17	1.00	8.21e <sup>-1</sup>	8.41e <sup>-1</sup>	1.22	8.41e <sup>-1</sup>	1.08	9.67e <sup>-1</sup>	8.03e <sup>-1</sup>	8.48e <sup>-1</sup>	<b>7.16e<sup>-1</sup></b>	1.05	7.63e <sup>-1</sup>	<u>7.51e<sup>-1</sup></u>	7.58e <sup>-1</sup>	7.96e <sup>-1</sup>	8.16e <sup>-1</sup>	8.47e <sup>-1</sup>	T.FM
	Rank	1.90e <sup>1</sup>	1.88e <sup>1</sup>	9.83	1.07e <sup>1</sup>	1.88e <sup>1</sup>	1.10e <sup>1</sup>	2.12e <sup>1</sup>	1.62e <sup>1</sup>	9.17	1.15e <sup>1</sup>	<b>6.67</b>	2.03e <sup>1</sup>	9.50	<b>8.67</b>	9.00	1.00e <sup>1</sup>	7.00	6.50	Moi.L
Energy	MASE	1.56	1.00	1.01	1.36	1.78	1.01	1.17	1.18	9.83e <sup>-1</sup>	1.11	1.02	9.93e <sup>-1</sup>	9.45e <sup>-1</sup>	<u>9.24e<sup>-1</sup></u>	<b>9.19e<sup>-1</sup></b>	1.04	9.87e <sup>-1</sup>	1.03	Chr.L
	CRPS	1.53	1.00	8.33e <sup>-1</sup>	1.70	1.07	6.30e <sup>-1</sup>	7.51e <sup>-1</sup>	9.35e <sup>-1</sup>	<b>6.12e<sup>-1</sup></b>	6.95e <sup>-1</sup>	6.73e <sup>-1</sup>	7.82e <sup>-1</sup>	6.48e <sup>-1</sup>	6.31e <sup>-1</sup>	6.28e <sup>-1</sup>	6.68e <sup>-1</sup>	6.15e <sup>-1</sup>	6.27e <sup>-1</sup>	P.TST
	Rank	2.51e <sup>1</sup>	2.13e <sup>1</sup>	1.67e <sup>1</sup>	2.32e <sup>1</sup>	2.00e <sup>1</sup>	9.56	1.41e <sup>1</sup>	2.01e <sup>1</sup>	7.69	9.44	1.09e <sup>1</sup>	1.75e <sup>1</sup>	1.12e <sup>1</sup>	9.28	9.19	9.71	6.66	7.56	Moi.B
Healthcare	MASE	1.16	1.00	7.84e <sup>-1</sup>	9.51e <sup>-1</sup>	7.65e <sup>-1</sup>	6.60e <sup>-1</sup>	8.03e <sup>-1</sup>	6.91e <sup>-1</sup>	6.86e <sup>-1</sup>	7.74e <sup>-1</sup>	6.98e <sup>-1</sup>	7.49e <sup>-1</sup>	<u>6.07e<sup>-1</sup></u>	6.45e <sup>-1</sup>	<b>5.99e<sup>-1</sup></b>	9.51e <sup>-1</sup>	6.75e <sup>-1</sup>	6.91e <sup>-1</sup>	Chr.L
	CRPS	1.19	1.00	5.70e <sup>-1</sup>	8.03e <sup>-1</sup>	7.23e <sup>-1</sup>	5.12e <sup>-1</sup>	9.12e <sup>-1</sup>	7.13e <sup>-1</sup>	5.76e <sup>-1</sup>	6.28e <sup>-1</sup>	6.52e <sup>-1</sup>	6.81e <sup>-1</sup>	4.96e <sup>-1</sup>	<u>4.85e<sup>-1</sup></u>	<b>4.46e<sup>-1</sup></b>	7.72e <sup>-1</sup>	5.14e <sup>-1</sup>	5.28e <sup>-1</sup>	Chr.L
	Rank	2.26e <sup>1</sup>	1.98e <sup>1</sup>	9.60	1.52e <sup>1</sup>	1.26e <sup>1</sup>	9.40	1.74e <sup>1</sup>	1.72e <sup>1</sup>	1.06e <sup>1</sup>	1.26e <sup>1</sup>	9.60	1.60e <sup>1</sup>	7.00	6.00	4.60	1.63e <sup>1</sup>	5.80	7.20	Chr.L
Nature	MASE	9.62e <sup>-1</sup>	1.00	1.02	1.06	1.64	8.71e <sup>-1</sup>	1.37	9.33e <sup>-1</sup>	9.16e <sup>-1</sup>	8.51e <sup>-1</sup>	8.80e <sup>-1</sup>	8.60e <sup>-1</sup>	8.51e <sup>-1</sup>	8.23e <sup>-1</sup>	8.13e <sup>-1</sup>	7.97e <sup>-1</sup>	<u>7.80e<sup>-1</sup></u>	<b>7.56e<sup>-1</sup></b>	Moi.L
	CRPS	1.33	1.00	6.58e <sup>-1</sup>	9.10e <sup>-1</sup>	5.35e <sup>-1</sup>	3.48e <sup>-1</sup>	5.61e <sup>-1</sup>	5.32e <sup>-1</sup>	3.47e <sup>-1</sup>	3.42e <sup>-1</sup>	3.33e <sup>-1</sup>	4.06e <sup>-1</sup>	3.83e <sup>-1</sup>	3.66e <sup>-1</sup>	3.64e <sup>-1</sup>	3.73e <sup>-1</sup>	<u>3.15e<sup>-1</sup></u>	<b>3.11e<sup>-1</sup></b>	Moi.L
	Rank	2.75e <sup>1</sup>	2.67e <sup>1</sup>	2.11e <sup>1</sup>	2.37e <sup>1</sup>	1.73e <sup>1</sup>	1.05e <sup>1</sup>	1.95e <sup>1</sup>	2.13e <sup>1</sup>	1.03e <sup>1</sup>	8.93	8.13	1.65e <sup>1</sup>	1.40e <sup>1</sup>	1.29e <sup>1</sup>	1.23e <sup>1</sup>	9.21	5.20	5.27	Moi.B
Sales	MASE	1.00	1.00	8.13e <sup>-1</sup>	8.73e <sup>-1</sup>	7.07e <sup>-1</sup>	7.16e <sup>-1</sup>	9.81e <sup>-1</sup>	7.04e <sup>-1</sup>	<b>6.90e<sup>-1</sup></b>	6.99e <sup>-1</sup>	7.00e <sup>-1</sup>	8.17e <sup>-1</sup>	7.33e <sup>-1</sup>	7.26e <sup>-1</sup>	7.24e <sup>-1</sup>	7.31e <sup>-1</sup>	<u>6.95e<sup>-1</sup></u>	7.10e <sup>-1</sup>	P.TST
	CRPS	8.96e <sup>-1</sup>	1.00	4.58e <sup>-1</sup>	4.80e <sup>-1</sup>	3.52e <sup>-1</sup>	3.52e <sup>-1</sup>	4.84e <sup>-1</sup>	4.14e <sup>-1</sup>	3.48e <sup>-1</sup>	3.51e <sup>-1</sup>	<b>3.44e<sup>-1</sup></b>	4.92e <sup>-1</sup>	3.66e <sup>-1</sup>	3.63e <sup>-1</sup>	3.62e <sup>-1</sup>	3.61e <sup>-1</sup>	<u>3.47e<sup>-1</sup></u>	3.63e <sup>-1</sup>	T.FM
	Rank	2.80e <sup>1</sup>	2.80e <sup>1</sup>	1.98e <sup>1</sup>	2.10e <sup>1</sup>	8.75	1.10e <sup>1</sup>	2.05e <sup>1</sup>	1.45e <sup>1</sup>	5.00	7.00	3.00	2.15e <sup>1</sup>	1.22e <sup>1</sup>	1.05e <sup>1</sup>	1.00e <sup>1</sup>	1.00e <sup>1</sup>	3.25	6.75	T.FM
Transport	MASE	1.26	1.00	9.74e <sup>-1</sup>	1.08	7.45e <sup>-1</sup>	6.79e <sup>-1</sup>	7.90e <sup>-1</sup>	7.31e <sup>-1</sup>	7.09e <sup>-1</sup>	7.07e <sup>-1</sup>	7.41e <sup>-1</sup>	7.39e <sup>-1</sup>	7.37e <sup>-1</sup>	7.12e <sup>-1</sup>	7.14e <sup>-1</sup>	7.26e <sup>-1</sup>	<u>6.34e<sup>-1</sup></u>	<b>6.07e<sup>-1</sup></b>	Moi.L
	CRPS	2.07	1.00	7.63e <sup>-1</sup>	1.33	4.84e <sup>-1</sup>	4.43e <sup>-1</sup>	5.31e <sup>-1</sup>	5.93e <sup>-1</sup>	4.61e <sup>-1</sup>	4.60e <sup>-1</sup>	5.10e <sup>-1</sup>	6.01e <sup>-1</sup>	5.30e <sup>-1</sup>	5.12e <sup>-1</sup>	5.12e <sup>-1</sup>	4.98e <sup>-1</sup>	<u>4.12e<sup>-1</sup></u>	<b>3.93e<sup>-1</sup></b>	Moi.L
	Rank	2.84e <sup>1</sup>	2.43e <sup>1</sup>	2.18e <sup>1</sup>	2.61e <sup>1</sup>	8.73	6.60	1.39e <sup>1</sup>	1.74e <sup>1</sup>	8.07	7.93	1.06e <sup>1</sup>	1.81e <sup>1</sup>	1.39e <sup>1</sup>	1.08e <sup>1</sup>	1.11e <sup>1</sup>	1.07e <sup>1</sup>	5.40	5.67	Moi.B
Web/CloudOps	MASE	1.13	1.00	9.57e <sup>-1</sup>	5.21e <sup>-1</sup>	8.50e <sup>-1</sup>	6.62e <sup>-1</sup>	6.29e <sup>-1</sup>	5.33e <sup>-1</sup>	<b>4.82e<sup>-1</sup></b>	4.88e <sup>-1</sup>	1.42	2.72e <sup>-1</sup>	6.78e <sup>-1</sup>	6.76e <sup>-1</sup>	6.75e <sup>-1</sup>	7.73e <sup>-1</sup>	7.62e <sup>-1</sup>	6.91e <sup>-1</sup>	P.TST
	CRPS	1.07	1.00	9.04e <sup>-1</sup>	6.08e <sup>-1</sup>	6.33e <sup>-1</sup>	5.03e <sup>-1</sup>	5.68e <sup>-1</sup>	5.70e <sup>-1</sup>	4.37e <sup>-1</sup>	<u>4.54e<sup>-1</sup></u>	7.39e <sup>-1</sup>	6.03e <sup>-1</sup>	6.29e <sup>-1</sup>	6.51e <sup>-1</sup>	6.47e <sup>-1</sup>	6.49e <sup>-1</sup>	6.28e <sup>-1</sup>	6.10e <sup>-1</sup>	P.TST
	Rank	2.19e <sup>1</sup>	2.18e <sup>1</sup>	1.99e <sup>1</sup>	1.66e <sup>1</sup>	1.48e <sup>1</sup>	6.95	1.22e <sup>1</sup>	1.29e <sup>1</sup>	4.75	5.85	1.84e <sup>1</sup>	1.35e <sup>1</sup>	1.29e <sup>1</sup>	1.45e <sup>1</sup>	1.48e <sup>1</sup>	1.35e <sup>1</sup>	1.22e <sup>1</sup>	1.13e <sup>1</sup>	P.TST

sul et al., 2023a), **Timer** (Liu et al., 2024), **TTM** (Ekambaram et al., 2024)], **VisionTS** (Chen et al., 2024). These models all provide publicly accessible model parameters for direct use. However, it is important to note that pre-training datasets of **TimesFM**, **Chronos**, and **Moirai** exhibit partial data leakage issues for GIFT-Eval. To keep comparison across models fair, in the main paper we report results with public checkpoints for each model. However, since **Moirai** provides pretraining code, here we pretrain a series of **Moirai** models using GIFT-Eval’s pretraining split to demonstrate its utility. We empirically investigate the impact of data leakage in Appendix F.3. Further details on model-specific hyperparameters and tuning can be found in Appendix A.

For readability concerns, we omit results from **Auto\_ETS**, **DLinear** and **Crossformer** models in the main tables, however, the reader may refer to Appendix F for results with all models available. For the same space concerns, we use abbreviations to replace each model in the tables. Here is a list of model→abbreviation pairs for reference: Naive: **Nv.**, Seasonal Naive: **S.Nv.**, Auto\_Arima: **A.Ar.**, Auto\_Theta: **A.Th.**, Auto\_ETS: **A.ETS**, DeepAR: **D.AR**, TFT: **TFT**, TiDE: **TiDE**, N-BEATS: **N-B.**, PatchTST: **P.TST**, iTransformer: **iTr.**, DLinear: **DLin.**, Crossformer: **C.former**, **R**: [Lag-Llama: **L-Llama**, **Timer**: **Timer**, **TTM**: **TTM**], TimesFM: **T.FM**, VisionTS: **V.TS**, Chronos: **Chr.**, Chronos<sub>Small</sub>: **Chr.s**, Chronos<sub>Base</sub>: **Chr.B**, Chronos<sub>Large</sub>: **Chr.L**, Moirai: **Moi.**, Moirai<sub>Small</sub>: **Moi.s**, Moirai<sub>Base</sub>: **Moi.B**, Moirai<sub>Large</sub>: **Moi.L**.

**Evaluation setting** Performance is assessed using two metrics: the Mean Absolute Scaled Error (MASE) for point forecasts and the Continuous Ranked Probability Score (CRPS) (Gneiting & Raftery, 2007) for probabilistic forecasts (definition of both metrics are in Appendix C), see Appendix F.2 for results with more metrics. To standardize comparison across benchmarks, both metrics are normalized against the Seasonal Naive baseline. To avoid skew from any single dataset, we employ a ‘Rank’ metric that assigns a numerical ranking to each model across all 97 configurations judging by their CRPS score. The average of these ranks is then reported as the final Rank for each model.

## 4.1 RESULTS

We present results across five distinct parts. The first four parts aggregate the results by the key characteristics that guided the development of our benchmark: domain, prediction length, frequency, and number of variates, then conclude the section with aggregation of results across all configurations. For results on all datasets, frequency and prediction length combinations see Tables 22 to 24.

**Domain | Table 2** The results across various domains demonstrate that foundation models consistently outperform both statistical and deep learning models. Notably, the foundation models achieve top performance in most areas, except in the Web/CloudOps domain. As discussed in Section 3.2 Web/CloudOps is one of the domains to exhibit the highest lumpiness. This pattern suggests that foundation models may struggle with time series possessing such characteristics. In contrast, deep learning models like **PatchTST** and **iTransformer** excel in these challenging domains, possibly indicating a shortfall of the training data used for foundation models in these areas. The comparison of different foundation models yields inconsistent conclusions across various domains. We believe

Table 3: Results on GIFT-Eval aggregated by Prediction Length. The best results across each row are **bolded**, while the second best results are underlined.

Pred. Len.	Metric	Nv.	S.Nv.	A.Ar.	A.Th.	D.AR	TFT	TIDE	N-B.	P.TST	iTr.	T.FM	V.TS	Chr.s	Chr.b	Chr.l	Moi.s	Moi.b	Moi.l	Best
Long	MASE	1.40	1.00	9.85e <sup>-1</sup>	8.69e <sup>-1</sup>	1.10	5.89e <sup>-1</sup>	6.55e <sup>-1</sup>	6.44e <sup>-1</sup>	<u>5.37e<sup>-1</sup></u>	5.66e <sup>-1</sup>	9.90e <sup>-1</sup>	<b>5.22e<sup>-1</sup></b>	6.58e <sup>-1</sup>	6.34e <sup>-1</sup>	6.32e <sup>-1</sup>	6.44e <sup>-1</sup>	6.25e <sup>-1</sup>	6.04e <sup>-1</sup>	V.TS
	CRPS	1.89	1.00	8.05e <sup>-1</sup>	1.40	6.28e <sup>-1</sup>	<u>3.79e<sup>-1</sup></u>	4.48e <sup>-1</sup>	5.65e <sup>-1</sup>	<b>3.68e<sup>-1</sup></b>	3.91e <sup>-1</sup>	5.18e <sup>-1</sup>	4.56e <sup>-1</sup>	5.22e <sup>-1</sup>	5.04e <sup>-1</sup>	5.02e <sup>-1</sup>	4.45e <sup>-1</sup>	4.23e <sup>-1</sup>	4.22e <sup>-1</sup>	P.TST
	Rank	2.72e <sup>1</sup>	2.31e <sup>1</sup>	2.09e <sup>1</sup>	2.43e <sup>1</sup>	1.72e <sup>1</sup>	<u>6.48</u>	1.16e <sup>1</sup>	1.61e <sup>1</sup>	6.00	<b>6.00</b>	7.19	1.51e <sup>1</sup>	1.26e <sup>1</sup>	1.56e <sup>1</sup>	1.40e <sup>1</sup>	1.44e <sup>1</sup>	9.29	8.24	8.19
Medium	MASE	1.46	1.00	1.02	1.17	1.33	9.49e <sup>-1</sup>	9.86e <sup>-1</sup>	1.03	<u>8.56e<sup>-1</sup></u>	8.67e <sup>-1</sup>	1.44	<b>8.47e<sup>-1</sup></b>	1.04	1.04	1.03	1.03	1.03	0.92e <sup>-1</sup>	V.TS
	CRPS	1.87	1.00	8.33e <sup>-1</sup>	1.53	6.40e <sup>-1</sup>	<u>4.68e<sup>-1</sup></u>	4.68e <sup>-1</sup>	4.78e <sup>-1</sup>	<b>4.61e<sup>-1</sup></b>	4.70e <sup>-1</sup>	6.30e <sup>-1</sup>	5.83e <sup>-1</sup>	6.25e <sup>-1</sup>	6.30e <sup>-1</sup>	6.22e <sup>-1</sup>	5.53e <sup>-1</sup>	5.35e <sup>-1</sup>	5.23e <sup>-1</sup>	P.TST
	Rank	2.62e <sup>1</sup>	2.16e <sup>1</sup>	1.99e <sup>1</sup>	2.43e <sup>1</sup>	1.36e <sup>1</sup>	<u>5.90</u>	1.24e <sup>1</sup>	1.70e <sup>1</sup>	5.14	<b>5.14</b>	7.71	1.41e <sup>1</sup>	1.41e <sup>1</sup>	1.50e <sup>1</sup>	1.50e <sup>1</sup>	1.42e <sup>1</sup>	1.00e <sup>1</sup>	8.86	8.62
Short	MASE	1.14	1.00	9.35e <sup>-1</sup>	9.55e <sup>-1</sup>	1.20	8.83e <sup>-1</sup>	1.14	8.62e <sup>-1</sup>	8.32e <sup>-1</sup>	8.89e <sup>-1</sup>	8.23e <sup>-1</sup>	8.71e <sup>-1</sup>	7.79e <sup>-1</sup>	<u>7.68e<sup>-1</sup></u>	<b>7.61e<sup>-1</sup></b>	8.97e <sup>-1</sup>	8.19e <sup>-1</sup>	8.21e <sup>-1</sup>	Chr.l
	CRPS	1.09	1.00	7.35e <sup>-1</sup>	8.16e <sup>-1</sup>	7.95e <sup>-1</sup>	5.92e <sup>-1</sup>	7.95e <sup>-1</sup>	7.48e <sup>-1</sup>	5.71e <sup>-1</sup>	6.11e <sup>-1</sup>	5.77e <sup>-1</sup>	7.51e <sup>-1</sup>	5.52e <sup>-1</sup>	<u>5.42e<sup>-1</sup></u>	<b>5.38e<sup>-1</sup></b>	6.09e <sup>-1</sup>	5.48e <sup>-1</sup>	5.53e <sup>-1</sup>	Chr.l
	Rank	2.36e <sup>1</sup>	2.31e <sup>1</sup>	1.64e <sup>1</sup>	1.86e <sup>1</sup>	1.62e <sup>1</sup>	1.09e <sup>1</sup>	1.79e <sup>1</sup>	1.87e <sup>1</sup>	9.27	<b>9.27</b>	1.02e <sup>1</sup>	8.80	1.96e <sup>1</sup>	9.65	8.33	8.33	1.14e <sup>1</sup>	6.18	6.93

Table 4: Results on GIFT-Eval aggregated by frequency. The best results across each row are **bolded**, while second best results are underlined.

Freq.	Metric	Nv.	S.Nv.	A.Ar.	A.Th.	D.AR	TFT	TIDE	N-B.	P.TST	iTr.	T.FM	V.TS	Chr.s	Chr.b	Chr.l	Moi.s	Moi.b	Moi.l	Best
10S	MASE	1.98	1.00	1.00	<u>1.59e<sup>-1</sup></u>	3.76e <sup>-1</sup>	3.57e <sup>-1</sup>	3.23e <sup>-1</sup>	2.71e <sup>-1</sup>	2.24e <sup>-1</sup>	2.35e <sup>-1</sup>	7.87e <sup>-1</sup>	<u>2.16e<sup>-1</sup></u>	5.23e <sup>-1</sup>	5.23e <sup>-1</sup>	5.06e <sup>-1</sup>	7.95e <sup>-1</sup>	8.41e <sup>-1</sup>	5.72e <sup>-1</sup>	A.Th.
	CRPS	1.44	1.00	1.00	<b>3.15e<sup>-1</sup></b>	7.54e <sup>-1</sup>	6.72e <sup>-1</sup>	7.05e <sup>-1</sup>	5.98e <sup>-1</sup>	5.36e <sup>-1</sup>	5.10e <sup>-1</sup>	1.30	6.91e <sup>-1</sup>	7.93e <sup>-1</sup>	8.59e <sup>-1</sup>	8.18e <sup>-1</sup>	1.24	1.06	1.02	A.Th.
	Rank	1.93e <sup>1</sup>	1.13e <sup>1</sup>	1.03e <sup>1</sup>	1.00	1.23e <sup>1</sup>	8.83	1.12e <sup>1</sup>	7.17	5.00	<b>2.50</b>	2.53e <sup>1</sup>	1.08e <sup>1</sup>	1.12e <sup>1</sup>	1.33e <sup>1</sup>	1.23e <sup>1</sup>	2.26e <sup>1</sup>	1.95e <sup>1</sup>	1.78e <sup>1</sup>	A.Th.
5T	MASE	9.42e <sup>-1</sup>	1.00	1.00	9.84e <sup>-1</sup>	1.40	8.36e <sup>-1</sup>	9.61e <sup>-1</sup>	8.84e <sup>-1</sup>	7.87e <sup>-1</sup>	7.73e <sup>-1</sup>	2.38	8.19e <sup>-1</sup>	8.72e <sup>-1</sup>	8.62e <sup>-1</sup>	8.69e <sup>-1</sup>	7.39e <sup>-1</sup>	<u>6.89e<sup>-1</sup></u>	<b>6.69e<sup>-1</sup></b>	Moi.l
	CRPS	1.19	1.00	1.00	9.48e <sup>-1</sup>	7.49e <sup>-1</sup>	5.36e <sup>-1</sup>	6.31e <sup>-1</sup>	6.99e <sup>-1</sup>	5.22e <sup>-1</sup>	5.22e <sup>-1</sup>	6.73e <sup>-1</sup>	7.02e <sup>-1</sup>	6.82e <sup>-1</sup>	6.83e <sup>-1</sup>	6.87e <sup>-1</sup>	4.96e <sup>-1</sup>	<u>4.84e<sup>-1</sup></u>	<b>4.61e<sup>-1</sup></b>	Moi.l
	Rank	2.34e <sup>1</sup>	2.39e <sup>1</sup>	2.24e <sup>1</sup>	2.28e <sup>1</sup>	1.77e <sup>1</sup>	6.58	1.33e <sup>1</sup>	1.64e <sup>1</sup>	6.75	<b>7.75</b>	1.52e <sup>1</sup>	1.63e <sup>1</sup>	1.48e <sup>1</sup>	1.51e <sup>1</sup>	1.58e <sup>1</sup>	7.44	6.42	4.58	Moi.l
10T	MASE	1.28	1.00	1.00	1.62	1.55	<u>9.42e<sup>-1</sup></u>	1.27	1.21	1.19	1.09	1.27	<b>9.12e<sup>-1</sup></b>	1.20	1.09	1.07	1.00	1.15	1.13	V.TS
	CRPS	2.08	1.00	1.00	2.51	5.37e <sup>-1</sup>	<b>3.64e<sup>-1</sup></b>	5.68e <sup>-1</sup>	6.88e <sup>-1</sup>	<u>4.34e<sup>-1</sup></u>	4.43e <sup>-1</sup>	4.59e <sup>-1</sup>	4.42e <sup>-1</sup>	5.47e <sup>-1</sup>	4.75e <sup>-1</sup>	4.71e <sup>-1</sup>	4.91e <sup>-1</sup>	5.04e <sup>-1</sup>	5.14e <sup>-1</sup>	TFT
	Rank	2.67e <sup>1</sup>	2.22e <sup>1</sup>	2.12e <sup>1</sup>	2.80e <sup>1</sup>	1.47e <sup>1</sup>	3.67	1.65e <sup>1</sup>	1.82e <sup>1</sup>	9.50	<b>8.00</b>	1.00e <sup>1</sup>	9.33	1.55e <sup>1</sup>	1.05e <sup>1</sup>	9.67	1.10e <sup>1</sup>	1.28e <sup>1</sup>	1.39e <sup>1</sup>	TFT
15T	MASE	1.52	1.00	9.78e <sup>-1</sup>	1.03	1.76	9.66e <sup>-1</sup>	1.02	1.02	<u>8.77e<sup>-1</sup></u>	<u>8.28e<sup>-1</sup></u>	9.56e <sup>-1</sup>	9.05e <sup>-1</sup>	9.20e <sup>-1</sup>	8.87e <sup>-1</sup>	8.85e <sup>-1</sup>	9.49e <sup>-1</sup>	9.25e <sup>-1</sup>	9.77e <sup>-1</sup>	P.TST
	CRPS	2.20	1.00	9.52e <sup>-1</sup>	1.51	1.26	7.08e <sup>-1</sup>	7.92e <sup>-1</sup>	9.63e <sup>-1</sup>	<u>6.55e<sup>-1</sup></u>	<b>6.51e<sup>-1</sup></b>	7.68e <sup>-1</sup>	8.56e <sup>-1</sup>	7.73e <sup>-1</sup>	7.49e <sup>-1</sup>	7.46e <sup>-1</sup>	7.39e <sup>-1</sup>	6.91e <sup>-1</sup>	7.20e <sup>-1</sup>	iTr.
	Rank	2.73e <sup>1</sup>	2.03e <sup>1</sup>	1.91e <sup>1</sup>	2.38e <sup>1</sup>	1.97e <sup>1</sup>	8.67	1.37e <sup>1</sup>	2.00e <sup>1</sup>	<b>5.00</b>	4.67	1.07e <sup>1</sup>	1.73e <sup>1</sup>	1.29e <sup>1</sup>	1.08e <sup>1</sup>	1.06e <sup>1</sup>	9.00	6.17	9.58	iTr.
H	MASE	1.46	1.00	1.02	1.28	1.31	8.25e <sup>-1</sup>	9.59e <sup>-1</sup>	8.72e <sup>-1</sup>	7.74e <sup>-1</sup>	8.05e <sup>-1</sup>	8.24e <sup>-1</sup>	7.70e <sup>-1</sup>	7.73e <sup>-1</sup>	<u>7.63e<sup>-1</sup></u>	<u>7.63e<sup>-1</sup></u>	8.92e <sup>-1</sup>	7.78e <sup>-1</sup>	7.70e <sup>-1</sup>	Chr.b
	CRPS	1.67	1.00	7.43e <sup>-1</sup>	1.57	6.23e <sup>-1</sup>	4.28e <sup>-1</sup>	5.11e <sup>-1</sup>	6.00e <sup>-1</sup>	4.07e <sup>-1</sup>	4.24e <sup>-1</sup>	4.69e <sup>-1</sup>	5.25e <sup>-1</sup>	4.68e <sup>-1</sup>	4.62e <sup>-1</sup>	4.64e <sup>-1</sup>	5.13e <sup>-1</sup>	4.13e <sup>-1</sup>	4.07e <sup>-1</sup>	P.TST
	Rank	2.75e <sup>1</sup>	2.48e <sup>1</sup>	2.20e <sup>1</sup>	2.66e <sup>1</sup>	1.52e <sup>1</sup>	8.77	1.44e <sup>1</sup>	1.85e <sup>1</sup>	<b>6.97</b>	<b>8.32</b>	1.16e <sup>1</sup>	1.64e <sup>1</sup>	1.18e <sup>1</sup>	1.10e <sup>1</sup>	1.12e <sup>1</sup>	1.13e <sup>1</sup>	5.42	5.23	Moi.l
D	MASE	1.00	1.00	8.82e <sup>-1</sup>	9.36e <sup>-1</sup>	9.06e <sup>-1</sup>	7.25e <sup>-1</sup>	1.15	7.75e <sup>-1</sup>	7.49e <sup>-1</sup>	8.31e <sup>-1</sup>	7.46e <sup>-1</sup>	8.22e <sup>-1</sup>	7.37e <sup>-1</sup>	<u>7.14e<sup>-1</sup></u>	<u>7.10e<sup>-1</sup></u>	7.83e <sup>-1</sup>	7.47e <sup>-1</sup>	7.66e <sup>-1</sup>	Chr.b
	CRPS	7.94e <sup>-1</sup>	1.00	4.69e <sup>-1</sup>	5.43e <sup>-1</sup>	4.91e <sup>-1</sup>	<b>3.70e<sup>-1</sup></b>	6.10e <sup>-1</sup>	5.24e <sup>-1</sup>	3.92e <sup>-1</sup>	4.38e <sup>-1</sup>	4.13e <sup>-1</sup>	5.04e <sup>-1</sup>	3.97e <sup>-1</sup>	3.78e <sup>-1</sup>	<u>3.77e<sup>-1</sup></u>	3.97e <sup>-1</sup>	3.86e <sup>-1</sup>	3.96e <sup>-1</sup>	TFT
	Rank	2.48e <sup>1</sup>	2.67e <sup>1</sup>	1.45e <sup>1</sup>	1.91e <sup>1</sup>	1.49e <sup>1</sup>	8.87	1.45e <sup>1</sup>	1.82e <sup>1</sup>	1.94e <sup>1</sup>	9.73	1.18e <sup>1</sup>	7.47	1.97e <sup>1</sup>	1.14e <sup>1</sup>	9.20	9.07	9.10	7.13	8.27
W	MASE	1.00	1.00	9.46e <sup>-1</sup>	1.03	1.46	9.21e <sup>-1</sup>	1.29	1.08	9.29e <sup>-1</sup>	1.25	8.47e <sup>-1</sup>	1.04	<u>7.45e<sup>-1</sup></u>	7.62e <sup>-1</sup>	<b>7.37e<sup>-1</sup></b>	1.00	9.01e <sup>-1</sup>	9.31e <sup>-1</sup>	Chr.l
	CRPS	8.74e <sup>-1</sup>	1.00	7.31e <sup>-1</sup>	7.87e <sup>-1</sup>	9.94e <sup>-1</sup>	7.26e <sup>-1</sup>	9.56e <sup>-1</sup>	9.71e <sup>-1</sup>	6.66e <sup>-1</sup>	9.56e <sup>-1</sup>	6.02e <sup>-1</sup>	9.43e <sup>-1</sup>	5.36e <sup>-1</sup>	5.42e <sup>-1</sup>	<b>5.29e<sup>-1</sup></b>	6.95e <sup>-1</sup>	6.37e <sup>-1</sup>	6.34e <sup>-1</sup>	Chr.l
	Rank	1.81e <sup>1</sup>	2.20e <sup>1</sup>	1.32e <sup>1</sup>	1.60e <sup>1</sup>	1.69e <sup>1</sup>	1.44e <sup>1</sup>	1.70e <sup>1</sup>	2.00e <sup>1</sup>	1.02e <sup>1</sup>	1.62e <sup>1</sup>	6.12	2.10e <sup>1</sup>	6.75	<b>6.00</b>	5.62	1.12e <sup>1</sup>	6.88	6.88	Chr.l
M	MASE	1.20	1.00	<u>7.59e<sup>-1</sup></u>	9.32e <sup>-1</sup>	1.22	9.01e <sup>-1</sup>	1.10	8.51e <sup>-1</sup>	8.59e <sup>-1</sup>	9.07e <sup>-1</sup>	<u>8.00e<sup>-1</sup></u>	9.15e <sup>-1</sup>	8.27e <sup>-1</sup>	8.57e <sup>-1</sup>	8.12e <sup>-1</sup>	1.04	8.07e <sup>-1</sup>	8.17e <sup>-1</sup>	A.Ar.
	CRPS	1.52	1.00	7.50e <sup>-1</sup>	8.73e <sup>-1</sup>	1.03	8.40e <sup>-1</sup>	1.16	9.62e <sup>-1</sup>	8.32e <sup>-1</sup>	8.03e <sup>-1</sup>	<b>7.33e<sup>-1</sup></b>	1.03	8.18e <sup>-1</sup>	8.49e <sup>-1</sup>	8.07e <sup>-1</sup>	9.93e <sup>-1</sup>	<u>7.51e<sup>-1</sup></u>	7.75e <sup>-1</sup>	T.FM
	Rank	2.52e <sup>1</sup>	1.89e <sup>1</sup>	1.16e <sup>1</sup>	1.56e <sup>1</sup>	1.02e <sup>1</sup>	2.00e <sup>1</sup>	1.44e <sup>1</sup>	1.90e <sup>1</sup>	7.40	<b>4.80</b>	1.90e <sup>1</sup>	1.06e <sup>1</sup>	1.16e <sup>1</sup>	1.04e <sup>1</sup>	1.07e <sup>1</sup>	4.20	7.00	7.00	Moi.b
Q	MASE	9.25e <sup>-1</sup>	1.00	8.00e <sup>-1</sup>	7.44e <sup>-1</sup>	9.00e <sup>-1</sup>	8.12e <sup>-1</sup>	1.05	7.56e <sup>-1</sup>	8.25e <sup>-1</sup>	7.69e <sup>-1</sup>	8.75e <sup>-1</sup>	8.50e <sup>-1</sup>	7.75e <sup>-1</sup>	7.69e <sup>-1</sup>	7.69e <sup>-1</sup>	7.76e <sup>-1</sup>	<u>7.11e<sup>-1</sup></u>	7.11e <sup>-1</sup>	Moi.l
	CRPS	9.51e <sup>-1</sup>	1.00	8.23e <sup>-1</sup>	7.97e <sup>-1</sup>	8.41e <sup>-1</sup>	8.37e <sup>-1</sup>	1.02	9.72e <sup>-1</sup>	8.35e <sup>-1</sup>	7.97e <sup>-1</sup>	8.53e <sup>-1</sup>	1.05	8.46e <sup>-1</sup>	8.40e <sup>-1</sup>	8.40e <sup>-1</sup>	7.94e <sup>-1</sup>	<b>7.40e<sup>-1</sup></b>	7.40e <sup>-1</sup>	Moi.l
	Rank	1.80e <sup>1</sup>	2.00e <sup>1</sup>	9.00	6.00	1.40e <sup>1</sup>	1.10e <sup>1</sup>	2.10e <sup>1</sup>	1.90e <sup>1</sup>	1.00e <sup>1</sup>	<b>7.00</b>	1.60e <sup>1</sup>	2.20e <sup>1</sup>	1.50e <sup>1</sup>	1.20e <sup>1</sup>	1.30e <sup>1</sup>	4.50	1.00	2.00	Moi.b
A	MASE	1.00	1.00	9.35e <sup>-1</sup>	7.83e <sup>-1</sup>	8.56e <sup>-1</sup>	7.78e <sup>-1</sup>	1.26	7.93e <sup>-1</sup>	8.29e <sup>-1</sup>	8.49e <sup>-1</sup>	8.44e <sup>-1</sup>	9.65e <sup>-1</sup>	9.42e <sup>-1</sup>	9.17e <sup>-1</sup>	9.17e <sup>-1</sup>	<u>7.51e<sup>-1</sup></u>	7.58e <sup>-1</sup>	<b>7.49e<sup>-1</sup></b>	Moi.l
	CRPS	9.93e <sup>-1</sup>	1.00	9.42e <sup>-1</sup>	8.33e <sup>-1</sup>	8.19e <sup>-1</sup>	7.97e <sup>-1</sup>	1.12	9.71e <sup>-1</sup>	8.48e <sup>-1</sup>	8.48e <sup>-1</sup>	8.48e <sup>-1</sup>	1.15	1.01	9.78e <sup>-1</sup>	9.78e <sup>-1</sup>	7.64e <sup>-1</sup>	<u>7.62e<sup>-1</sup></u>	<b>7.57e<sup>-1</sup></b>	Moi.l
	Rank	1.90e <sup>1</sup>	2.00e <sup>1</sup>	1.40e <sup>1</sup>	1.00e <sup>1</sup>	8.00	6.00	2.20e <sup>1</sup>	1.60e <sup>1</sup>	1.20e <sup>1</sup>	1.10e <sup>1</sup>									



Table 5: Results on GIFT-Eval aggregated by number of variates. The best results across each row are **bolded**, while the second best results are underlined.

Num. Var.	Metric	Nv.	S.Nv.	A.Ar.	A.Th.	D.AR	TFT	TiDE	N-B.	P.TST	iTr.	T.FM	V.TS	Chr.s	Chr.B	Chr.L	Moi.s	Moi.B	Moi.L	Best
Multivariate	MASE	1.15	1.00	1.03	8.01e <sup>-1</sup>	1.50	8.40e <sup>-1</sup>	1.01	7.82e <sup>-1</sup>	<u>7.11e<sup>-1</sup></u>	7.37e <sup>-1</sup>	1.17	<b>6.95e<sup>-1</sup></b>	8.04e <sup>-1</sup>	7.94e <sup>-1</sup>	7.88e <sup>-1</sup>	8.44e <sup>-1</sup>	8.31e <sup>-1</sup>	8.11e <sup>-1</sup>	V.TS
	CRPS	1.26	1.00	8.37e <sup>-1</sup>	9.26e <sup>-1</sup>	8.02e <sup>-1</sup>	4.95e <sup>-1</sup>	6.59e <sup>-1</sup>	6.41e <sup>-1</sup>	<b>4.51e<sup>-1</sup></b>	4.78e <sup>-1</sup>	5.82e <sup>-1</sup>	5.85e <sup>-1</sup>	5.55e <sup>-1</sup>	5.55e <sup>-1</sup>	5.52e <sup>-1</sup>	5.44e <sup>-1</sup>	5.15e <sup>-1</sup>	5.25e <sup>-1</sup>	P.TST
	Rank	<u>2.40e<sup>-1</sup></u>	<u>2.26e<sup>-1</sup></u>	1.95e <sup>-1</sup>	2.08e <sup>-1</sup>	1.90e <sup>-1</sup>	8.95	1.55e <sup>-1</sup>	1.69e <sup>-1</sup>	6.56	<b>7.05</b>	1.37e <sup>-1</sup>	1.53e <sup>-1</sup>	1.24e <sup>-1</sup>	1.23e <sup>-1</sup>	1.25e <sup>-1</sup>	9.94	8.63	8.91	P.TST
Univariate	MASE	1.36	1.00	9.12e <sup>-1</sup>	1.15	1.02	8.08e <sup>-1</sup>	9.59e <sup>-1</sup>	8.92e <sup>-1</sup>	8.05e <sup>-1</sup>	8.57e <sup>-1</sup>	8.29e <sup>-1</sup>	8.45e <sup>-1</sup>	7.97e <sup>-1</sup>	<u>7.80e<sup>-1</sup></u>	<b>7.75e<sup>-1</sup></b>	8.93e <sup>-1</sup>	7.90e <sup>-1</sup>	7.80e <sup>-1</sup>	Chr.L
	CRPS	1.49	1.00	7.21e <sup>-1</sup>	1.16	6.62e <sup>-1</sup>	5.24e <sup>-1</sup>	6.46e <sup>-1</sup>	7.30e <sup>-1</sup>	5.35e <sup>-1</sup>	5.64e <sup>-1</sup>	5.69e <sup>-1</sup>	6.83e <sup>-1</sup>	5.64e <sup>-1</sup>	5.47e <sup>-1</sup>	5.43e <sup>-1</sup>	5.98e <sup>-1</sup>	5.16e <sup>-1</sup>	5.08e <sup>-1</sup>	Chr.L
	Rank	2.56e <sup>-1</sup>	2.29e <sup>-1</sup>	1.70e <sup>-1</sup>	2.13e <sup>-1</sup>	1.34e <sup>-1</sup>	8.76	1.52e <sup>-1</sup>	1.85e <sup>-1</sup>	8.56	9.80	9.46	1.81e <sup>-1</sup>	1.19e <sup>-1</sup>	9.94	9.69	1.17e <sup>-1</sup>	6.07	<u>6.50</u>	<b>Moi.B</b>

Table 6: Results on GIFT-Eval aggregated by all results. The best results across each row are **bolded**, while the second best results are underlined.

Metric	Nv.	S.Nv.	A.Ar.	A.Th.	D.AR	TFT	TiDE	N-B.	P.TST	iTr.	T.FM	V.TS	Chr.s	Chr.B	Chr.L	Moi.s	Moi.B	Moi.L	Best
MASE	1.26	1.00	9.64e <sup>-1</sup>	9.78e <sup>-1</sup>	1.21	8.22e <sup>-1</sup>	9.80e <sup>-1</sup>	8.42e <sup>-1</sup>	<b>7.62e<sup>-1</sup></b>	8.02e <sup>-1</sup>	9.67e <sup>-1</sup>	<u>7.75e<sup>-1</sup></u>	8.00e <sup>-1</sup>	7.86e <sup>-1</sup>	7.81e <sup>-1</sup>	8.74e <sup>-1</sup>	8.11e <sup>-1</sup>	7.97e <sup>-1</sup>	P.TST
CRPS	1.38	1.00	7.70e <sup>-1</sup>	1.05	7.21e <sup>-1</sup>	<u>5.11e<sup>-1</sup></u>	6.52e <sup>-1</sup>	6.89e <sup>-1</sup>	<b>4.96e<sup>-1</sup></b>	5.24e <sup>-1</sup>	5.75e <sup>-1</sup>	6.38e <sup>-1</sup>	5.60e <sup>-1</sup>	5.51e <sup>-1</sup>	5.47e <sup>-1</sup>	5.76e <sup>-1</sup>	5.16e <sup>-1</sup>	5.15e <sup>-1</sup>	P.TST
Rank	2.49e <sup>-1</sup>	2.28e <sup>-1</sup>	1.81e <sup>-1</sup>	2.11e <sup>-1</sup>	1.59e <sup>-1</sup>	8.85	1.53e <sup>-1</sup>	1.78e <sup>-1</sup>	7.67	8.58	1.13e <sup>-1</sup>	1.69e <sup>-1</sup>	1.21e <sup>-1</sup>	1.10e <sup>-1</sup>	1.09e <sup>-1</sup>	1.10e <sup>-1</sup>	7.21	<b>7.57</b>	<b>Moi.B</b>

across all evaluated metrics. *Moirai* outperforms other foundation models, as it is the only model that supports multivariate forecasting. On the other hand, in univariate scenarios, foundation models, especially the large variant of *Moirai*, demonstrate superior performance over their deep learning counterparts. This suggests that foundation models, with their broader pretraining on diverse data sets, are particularly adept at extracting and leveraging predictive signals from single streams of data.

**General | Table 6** The comprehensive aggregation of results across the entire benchmark offers insightful performance distinctions. *PatchTST* emerges as the most dominant model for MASE and CRPS metrics, with *Moirai*<sub>Large</sub> securing the first place within the Rank metric. We also present the number of times each model achieves the best or second best results in Table 7. *Moirai*<sub>Large</sub> appears most frequently as the best performer, and as the model that appears in top 2 most frequently. The discrepancy between the RANK and MASE or CRPS metrics suggests that certain datasets may disproportionately influence the metric-based results, which is not captured by the ranking-based outcomes. Thus *PatchTST* offers reliable results across diverse datasets, making it a strong generalist. In contrast, *Moirai*<sub>Large</sub> delivers better performance on particular cases.

Some recent works (Shi et al., 2024a;b; Ansari et al., 2024) have verified the scaling law in time series foundation models (*i.e.*, larger model performs better), however, GIFT-Eval does not consistently support this conclusion.

## 4.2 QUALITATIVE RESULTS / FAILURE CASES

In addition to the quantitative results discussed earlier, we present qualitative analyses by sharing forecasting samples across various datasets using both deep learning and foundation models. For the deep learning models, we selected four representatives: *PatchTST* and *iTransformer*, from recent transformer-based architectures, and *DeepAR* and *N-BEATS*, which are more traditional deep learning approaches. Regarding foundational models, we included *Moirai* to represent encoder-decoder architectures, *Chronos* as a decoder-only model, and *VisionTS* due to its unique method of representing the time series through image modality. By examining how these models perform on different datasets, we aim to provide deeper insights into their forecasting behaviors, strengths, and limitations.

The plots in Figures 2(a) and 2(b) show forecasts by deep learning models on the multivariate *Bizitobs\_l2c* dataset (hourly, medium-term) and the univariate *Solar* dataset (ten-minutely, medium-term). In Figure 2(a) the irregular patterns challenge the models, with only *PatchTST* getting close to capturing some of the regular spikes accurately. *DeepAR* and *N-BEATS* perform reasonably but miss key periodic spikes, while *iTransformer*, despite its multivariate capability, oversimplifies the data into a sinusoidal pattern. In Figure 2(b), traditional models handle seasonal data better but

Table 7: Best and second best counts for each model across GIFT-Eval dataset configurations (97) according to the Rank metric. The best results across each row are **bolded**.

	Moi.L	Moi.B	P.TST	iTr.	C.Former	TFT	T.FM	Chr.L	Moi.s	Chr.B	A.Th.	D.AR	A.Ar.	A.ETS	Chr.s	N-B.	TiDE	Nv.	S.Nv.	DLin.	Timer	TTM	L-Llama	V.TS
Best	16	12	8	7	15	11	8	7	3	2	6	1	1	0	0	0	0	0	0	0	0	0	0	0
Second Best	14	14	13	13	2	3	3	4	7	7	0	5	3	3	3	2	1	0	0	0	0	0	0	0
Total	<b>30</b>	<u>26</u>	21	20	17	14	11	11	10	9	6	6	4	3	3	2	1	0	0	0	0	0	0	0

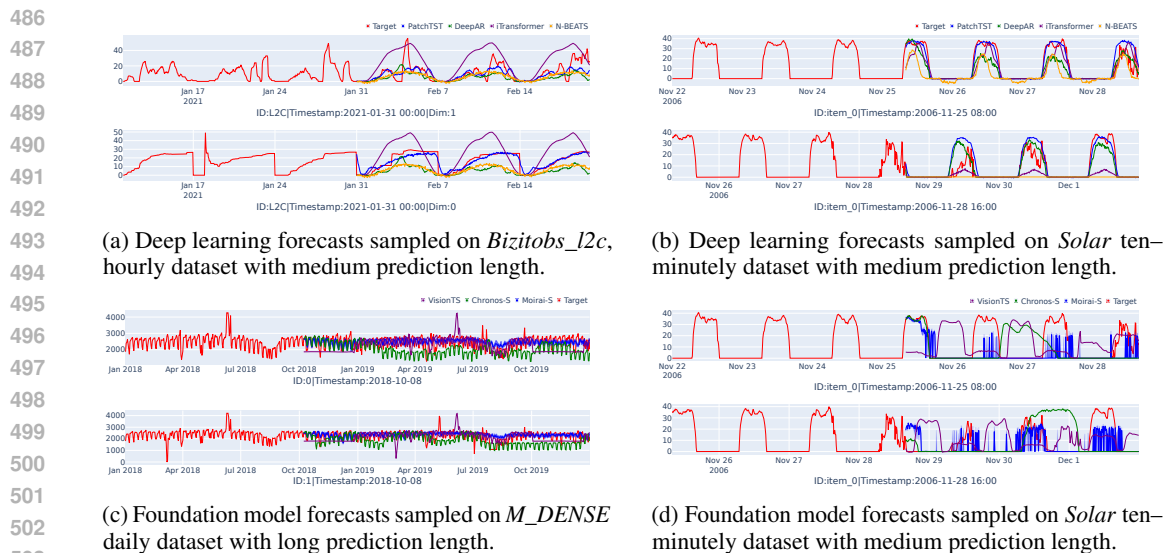


Figure 2: Qualitative plots showing forecasts from various deep learning and foundation models on several time series forecasting datasets.

still tend to underpredict, with N-BEATS producing a flat forecast in the second plot. PatchTST consistently outperforms others in both instances, showing robustness with both regular and irregular series, while iTransformer continues to underperform.

The plots in Figures 2(c) and 2(d) show forecasts by foundation models on two univariate datasets: *M\_DENSE* (daily, long-term) and *Solar* (ten-minutely, medium-term). Figure 2(c) displays varying performance among the foundation models. Chronos shows a clear degradation in performance as the prediction horizon extends, struggling to maintain accuracy over time, while VisionTS captures spikes but misaligns them. Moirai offers smoother, more conservative forecasts, which may result in less sensitivity to extreme events but provide more consistent alignment with the general trend. In Figure 2(d) VisionTS predicts seasonal peaks but with timing shifts. On the other hand, both Moirai and Chronos struggle to capture the well-spaced regularity of the data, missing key trends altogether. These poor results across all foundation models (see Figure 2(b) vs Figure 2(d)) mirror the quantitative findings in Section 4.1, *i.e.* deep learning models outperform foundation models at higher frequencies. For more qualitative examples see Appendix F.4

## 5 CONCLUSION

We introduce GIFT-Eval, a benchmark designed to evaluate time series forecasting models with diversity across four key characteristics: domain, frequency, number of variates, and prediction length. We ensure additional diversity by verifying six statistical features across temporal attributes, forecastability, and regularity. In addition to the train/test dataset, we also provide a pretraining dataset with no leakage into our evaluation set. With this, we aim to provide the necessary ground for fairly comparing different families of models, including foundation models, across a diverse benchmark. We conduct comprehensive experiments with R: [20 baselines] encompassing statistical, deep learning, and foundation models. Leveraging our detailed taxonomy, we provide insights into each model’s strengths relative to different characteristics. We also conduct a qualitative analysis highlighting failure cases in both deep learning and foundation models. GIFT-Eval is a comprehensive benchmark with fine-grained taxonomy that we hope will accelerate the development of new foundation time series forecasting models.

## REFERENCES

Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM*

- 540 *SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2019.  
541
- 542 Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan  
543 Gasthaus, Tim Januschowski, Danielle C. Maddix, Syama Rangapuram, David Salinas, Jasper  
544 Schulz, Lorenzo Stella, Ali Caner Turkmen, and Yuyang Wang. Gluonts: Probabilistic and neural  
545 time series modeling in python. *Journal of Machine Learning Research*, 21(116):1–6, 2020a. URL  
546 <http://jmlr.org/papers/v21/19-820.html>.
- 547 Alexander Alexandrov, Konstantinos Benidis, Michael Bohlke-Schneider, Valentin Flunkert, Jan  
548 Gasthaus, Tim Januschowski, Danielle C. Maddix, Syama Rangapuram, David Salinas, Jasper  
549 Schulz, Lorenzo Stella, Ali Caner Turkmen, and Yuyang Wang. GluonTS: Probabilistic and Neural  
550 Time Series Modeling in Python. *Journal of Machine Learning Research*, 21(116):1–6, 2020b.  
551 URL <http://jmlr.org/papers/v21/19-820.html>.
- 552 Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen,  
553 Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, et al.  
554 Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815*, 2024.  
555
- 556 George E. P. Box and David A. Pierce. Distribution of residual autocorrelations in autoregressive-  
557 integrated moving average time series models. *Journal of the American Statistical Association*, 65  
558 (332):1509–1526, 1970.
- 559 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé, Jared Kaplan, Harrison  
560 Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger,  
561 Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick  
562 Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter,  
563 Philippe Tillet, Felipe Petroski Such, David W. Cummings, Matthias Plappert, Fotios Chantzis,  
564 Elizabeth Barnes, Ariel Herbert-Voss, William H. Guss, Alex Nichol, Igor Babuschkin, Suchir  
565 Balaji, Shantanu Jain, Andrew Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa,  
566 Alec Radford, Matthew M. Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder,  
567 Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating  
568 large language models trained on code. *ArXiv*, abs/2107.03374, 2021. URL [https://api.  
569 semanticscholar.org/CorpusID:235755472](https://api.semanticscholar.org/CorpusID:235755472).
- 570 Mouxiang Chen, Lefei Shen, Zhuo Li, Xiaoyun Joy Wang, Jianling Sun, and Chenghao Liu. Visions:  
571 Visual masked autoencoders are free-lunch zero-shot time series forecasters. 2024. URL [https://api.  
572 semanticscholar.org/CorpusID:272310529](https://api.semanticscholar.org/CorpusID:272310529).
- 573 Abhimanyu Das, Weihao Kong, Andrew B. Leach, Shaan Mathur, Rajat Sen, and Rose Yu. Long-  
574 term forecasting with tide: Time-series dense encoder. *ArXiv*, abs/2304.08424, 2023a. URL  
575 <https://api.semanticscholar.org/CorpusID:258180439>.
- 576 Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation  
577 model for time-series forecasting. *ArXiv*, abs/2310.10688, 2023b. URL [https://api.  
578 semanticscholar.org/CorpusID:264172792](https://api.semanticscholar.org/CorpusID:264172792).
- 579 Vijay Ekambaram, Arindam Jati, Nam H. Nguyen, Pankaj Dayama, Chandra Reddy, Wesley M.  
580 Gifford, and Jayant Kalagnanam. Tiny time mixers (ttms): Fast pre-trained models for enhanced  
581 zero/few-shot forecasting of multivariate time series. *ArXiv*, abs/2401.03955, 2024. URL [https://api.  
582 semanticscholar.org/CorpusID:266844130](https://api.semanticscholar.org/CorpusID:266844130).
- 583 Patrick Emami, Abhijeet Sahu, and Peter Graf. Buildingsbench: A large-scale dataset of 900k  
584 buildings and benchmark for short-term load forecasting. In *Thirty-seventh Conference on  
585 Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=c5rqd6PZn6>.
- 586 Valentin Flunkert, David Salinas, and Jan Gasthaus. Deepar: Probabilistic forecasting with  
587 autoregressive recurrent networks. *ArXiv*, abs/1704.04110, 2017. URL [https://api.  
588 semanticscholar.org/CorpusID:12199225](https://api.semanticscholar.org/CorpusID:12199225).
- 589 Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and  
590 Marinka Zitnik. Units: A unified multi-task time series model. 2024. URL [https://api.  
591 semanticscholar.org/CorpusID:268201715](https://api.semanticscholar.org/CorpusID:268201715).
- 592  
593

- 594 F. Garza, M. M. Canseco, C. Challu, and K. G. Olivares. Statsforecast: Lightning fast forecasting  
595 with statistical and econometric models. Presented at PyCon Salt Lake City, Utah, US, 2022. URL:  
596 <https://github.com/Nixtla/statsforecast>.  
597
- 598 Federico Garza, Kin Gutierrez, Cristian Challu, Jose Moralez, Ricardo Olivares, and Max Mergen-  
599 thaler. tsfeatures: Time series feature extraction in python, 2024. URL <https://github.com/Nixtla/tsfeatures>. Accessed: 2024-09-24.  
600
- 601 Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation.  
602 *Journal of the American Statistical Association*, 102(477):359–378, 2007.  
603
- 604 Rakshitha Wathsadini Godahewa, Christoph Bergmeir, Geoffrey I. Webb, Rob Hyndman, and Pablo  
605 Montero-Manso. Monash time series forecasting archive. In *Thirty-fifth Conference on Neural  
606 Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL <https://openreview.net/forum?id=wEclmgAjU->.  
607
- 608 Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski.  
609 Moment: A family of open time-series foundation models. *ArXiv*, abs/2402.03885, 2024. URL  
610 <https://api.semanticscholar.org/CorpusID:267500205>.  
611
- 612 John Haslett and Adrian E. Raftery. Space-time modelling with long-memory dependence: Assessing  
613 ireland’s wind power resource. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*,  
614 38(1):1–50, 1989. doi: 10.2307/2347679. URL <https://doi.org/10.2307/2347679>.  
615
- 616 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and  
617 Jacob Steinhardt. Measuring massive multitask language understanding. *ArXiv*, abs/2009.03300,  
618 2020. URL <https://api.semanticscholar.org/CorpusID:221516475>.  
619
- 620 Addison Howard, Haruka Yui, Mark McDonald, and Will Cukierski. Recruit  
621 restaurant visitor forecasting. [https://kaggle.com/competitions/  
recruit-restaurant-visitor-forecasting](https://kaggle.com/competitions/recruit-restaurant-visitor-forecasting), 2017. Kaggle.  
622
- 623 Rob Hyndman, Anne B. Koehler, J. Keith Ord, and Ralph D. Snyder. *Forecasting with Exponential  
624 Smoothing: The State Space Approach*. Springer, 2008.
- 625 Rob J. Hyndman and George Athanasopoulos. *Forecasting: Principles and Practice*. OTexts, 2018.  
626
- 627 Guokun Lai, Wei-Cheng Chang, Yiming Yang, and Hanxiao Liu. Modeling long- and short-  
628 term temporal patterns with deep neural networks. *The 41st International ACM SIGIR Con-  
629 ference on Research & Development in Information Retrieval*, 2017. URL [https://api.  
630 semanticscholar.org/CorpusID:4922476](https://api.semanticscholar.org/CorpusID:4922476).
- 631 Bryan Lim, Sercan Ö. Arik, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers  
632 for interpretable multi-horizon time series forecasting. *ArXiv*, abs/1912.09363, 2019. URL  
633 <https://api.semanticscholar.org/CorpusID:209414891>.  
634
- 635 Xu Liu, Yutong Xia, Yuxuan Liang, Junfeng Hu, Yiwei Wang, Lei Bai, Chao Huang, Zhenguang  
636 Liu, Bryan Hooi, and Roger Zimmermann. Largest: A benchmark dataset for large-scale traffic  
637 forecasting. *arXiv preprint arXiv:2306.08259*, 2023a.
- 638 Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng Long.  
639 itransformer: Inverted transformers are effective for time series forecasting. *ArXiv*, abs/2310.06625,  
640 2023b. URL <https://api.semanticscholar.org/CorpusID:263830644>.  
641
- 642 Yong Liu, Haoran Zhang, Chenyu Li, Xiangdong Huang, Jianmin Wang, and Mingsheng Long. Timer:  
643 Generative pre-trained transformers are large time series models. In *International Conference  
644 on Machine Learning*, 2024. URL [https://api.semanticscholar.org/CorpusID:  
645 267412273](https://api.semanticscholar.org/CorpusID:267412273).
- 646 Markus Löning, A. Bagnall, Sajaysurya Ganesh, Viktor Kazakov, Jason Lines, and Franz J. Király.  
647 sktime: A unified interface for machine learning with time series. *ArXiv*, abs/1909.07872, 2019.  
URL <https://api.semanticscholar.org/CorpusID:202583848>.

- 648 Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. The m4 competition:  
649 Results, findings, conclusion and way forward. *International Journal of Forecasting*, 2018. URL  
650 <https://api.semanticscholar.org/CorpusID:158696437>.  
651
- 652 Paolo Mancuso, Veronica Piccialli, and Antonio Maria Sudoso. A machine learning approach  
653 for forecasting hierarchical time series. *Expert Syst. Appl.*, 182:115102, 2020. URL <https://api.semanticscholar.org/CorpusID:219177009>.  
654
- 655 Philipp Moritz, Robert Nishihara, Stephanie Wang, Alexey Tumanov, Richard Liaw, Eric Liang,  
656 William Paul, Michael I. Jordan, and Ion Stoica. Ray: A distributed framework for emerging ai  
657 applications. *ArXiv*, abs/1712.05889, 2017. URL <https://api.semanticscholar.org/CorpusID:34552495>.  
658
- 659 Soukayna Mouatadid, Paulo Orenstein, Genevieve Elaine Flaspohler, Miruna Oprescu, Judah Cohen,  
660 Franklyn Wang, Sean Edward Knight, Maria Geogdzhayeva, Samuel James Levang, Ernest  
661 Fraenkel, and Lester Mackey. SubseasonalclimateUSA: A dataset for subseasonal forecasting  
662 and benchmarking. In *Thirty-seventh Conference on Neural Information Processing Systems  
663 Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=pWkrU6raMt>.  
664
- 665 Tung Nguyen, Jason Kyle Jewik, Hritik Bansal, Prakhar Sharma, and Aditya Grover. Climatelearn:  
666 Benchmarking machine learning for weather and climate modeling. In *Thirty-seventh Conference  
667 on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=RZJEkLF1Px>.  
668
- 669 Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64  
670 words: Long-term forecasting with transformers. *ArXiv*, abs/2211.14730, 2022. URL <https://api.semanticscholar.org/CorpusID:254044221>.  
671
- 672 Boris N. Oreshkin, Dmitri Carpvov, Nicolas Chapados, and Yoshua Bengio. N-beats: Neural basis  
673 expansion analysis for interpretable time series forecasting. *ArXiv*, abs/1905.10437, 2019. URL  
674 <https://api.semanticscholar.org/CorpusID:166228758>.  
675
- 676 Santosh Palaskar, Vijay Ekambaram, Arindam Jati, Neelamadhav Gantayat, Avirup Saha, Seema  
677 Nagar, Nam Nguyen, Pankaj Dayama, Renuka Sindhgatta, Prateeti Mohapatra, Harshit Kumar,  
678 Jayant Kalagnanam, Nandyala Hemachandra, and Narayan Rangaraj. Automixer for improved  
679 multivariate time-series forecasting on business and it observability data. *Proceedings of the AAAI  
680 Conference on Artificial Intelligence*, 38:22962–22968, 03 2024. doi: 10.1609/aaai.v38i21.30336.  
681
- 682 Youngsuk Park, Danielle C. Maddix, Francois-Xavier Aubet, Kelvin K. Kan, Jan Gasthaus, and  
683 Yuyang Wang. Learning quantile functions without quantile crossing for distribution-free time  
684 series forecasting. In *International Conference on Artificial Intelligence and Statistics*, 2021. URL  
685 <https://api.semanticscholar.org/CorpusID:244103049>.  
686
- 687 Xiangfei Qiu, Jilin Hu, Lekui Zhou, Xingjian Wu, Junyang Du, Buang Zhang, Chenjuan Guo, Aoying  
688 Zhou, Christian S. Jensen, Zhenli Sheng, and Bin Yang. Tfb: Towards comprehensive and fair  
689 benchmarking of time series forecasting methods. *Proc. VLDB Endow.*, 17:2363–2377, 2024. URL  
690 <https://api.semanticscholar.org/CorpusID:268793935>.  
691
- 692 Kashif Rasul. PytorchTS, 2021. URL [https://github.com/zalandoresearch/  
693 pytorch-ts](https://github.com/zalandoresearch/pytorch-ts).
- 694 Kashif Rasul, Abdul-Saboor Sheikh, Ingmar Schuster, Urs M. Bergmann, and Roland Voll-  
695 graf. Multi-variate probabilistic time series forecasting via conditioned normalizing flows.  
696 *ArXiv*, abs/2002.06103, 2020. URL [https://api.semanticscholar.org/CorpusID:  
697 211126472](https://api.semanticscholar.org/CorpusID:211126472).  
698
- 699 Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. Autoregressive denoising  
700 diffusion models for multivariate probabilistic time series forecasting. In *International Conference  
701 on Machine Learning*, 2021. URL [https://api.semanticscholar.org/CorpusID:  
231719657](https://api.semanticscholar.org/CorpusID:231719657).

- 702 Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Arian Khorasani, George Adamopoulos,  
703 Rishika Bhagwatkar, Marin Bilovs, Hena Ghonia, Nadhir Vincent Hassen, Anderson Schneider,  
704 Sahil Garg, Alexandre Drouin, Nicolas Chapados, Yuriy Nevmyvaka, and Irina Rish. Lag-llama:  
705 Towards foundation models for probabilistic time series forecasting. 2023a. URL <https://api.semanticscholar.org/CorpusID:263909560>.  
706
- 707 Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Arian Khorasani, George Adamopoulos,  
708 Rishika Bhagwatkar, Marin Bilovs, Hena Ghonia, Nadhir Vincent Hassen, Anderson Schneider,  
709 Sahil Garg, Alexandre Drouin, Nicolas Chapados, Yuriy Nevmyvaka, and Irina Rish. Lag-llama:  
710 Towards foundation models for probabilistic time series forecasting. 2023b. URL <https://api.semanticscholar.org/CorpusID:263909560>.  
711
- 712 Neal Richardson, Ian Cook, Neal Crane, Dewey Dunnington, Romain François, Jonathan Keane,  
713 Diana Moldovan-Grunfeld, Jeroen Ooms, Julia Wujciak-Jens, and Apache Arrow. arrow: Integra-  
714 tion to 'apache' 'arrow', 2023. URL <https://github.com/apache/arrow/>. R package  
715 version 14.0.2, <https://arrow.apache.org/docs/r/>.  
716
- 717 Zezhi Shao, Fei Wang, Yongjun Xu, Wei Wei, Chengqing Yu, Zhao Zhang, Di Yao, Guangyin Jin, Xin  
718 Cao, Gao Cong, Christian S. Jensen, and Xueqi Cheng. Exploring progress in multivariate time se-  
719 ries forecasting: Comprehensive benchmarking and heterogeneity analysis. *ArXiv*, abs/2310.06119,  
720 2023. URL <https://api.semanticscholar.org/CorpusID:263829289>.  
721
- 722 Siqi Shen, Vincent Beek, and Alexandru Iosup. Statistical characterization of business-critical  
723 workloads hosted in cloud datacenters. *Proceedings - 2015 IEEE/ACM 15th International Sym-*  
724 *posium on Cluster, Cloud, and Grid Computing, CCGrid 2015*, pp. 465–474, 07 2015. doi:  
725 10.1109/CCGrid.2015.60.
- 726 Jingzhe Shi, Qinwei Ma, Huan Ma, and Lei Li. Scaling law for time series forecasting.  
727 *ArXiv*, abs/2405.15124, 2024a. URL <https://api.semanticscholar.org/CorpusID:270045141>.  
728
- 729 Xiaoming Shi, Shiyu Wang, Yuqi Nie, Dianqi Li, Zhou Ye, Qingsong Wen, and Ming Jin. Time-  
730 moe: Billion-scale time series foundation models with mixture of experts. 2024b. URL <https://api.semanticscholar.org/CorpusID:272832214>.  
731
- 732 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam  
733 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska,  
734 Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W.  
735 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda  
736 Askill, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders An-  
737 dreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmuller, Andrew M. Dai, Andrew La,  
738 Andrew Kyle Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta,  
739 Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul  
740 Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat,  
741 Aykut Erdem, Ayla Karakacs, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bo-  
742 janowski, Batuhan Ozyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno  
743 Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Stephen Howald, Bryan Orinion, Cameron Diao,  
744 Cameron Dour, Catherine Stinson, Cedrick Argueta, C'esar Ferri Ram'irez, Chandan Singh,  
745 Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites,  
746 Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera,  
747 Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Daniel H  
748 Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy,  
749 Daniel Mosegu'i Gonz'alez, Danielle R. Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito,  
750 Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis  
751 Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuweke Hupkes, Diganta Misra,  
752 Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova,  
753 Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth P. Donoway, Ellie  
754 Pavlick, Emanuele Rodolà, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A.  
755 Chi, Ethan Dyer, Ethan J. Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii,  
Fanyue Xia, Fatemeh Siar, Fernando Mart'inez-Plumed, Francesca Happ'e, François Chollet,

- 756 Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Gi-  
757 ambattista Parascandolo, Giorgio Mariani, Gloria Xinyue Wang, Gonzalo Jaimovitch-L’opez,  
758 Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi,  
759 Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schutze, Hiromu Yakura, Hongming Zhang,  
760 Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, John Kernion, Jacob Hilton,  
761 Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou,  
762 Jan Koco’ n, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Narain  
763 Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher,  
764 Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Oluwadara Alabi, Jiacheng Xu,  
765 Jiaming Song, Jillian Tang, Jane W Waweru, John Burden, John Miller, John U. Balis, Jonathan  
766 Batchelder, Jonathan Berant, Jorg Frohberg, Jos Rozen, José Hernández-Orallo, Joseph Boude-  
767 man, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil  
768 Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert,  
769 Kautubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Wallace Mathewson, Kristen Chiafullo,  
770 Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao,  
771 Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca  
772 Moschella, Luca Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Col’ on, Luke Metz,  
773 Lutfi Kerem cSenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal  
774 Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ram’ irez  
775 Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leav-  
776 itt, Matthias Hagen, Mátyás Schubert, Medina Baitemirova, Melody Arnaud, Melvin McElrath,  
777 Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube,  
778 Michal Swkedrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimeo  
779 Xu, Mirac Suzgun, Mitch Walker, Monica Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva,  
780 Mozhdah Gheini, T MukundVarma, Nanyun Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari  
781 Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia,  
782 Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha Iyer, Noah Constant, Noah  
783 Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo  
784 Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoor-  
785 molabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, P Hwang,  
786 P. Milkowski, Piyush S. Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen,  
787 Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphael Mil-  
788 liere, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert  
789 Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan Le Bras, Rosanne Liu, Rowan Jacobs,  
790 Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang,  
791 Sahib Singh, Saif Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel  
792 Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A.  
793 Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebas-  
794 tian Schuster, Sepideh Sadeghi, Shadi S. Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi,  
795 Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam  
796 Upadhyay, Shyamolima Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy,  
797 Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Bradley Torene, Sriharsha Hatwar, Stanislas  
798 Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T  
799 Piantadosi, Stuart M. Shieber, Summer Mishserghi, Svetlana Kiritchenko, Swaroop Mishra, Tal  
800 Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsunori Hashimoto, Te-Lin Wu, Théo Desbordes,  
801 Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev,  
802 Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz,  
803 Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Venkatesh  
804 Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William  
805 Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu,  
806 Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi,  
807 Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yu Hou, Yuntao Bai, Zachary  
808 Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation  
809 game: Quantifying and extrapolating the capabilities of language models. *ArXiv*, abs/2206.04615,  
2022. URL <https://api.semanticscholar.org/CorpusID:263625818>.
- 808 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csd: Conditional score-based  
809 diffusion models for probabilistic time series imputation. In *Neural Information Processing  
Systems*, 2021. URL <https://api.semanticscholar.org/CorpusID:235765577>.

- 810 Sean J. Taylor and Benjamin Letham. Forecasting at scale. *The American Statistician*, 72(1):  
811 37–45, 2018. doi: 10.1080/00031305.2017.1380080. URL [https://doi.org/10.1080/  
812 00031305.2017.1380080](https://doi.org/10.1080/00031305.2017.1380080).
- 813 Artur Trindade. ElectricityLoadDiagrams20112014. UCI Machine Learning Repository, 2015. DOI:  
814 <https://doi.org/10.24432/C58C86>.
- 815 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.  
816 Glue: A multi-task benchmark and analysis platform for natural language understand-  
817 ing. In *BlackboxNLP@EMNLP*, 2018. URL [https://api.semanticscholar.org/  
818 CorpusID:5034059](https://api.semanticscholar.org/CorpusID:5034059).
- 819 Jingyuan Wang, Jiawei Jiang, Wenjun Jiang, Chengkai Han, and Wayne Xin Zhao. Towards effi-  
820 cient and comprehensive urban spatial-temporal prediction: A unified library and performance  
821 benchmark. *arXiv preprint arXiv:2304.14343*, 2023a.
- 822 Zhixian Wang, Qingsong Wen, Chaoli Zhang, Liang Sun, Leandro Von Krannichfeldt, and Yi Wang.  
823 Benchmarks and custom package for electrical load forecasting. *arXiv preprint arXiv:2307.07191*,  
824 2023b.
- 825 Gerald Woo, Chenghao Liu, Akshat Kumar, and Doyen Sahoo. Pushing the limits of pre-training for  
826 time series forecasting in the cloudops domain. *arXiv preprint arXiv:2310.05063*, 2023.
- 827 Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo.  
828 Unified training of universal time series forecasting transformers. *ArXiv*, abs/2402.02592, 2024.  
829 URL <https://api.semanticscholar.org/CorpusID:267411817>.
- 830 Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers  
831 with auto-correlation for long-term series forecasting. In *Neural Information Processing Systems*,  
832 2021. URL <https://api.semanticscholar.org/CorpusID:235623791>.
- 833 Haixu Wu, Teng Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet:  
834 Temporal 2d-variation modeling for general time series analysis. *ArXiv*, abs/2210.02186, 2022.  
835 URL <https://api.semanticscholar.org/CorpusID:252715491>.
- 836 Ailing Zeng, Mu-Hwa Chen, L. Zhang, and Qiang Xu. Are transformers effective for time  
837 series forecasting? In *AAAI Conference on Artificial Intelligence*, 2022. URL <https://api.semanticscholar.org/CorpusID:249097444>.
- 838 Jiawen Zhang, Xumeng Wen, Zhenwei Zhang, Shun Zheng, Jia Li, and Jiang Bian. Probts: Bench-  
839 marking point and distributional forecasting across diverse prediction horizons. 2023. URL  
840 <https://api.semanticscholar.org/CorpusID:270559735>.
- 841 Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency  
842 for multivariate time series forecasting. In *The Eleventh International Conference on Learning  
843 Representations*, 2023. URL <https://openreview.net/forum?id=vSVLM2j9eie>.
- 844 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wan  
845 Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting.  
846 *ArXiv*, abs/2012.07436, 2020. URL [https://api.semanticscholar.org/CorpusID:  
847 229156802](https://api.semanticscholar.org/CorpusID:229156802).
- 848 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.  
849 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings  
850 of the AAAI conference on artificial intelligence*, volume 35, pp. 11106–11115, 2021.

## 859 A EXPERIMENTAL SETUP DETAILS

860 **Statistical models** We utilize the statsforecast (Garza et al., 2022) library to implement all five  
861 statistical baselines: Naive, Seasonal Naive, Auto\_Theta, Auto\_ETS and Auto\_Arima.  
862 Inference is performed on a CPU server equipped with 96 cores. For each dataset, a time limit of  
863



one day is set for the statistical model to complete its run. For any model that times out we halt it and replace its results with those from the `Seasonal Naive` model as a fallback. Given that some datasets in our benchmark are particularly long, we impose a maximum size constraint on each statistical baseline (set to 1000 with our time constraints), truncating the time series to this maximum size.

Table 8: Hyperparameter search range for deep learning baselines.

	TIDE						
Parameters	num_layers_encoder	num_layers_decoder	hidden_dim	temporal_hidden_dim	decoder_output_dim	dropout_rate	lr
Search Range	[1,2]	[1,2]	[256,512,1024]	[64,128]	[8,16,32]	[0.0, 0.5]	[1e-5:1e-1]
	N-BEATS			PatchTST			
Parameters	loss_function	hidden_layer_units	share_weights_in_stack	nb_blocks_per_stack	lr	d_model	num_encoder_layers
Search Range	["mase", "mape", "smape"]	[256, 512, 1024, 2048]	[True, False]	[3, 4]	[1e-5:1e-1]	[128, 256, 512]	[2, 3, 4]
	iTransformer			DeepAR		DLinear	
Parameters	d_model	num_encoder_layers	lr	hidden_size	num_layers	lr	lr
Search Range	[128, 256, 512]	[2, 3, 4]	[1e-5:1e-1]	[20,25,...,80]	[1,2,3,4]	[1e-5:1e-1]	[1e-5:1e-1]
	Crossformer			TF2			
Parameters	d_model	n_heads	lr	num_heads	hidden_dim	lr	
Search Range	[64,128,256]	[2,4,8]	[1e-5:5e-3]	[2,4,8]	[16,32,64]	[1e-5:1e-1]	

**Deep learning models** For all deeplearning models we either used models readily available in gluonts library (Alexandrov et al., 2020b) or we write our own wrappers. Where feasible we also add a probabilistic forecasting head to the models. Where direct probabilistic outputs are not feasible, we generate probabilistic evaluations by converting point forecasts into sample forecasts using a single sample. To identify the optimal hyperparameters, we conducted a comprehensive search across all 97 runs included in GIFT-Eval. We employed the ray library (Moritz et al., 2017) to parallelize the search on a single GPU and used the optuna (Akiba et al., 2019) library to extend this parallelization across multiple GPU servers. We search for 15 trials for each deep learning model per each of the 97 runs. Table 8 lists the range of parameters we search for each model. On top of the listed parameters for each model, we also search for weight decay on all runs in the range:  $[1e - 8 : 1e - 2]$ , **R**: [and for context length in range  $[1, 2, 4, 8] \times prediction\_length$ .] For the Crossformer model on the long term setting of *Jena Weather* dataset with both ten-minutely and hourly frequencies, we had to limit the search for `d_model` and `n_heads`, fixing them at 32 and 1, respectively. This adjustment was necessary because the model’s attention mechanism operates across multiple variates, leading to an OOM (Out of Memory) error due to the high number of variates present in this dataset.

**Foundation models** For all foundation models we use their public versions available online and conduct zero-shot evaluation on our benchmark’s test-split. Since Moirai (Woo et al., 2024) provides multi-patch size projections and varying context lengths. We adopt the similar approach by defining a frequency-to-patch size mapping as follows:

- Yearly, Quarterly: 8
- Monthly: 8
- Weekly, Daily: 16
- Hourly: 32
- Minute-level: 32
- Second-level: 64

**R**: [We set context length to 4000]. We used the public available Moirai models from the corresponding HuggingFace repos, i.e., `MoiraiSmall` - <https://huggingface.co/Salesforce/moirai-1.1-R-small>, `MoiraiBase` - <https://huggingface.co/Salesforce/moirai-1.1-R-base>, `MoiraiLarge` - <https://huggingface.co/Salesforce/moirai-1.1-R-large>.

For Chronos, we mainly follow their official implementation<sup>3</sup> for evaluation: with the number of samples as 20. The models are loaded from the corresponding HuggingFace repos, e.g., `ChronosTiny` - <https://huggingface.co/amazon/chronos-t5-tiny>, `ChronosSmall` - <https://huggingface.co/amazon/chronos-t5-small>, `ChronosBase` - <https://huggingface.co/amazon/chronos-t5-base>.

<sup>3</sup><https://github.com/amazon-science/chronos-forecasting/blob/main/scripts/evaluation/evaluate.py>

For TimesFM, we follow their official implementation<sup>4</sup> for evaluation. We set the context length for evaluation as 512 as mentioned in their paper since the maximum context length in training is 512. Following their default setting in their example, we keep the input patch length as 32, the output patch length as 128, the number of layers as 20, and the model dimension as 1280. TimesFM comes with only one model size, i.e., timesfm-1.0-200m, and we load the model from <https://huggingface.co/google/timesfm-1.0-200m>.

For VisionTS, we follow their official implementation<sup>5</sup> for evaluation. We set the context length as 2000, the norm constant as 0.4, the alignment constant as 0.4 according to their default settings. We use their implementation for seasonality detection to generate a candidate list and search an optimal seasonality parameter with the validation data.

**Additional parameters and computational resources.** All experiments are conducted on eight NVIDIA A100 GPUs. For models that has gone through training the loss function and optimizer are set following their original implementation. Additionally we set the batch size to 128 and, number of batches per epoch to 100, and finally number of epochs to 50.

## B DETAILS OF TIME SERIES FEATURES

This section gives a detailed view of the time series features we used to analyze the test portion of our data in Section 3.2. We use tsfeatures library (Garza et al., 2024) to calculate each metric. Given the scale of our dataset, we limit each time series history to the most recent 500 data points before computing the respective features. The prediction length remains faithful to the original values specified for the dataset and is not clipped. **R: [ Table 9 shows specific time series features of each dataset where we computed specific we classified them based on whether each feature (e.g., trend, seasonality, entropy) was lower or higher than the median value across all datasets].**

We also acknowledge that for some overly short time series, tsfeatures may output NaN (Not a Number) values for certain features—for example, the seasonal strength of some yearly time series data. In such cases, we exclude these NaN values during aggregation. Below we provide specific details for each feature used:

**Trend** Using the STL (Seasonal and Trend decomposition using Loess) method, a time series  $x_t$  is decomposed into a trend component  $f_t$ , multiple seasonal components  $s_{i,t}$  for  $i = 1, \dots, M$ , and a remainder component  $e_t$ :

$$x_t = f_t + s_{1,t} + \dots + s_{M,t} + e_t,$$

where  $M$  is the number of seasonal periods. The strength of the trend is quantified by comparing the variance of the remainder component  $e_t$  to the combined variance of the trend and remainder components. Specifically, the strength of the trend is defined as:

$$\text{trend} = 1 - \frac{\text{Var}(e_t)}{\text{Var}(f_t + e_t)}.$$

If the calculated value of trend is less than 0, it is set to 0; if it is greater than 1, it is set to 1. This measure indicates the proportion of the variability in the time series that is explained by the trend component, with values closer to 1 signifying a stronger trend.

**Seasonal Strength** Following the same decomposition above the strength of each seasonal component is quantified by comparing the variance of the remainder  $e_t$  to the combined variance of the seasonal component  $s_{i,t}$  and the remainder.

For each seasonal component  $s_{i,t}$ , the strength of seasonality is defined as:

<sup>4</sup>[https://github.com/google-research/timesfm/blob/master/experiments/long\\_horizon\\_benchmarks/run\\_eval.py](https://github.com/google-research/timesfm/blob/master/experiments/long_horizon_benchmarks/run_eval.py)

<sup>5</sup>[https://github.com/Keytoyce/VisionTS/blob/main/eval\\_gluonts/run.py](https://github.com/Keytoyce/VisionTS/blob/main/eval_gluonts/run.py)

Table 9: R: [Time Series features classification across all datasets in GIFT-Eval.]

	dataset	frequency	trend	seas. str.	entropy	hurst	lumpiness	stability
972								
973								
974								
975								
976	m4_yearly	A	high	low	high	high	low	high
977	bitbrains_fast_storage	5T	high	high	low	high	high	low
978	bitbrains_fast_storage	H	low	low	high	low	high	low
979	bitbrains_rnd	5T	high	high	low	low	low	low
980	bitbrains_rnd	H	high	high	low	low	high	low
981	bizitobs_application	10S	high	low	low	high	low	high
982	bizitobs_l2c	5T	high	low	low	low	low	high
983	bizitobs_l2c	H	low	low	high	low	high	high
984	bizitobs_service	10S	low	low	high	low	low	high
985	car_parts	M	low	low	high	low	high	low
986	covid_deaths	D	high	low	low	high	low	high
987	electricity	15T	low	high	low	high	low	low
988	electricity	D	high	high	low	high	low	high
989	electricity	H	high	high	low	high	low	low
990	electricity	W	high	high	low	low	low	high
991	ett1	15T	low	low	high	low	low	high
992	ett1	D	low	low	high	low	high	high
993	ett1	H	low	high	high	low	high	low
994	ett1	W	low	low	high	low	high	high
995	ett2	15T	high	low	low	high	low	high
996	ett2	D	high	low	high	low	high	high
997	ett2	H	high	low	low	high	low	high
998	ett2	W	high	low	high	low	high	high
999	hierarchical_sales	D	high	high	low	low	low	low
1000	hierarchical_sales	W	low	low	high	low	high	low
1001	hospital	M	low	low	high	low	low	low
1002	jena_weather	10T	high	high	low	low	low	low
1003	jena_weather	D	low	low	high	high	high	high
1004	jena_weather	H	high	high	low	low	low	low
1005	kdd_cup_2018	D	high	high	low	low	high	low
1006	kdd_cup_2018	H	high	high	low	low	low	low
1007	loop_seattle	5T	low	low	high	low	high	low
1008	loop_seattle	D	low	high	high	low	high	low
1009	loop_seattle	H	low	high	high	low	high	low
1010	m_dense	D	low	high	high	low	low	low
1011	m_dense	H	low	high	high	low	low	low
1012	m4_daily	D	high	low	low	high	low	high
1013	m4_hourly	H	low	high	low	low	low	low
1014	m4_monthly	M	low	low	high	high	low	high
1015	m4_quarterly	Q	high	low	high	high	low	high
1016	m4_weekly	W	high	low	high	high	high	high
1017	restaurant	D	high	high	low	low	high	low
1018	saugeen	D	high	low	high	low	high	high
1019	saugeen	M	low	high	high	low	high	low
1020	saugeen	W	low	low	high	low	high	high
1021	solar	10T	low	low	low	low	high	low
1022	solar	D	low	low	high	high	high	low
1023	solar	H	low	high	low	low	low	low
1024	solar	W	low	low	high	high	low	high
1025	sz_taxi	15T	low	low	high	high	high	low
	sz_taxi	H	low	low	high	high	high	low
	temperature_rain	D	high	high	low	high	high	high
	us_births	D	high	high	high	high	low	low
	us_births	M	high	high	high	high	low	high
	us_births	W	high	low	low	high	low	high

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

$$\text{seasonal\_strength}_i = 1 - \frac{\text{Var}(e_t)}{\text{Var}(s_{i,t} + e_t)}.$$

If the calculated value of  $\text{seasonal\_strength}_i$  is less than 0, it is set to 0; if it is greater than 1, it is set to 1. For non-seasonal time series,  $\text{seasonal\_strength} = 0$ . This measure indicates the proportion of the variability in the time series that is explained by the  $i$ -th seasonal component, with values closer to 1 signifying stronger seasonality for that component.

**Entropy** Entropy is defined as the Shannon entropy of the normalized spectral density estimate  $\hat{f}(\lambda)$ :

$$\text{Entropy} = - \int_{-\pi}^{\pi} \hat{f}(\lambda) \log \hat{f}(\lambda) d\lambda,$$

where  $\hat{f}(\lambda)$  is an estimate of the spectral density of the data. A lower spectral entropy indicates a higher signal-to-noise ratio, meaning the time series has more predictable patterns and is easier to forecast. Conversely, a higher spectral entropy suggests that the series is more complex and difficult to predict.

**Hurst Exponent** The *Hurst exponent* (*hurst*) is computed as 0.5 plus the maximum likelihood estimate of the fractional differencing order  $d$  by Haslett & Raftery (1989). The addition of 0.5 ensures consistency with the traditional Hurst coefficient. The values of the Hurst exponent vary between 0 and 1, with higher values indicating a smoother trend, less volatility, and less roughness.

**Stability** Stability measures the variability of the mean values across all tiles. It is defined as the variance of the means of the tiled windows. If the time series is divided into  $N$  tiles and  $\bar{x}_i$  represents the mean of the  $i$ -th tile, then the stability is calculated as:

$$\text{Stability} = \text{Var}(\bar{x}_1, \bar{x}_2, \dots, \bar{x}_N).$$

A lower stability indicates that the means are consistent across tiles, suggesting a stable time series. A higher stability implies significant differences in means, indicating potential shifts or trends in the data.

**Lumpiness** Lumpiness assesses the variability of the variances across all tiles. It is defined as the variance of the variances of the tiled windows. Let  $s_i^2$  denote the variance of the  $i$ -th tile. Lumpiness is then computed as:

$$\text{Lumpiness} = \text{Var}(s_1^2, s_2^2, \dots, s_N^2).$$

A higher lumpiness suggests that the variability within the tiles differs significantly, indicating that the time series may have periods of high and low volatility. A lower lumpiness means the variances are similar across tiles, pointing to a more homogeneous time series in terms of variability.

## C EVALUATION METRICS

We use two metrics to evaluate performance of forecasters: Mean Absolute Scaled Error (MASE) for point forecasting ability and Continuous Ranked Probability Score (CRPS) for probabilistic forecasting. For both metrics we use gluonts library implementation to calculate final values (Alexandrov et al., 2020a).

**MASE** **R:** [ MASE (Mean Absolute Scaled Error) is an evaluation metric commonly used in time series analysis to assess forecast accuracy. Unlike metrics such as MAPE, MASE addresses issues of scale dependence and sensitivity to outliers. It is defined as the mean absolute error of the forecast  $\hat{Y}_t$ ,

1080 scaled by the mean absolute error of a naïve benchmark forecast, typically a one-step-ahead lag of  
 1081 the actual values. The formula for MASE is: ]

$$1082$$

$$1083 \text{MASE} = \frac{\frac{1}{n} \sum_{t=1}^n |Y_t - \hat{Y}_t|}{\frac{1}{n-1} \sum_{t=2}^n |Y_t - Y_{t-1}|},$$

$$1084$$

$$1085$$

$$1086$$

1087 where:

- 1088 •  $Y_t$  is the actual value at time  $t$ ,
- 1089 •  $\hat{Y}_t$  is the forecasted value at time  $t$ ,
- 1090 •  $n$  is the number of observations.

1091

1092 MASE is scale-independent, making it suitable for comparing forecast accuracy across different  
 1093 time series. A MASE value less than 1 indicates that the forecast performs better than the naïve  
 1094 benchmark, while a value greater than 1 indicates worse performance. It is particularly useful in  
 1095 scenarios with varying scales or when evaluating the effectiveness of forecasts relative to a simple  
 1096 baseline.

1097

1098 **CRPS** The *Continuous Ranked Probability Score* (CRPS) is a metric used in probabilistic forecast-  
 1099 ing to evaluate the accuracy of predicted cumulative distribution functions (CDFs) against observed  
 1100 values. Given a predicted distribution with CDF  $F$  and a ground truth value  $y$ , the CRPS is defined  
 1101 as:

$$1102$$

$$1103 \text{CRPS}(F, y) = \int_0^1 2\Lambda_\alpha(F^{-1}(\alpha), y) d\alpha,$$

$$1104$$

1105 where the quantile loss  $\Lambda_\alpha(q, y)$  is defined as:

$$1106$$

$$1107 \Lambda_\alpha(q, y) = (\alpha - \mathbf{1}\{y < q\})(y - q).$$

$$1108$$

1109

1110 In practice, computing the CRPS integral can be computationally intensive. To address this, we  
 1111 approximate the CRPS using a discrete sum over a finite set of quantile levels. This approximation,  
 1112 often referred to as the mean weighted quantile loss (Park et al., 2021), is given by:

$$1113$$

$$1114 \text{CRPS} \approx \frac{1}{K} \sum_{k=1}^K \text{wQL}[\alpha_k],$$

$$1115$$

1116 where  $K$  is the number of quantile levels, and  $\{\alpha_1, \alpha_2, \dots, \alpha_K\}$  are the selected quantile levels (e.g.,  
 1117  $\alpha_k = 0.1k$  for  $k = 1, 2, \dots, 9$  when  $K = 9$ ).

1118 The weighted quantile loss  $\text{wQL}[\alpha]$  for each quantile level  $\alpha$  is calculated as:

$$1119$$

$$1120 \text{wQL}[\alpha] = 2 \frac{\sum_t \Lambda_\alpha(\hat{q}_t(\alpha), y_t)}{\sum_t |y_t|},$$

$$1121$$

1122 where:

- 1123 •  $\hat{q}_t(\alpha)$  is the predicted  $\alpha$ -quantile at time step  $t$ ,
- 1124 •  $y_t$  is the actual observed value at time  $t$ ,
- 1125 •  $\Lambda_\alpha(\hat{q}_t(\alpha), y_t)$  is the quantile loss at time  $t$  for quantile level  $\alpha$ .

## D GIFT-EVAL TEST DATASETS

In this section we provide comprehensive list of datasets used in test portion of GIFT-Eval along with original sources, for details regarding the pretraining portion see Appendix E. We utilize 10 open domain sources to curate the benchmark, here we list each one along with its properties in detail. We incorporate Jena Weather<sup>6</sup> dataset following **Autoformer** (Wu et al., 2021). We process BizITObs Application, Service, and L2C<sup>7</sup> following the pipeline in **AutoMixer** (Palaskar et al., 2024). These datasets consist of business and IT observability data, fusing both business KPIs and IT event channels into multivariate time series data. Within the same domain we also process Bitbrains datasets from **Grid Workloads Archive** (Shen et al., 2015). The Restaurant data is borrowed from **Recruit Restaurant Forecasting Competition** (Howard et al., 2017), The task associated with this dataset is to use reservation and visitation data to predict the total number of visitors to a restaurant for future dates. From **Informer** (Zhou et al., 2021) we utilize ETT1 and ETT2 datasets, which denote electricity transformer temperature and serve as an indicator used in the electricity power long-term deployment. Datasets for Transport domain are extracted from **LibCity** (Wang et al., 2023a), which provides a collection of urban time series datasets. We utilize the solar dataset from **LSTNet** (Lai et al., 2017) where the task is to predict solar plant energy outputs. The second dataset for Sales data is by Mancuso et al. (2020). **Monash** (Godahewa et al., 2021) is a large collection of diverse time series datasets across many domains, we choose a subset of these datasets making sure there is no leak from pretrain to test split. Finally, from **UCI ML Archive** (Trindade, 2015) we use the electricity dataset which contains electricity consumption of 370 individual clients. Table 14 lists all datasets, along with their source, frequency, prediction length and number of variates setup and presents various statistics from number of series, to series length, and also number of observations. We use last 10% of each timeseries in the test portion of our data for testing and keep the rest for training.

Tables 10 to 13 present detailed statistics on the number of time series and total observations within each characteristic category of the test benchmark. Specifically, these tables break down the data by domain (Table 11), frequency (Table 12), prediction length (Table 10), and variate count (Table 13), offering a quantitative overview of the dataset’s composition.

Table 10: GIFT-Eval Test data statistics aggregated by prediction length.

Pred. Length	6	8	12	13	14	18	30	48	60	480	600	720	900
# Series	22,974	24,629	3,443	359	4,227	48,000	34,398	6,194	22	3,874	22	3,874	22
# Obs	845,109	2,525,512	201,042	371,579	10,023,836	11,246,411	1,447,848	131,125,706	194,369	129,375,020	194,369	129,375,020	194,369

Table 11: GIFT-Eval Test data statistics aggregated by domain.

Domain	Econ/Fin	Energy	Healthcare	Nature	Sales	Transport	Web/CloudOps	Grand Total
# Series	99,974	2,036	1,036	32,618	3,717	1,341	3,524	144,246
# Obs	25,266,415	74,119,755	129,408	3,154,921	671,707	38,028,955	16,610,251	157,981,412

Table 12: GIFT-Eval Test data statistics aggregated by frequency.

Frequency	10S	10T	15T	5T	A	D	H	M	Q	W	Grand Total
# Series	22	138	528	2,074	22,974	38,625	3,454	51,443	24,000	988	144,246
# Obs	194,369	7,253,424	52,498,336	49,105,728	845,109	11,471,684	22,268,218	11,447,453	2,406,108	490,983	157,981,412

Table 13: GIFT-Eval Test data statistics aggregated by number of variates.

# Variates	1	2	7	21	Grand Total
# Series	140,711	3,522	10	3	144,246
# Obs	141,133,451	16,575,619	210,488	61,854	157,981,412

<sup>6</sup><https://www.bgc-jena.mpg.de/wetter/>

<sup>7</sup><https://github.com/BizITObs/BizITObservabilityData/tree/main>

Table 14: R: [Individual statistics of GIFT-Eval benchmark across all datasets.]

Dataset	Source	Domain	Frequency	# Series	Series Length			Target Variables	Feat Dynamic	Short-term		Mid-term		Long-term		
					Avg	Min	Max			Feat Length	Windows	Feat Length	Windows	Feat Length	Windows	
Japan Weather	Autoformer (https://www.hkg-jsta.mpg.de/wetter/)	Nature	HT	1	52,704	52,704	52,704	21	0	48	20	480	11	720	8	
Japan Weather	Autoformer (https://www.hkg-jsta.mpg.de/wetter/)	Nature	H	1	8,764	8,764	8,764	21	0	48	19	480	2	720	2	
Japan Weather	Autoformer (https://www.hkg-jsta.mpg.de/wetter/)	Nature	D	1	366	366	366	21	0	30	2	480	2	720	2	
BiUF0ms - Application	Autoformer (https://www.hkg-jsta.mpg.de/wetter/)	WebCloudOps	HT	1	8,835	8,835	8,835	2	35	60	15	600	2	900	1	
BiUF0ms - Service	Autoformer (https://www.hkg-jsta.mpg.de/wetter/)	WebCloudOps	HT	21	8,835	8,835	8,835	185,535	2	34	60	15	600	2	900	1
BiUF0ms - L2C	Autoformer (https://www.hkg-jsta.mpg.de/wetter/)	WebCloudOps	HT	1	31,968	31,968	31,968	7	2	48	20	480	7	720	5	
BiUF0ms - L2C	Autoformer (https://www.hkg-jsta.mpg.de/wetter/)	WebCloudOps	H	1	2,664	2,664	2,664	7	2	48	6	480	1	720	1	
Bitcoins - Fast Storage	Grid Workloads Archive	WebCloudOps	ST	1,250	8,640	8,640	8,640	10,000,000	2	5	48	18	480	2	720	2
Bitcoins - Fast Storage	Grid Workloads Archive	WebCloudOps	H	1,250	721	721	721	901,250	2	5	48	2	480	2	720	2
Bitcoins - rd	Grid Workloads Archive	WebCloudOps	ST	500	8,640	8,640	8,640	4,320,000	2	5	48	18	480	2	720	2
Bitcoins - rd	Grid Workloads Archive	WebCloudOps	H	500	720	720	720	300,000	2	5	48	2	480	2	720	2
Restaurant	https://www.kaggle.com/competitions/yelp-restaurant-visitor-forecasting/overview	Sales	D	807	358	67	478	289,303	1	0	30	1	480	15	720	10
EFT1	Informa	Energy	HT	69,680	69,680	69,680	69,680	7	0	48	20	480	4	720	3	
EFT1	Informa	Energy	H	1	17,420	17,420	17,420	17,420	7	0	48	20	480	4	720	3
EFT1	Informa	Energy	D	1	728	728	728	728	7	0	30	1	480	15	720	10
EFT2	Informa	Energy	WTHU	1	103	103	103	103	7	0	8	2	480	15	720	10
EFT2	Informa	Energy	HT	1	69,680	69,680	69,680	69,680	7	0	48	20	480	4	720	3
EFT2	Informa	Energy	H	1	17,420	17,420	17,420	17,420	7	0	48	20	480	4	720	3
EFT2	Informa	Energy	D	1	728	728	728	728	7	0	30	1	480	15	720	10
Loop Seattle	LibCity	Transport	WTHU	1	103	103	103	103	7	0	8	2	480	15	720	10
Loop Seattle	LibCity	Transport	ST	323	105,120	105,120	105,120	33,053,560	1	0	48	20	480	20	720	15
Loop Seattle	LibCity	Transport	H	323	8,760	8,760	8,760	2,829,480	1	0	48	19	280	2	720	2
Loop Seattle	LibCity	Transport	D	1	323	365	365	117,895	1	0	30	2	480	1	720	1
SZ-Taxi	LibCity	Transport	HT	156	2,976	2,976	2,976	464,256	1	0	48	2	480	4	720	3
SZ-Taxi	LibCity	Transport	H	156	744	744	744	116,064	1	0	48	2	480	4	720	3
M_BENSE	LibCity	Transport	H	30	17,520	17,520	17,520	25,660	1	7	48	20	480	4	720	3
M_BENSE	LibCity	Transport	D	30	730	730	730	21,900	1	7	30	3	480	11	720	8
Solar	ESTNet	Energy	HT	137	8,760	8,760	8,760	1,200,120	1	0	48	19	480	2	720	2
Solar	ESTNet	Energy	D	137	365	365	365	50,008	1	0	30	2	480	2	720	2
Solar	ESTNet	Energy	W-FRI	137	52	52	52	7,124	1	0	8	1	480	1	720	1
Hierarchical Sales	Monash et al.	Sales	D	118	1,825	1,825	1,825	245,590	1	0	30	1	480	1	720	1
M4 Yearly	Monash et al.	Sales	W-WED	118	260	260	260	30,660	1	0	8	4	480	1	720	1
M4 Quarterly	Monash	Q-DIE	118	260	260	260	30,660	1	0	8	4	480	1	720	1	
M4 Monthly	Monash	Q-DIE	24,000	100	24	874	2,406,108	1	0	8	1	480	1	720	1	
M4 Weekly	Monash	W-SUN	859	1,035	93	2,610	371,579	1	0	13	1	480	1	720	1	
M4 Daily	Monash	W-SUN	859	1,035	93	2,610	371,579	1	0	13	1	480	1	720	1	
M4 Hourly	Monash	W-SUN	859	1,035	93	2,610	371,579	1	0	13	1	480	1	720	1	
COVID-Deaths	Monash	Healthcare	M	414	902	748	1,008	373,372	1	0	48	2	480	2	720	2
US Births	Monash	Healthcare	M	767	84	84	84	64,428	1	0	12	1	480	1	720	1
US Births	Monash	Healthcare	D	266	212	212	212	56,392	1	0	8	14	480	1	720	1
US Births	Monash	Healthcare	D	1	7,305	7,305	7,305	7,305	1	0	30	20	480	1	720	1
Saugen	Monash	Healthcare	W-FRI	1	1,043	1,043	1,043	1,043	1	0	8	14	480	1	720	1
Saugen	Monash	Healthcare	M	1	240	240	240	240	1	0	12	2	480	1	720	1
Saugen	Monash	Healthcare	D	1	23,741	23,741	23,741	23,741	1	0	30	20	480	1	720	1
Temperature Rain	Monash	Nature	WTHU	1	3,391	3,391	3,391	3,391	1	0	8	20	480	1	720	1
KDD Cup 2018	Monash	Nature	M	1	780	780	780	780	1	0	12	7	480	2	720	2
KDD Cup 2018	Monash	Nature	D	32,072	725	725	725	780	1	0	30	3	480	2	720	2
KDD Cup 2018	Monash	Nature	D	270	10,938	9,504	10,920	2,942,354	1	0	48	20	480	2	720	2
Electricity	UCT ML Archive	Energy	HT	370	140,256	140,256	140,256	51,894,720	1	0	48	20	480	20	720	20
Electricity	UCT ML Archive	Energy	H	370	35,064	35,064	35,064	12,973,680	1	0	48	20	480	8	720	5
Electricity	UCT ML Archive	Energy	D	370	1,461	1,461	1,461	540,570	1	0	30	3	480	1	720	1
Electricity	UCT ML Archive	Energy	W-FRI	370	208	208	208	76,960	1	0	8	3	480	1	720	1

## E GIFT-EVAL PRE-TRAINING DATASETS

The pre-training split of GIFT-Eval is constructed based on LOTSA (Woo et al., 2024), and we excluded certain datasets from it to form part of the evaluation set, making it more diverse and balanced.

Here is a brief discussion on each of the used sources: **BuildingsBench** (Emami et al., 2023) compiled datasets on residential and commercial building energy consumption. **ClimateLearn** (Nguyen et al., 2023) offered time series of various climate-related variables, including temperature, humidity, and multiple pressure levels. **CloudOps TSF** (Woo et al., 2023) introduced large-scale CloudOps time series datasets that capture key variables such as CPU and memory utilization. **GluonTS** (Alexandrov et al., 2020a) provided a variety of datasets commonly used in time series forecasting. **LargeST** (Liu et al., 2023a) sourced from the California Department of Transportation Performance Measurement System (PeMS) to date, which is widely used for traffic forecasting. **LibCity** (Wang et al., 2023a) provided a collection urban spatio-temporal datasets. **SubseasonalClimateUSA** (Mouatadid et al., 2023) provided climate time series data at daily level. **ProEnFo** (Wang et al., 2023b) introduced a range of datasets for load forecasting which include various covariates such as temperature, humidity, and wind speed. **Monash** (Godaheewa et al., 2021) is a large collection of diverse time series datasets, the most popular source for building time series foundation models. **LOTSA\_Others** (Woo et al., 2024) are complementary datasets collected by LOTSA to enhance the diversity.

The complete list of pre-training datasets and their respective sources, key properties are provided in Table 15.

## F FINEGRAINED RESULTS

### F.1 R: [RESULTS AGGREGATED BY TIME SERIES FEATURES]

R: [ Table 16 shows results aggregated by time series features. We use the determined features of each dataset as depicted in Table 9 and aggregated across both high and low property of each feature in order get these results. The Table reveals a few key insights about model performance under varying features:]

R: [Temporal Strength (Trend and Seasonality): When datasets exhibit low trend or low seasonality, deep learning models such as PatchTST, TFT, and iTransformer generally perform better compared against foundation models. This may suggest that deep learning models are effective in handling scenarios with less pronounced temporal dynamics. In contrast, at higher levels of trend or seasonality, foundation models like MoiraiLarge often show better performance, potentially indicating their strength in capturing more regular patterns.]

Table 15: Pretraining datasets and their key properties.

Dataset	Source	Domain	Frequency	# Time Series	# Targets	# Covariates	# Obs.
BDG-2 Panther	BuildingsBench (Emami et al., 2023)	Energy	H	105	1	0	919,800
BDG-2 Fox	BuildingsBench (Emami et al., 2023)	Energy	H	135	1	0	2,324,568
BDG-2 Rat	BuildingsBench (Emami et al., 2023)	Energy	H	280	1	0	4,728,288
BDG-2 Bear	BuildingsBench (Emami et al., 2023)	Energy	H	91	1	0	1,482,312
Low Carbon London	BuildingsBench (Emami et al., 2023)	Energy	H	713	1	0	9,543,348
SMART	BuildingsBench (Emami et al., 2023)	Energy	H	5	1	0	95,709
IDEAL	BuildingsBench (Emami et al., 2023)	Energy	H	219	1	0	1,265,672
Sceaux	BuildingsBench (Emami et al., 2023)	Energy	H	1	1	0	34,223
Borealis	BuildingsBench (Emami et al., 2023)	Energy	H	15	1	0	83,269
Buildings900K	BuildingsBench (Emami et al., 2023)	Energy	H	1,792,328	1	0	15,702,590,000
CMIP6	ClimateLearn (Nguyen et al., 2023)	Climate	6H	1,351,680	53	0	1,973,453,000
ERA5	ClimateLearn (Nguyen et al., 2023)	Climate	H	245,760	45	0	2,146,959,000
Azure VM Traces 2017	CloudOpsTSF (Woo et al., 2023)	CloudOps	5T	159,472	1	2	885,522,908
Borg Cluster Data 2011	CloudOpsTSF (Woo et al., 2023)	CloudOps	5T	143,386	2	5	537,552,854
Alibaba Cluster Trace 2018	CloudOpsTSF (Woo et al., 2023)	CloudOps	5T	58,409	2	6	95,192,530
Taxi	GluonTS (Alexandrov et al., 2020a)	Transport	30T	67,984	1	0	54,999,060
Uber TLC Daily	GluonTS (Alexandrov et al., 2020a)	Transport	D	262	1	0	47,087
Uber TLC Hourly	GluonTS (Alexandrov et al., 2020a)	Transport	H	262	1	0	1,129,444
Wiki-Rolling	GluonTS (Alexandrov et al., 2020a)	Web	D	47,675	1	0	40,619,100
M5	GluonTS (Alexandrov et al., 2020a)	Sales	D	30,490	1	0	58,327,370
LargeST	LargeST (Liu et al., 2023a)	Transport	5T	42,333	1	0	4,452,510,528
PEMS03	LibCity (Wang et al., 2023a)	Transport	5T	358	1	0	9,382,464
PEMS04	LibCity (Wang et al., 2023a)	Transport	5T	307	3	0	5,216,544
PEMS07	LibCity (Wang et al., 2023a)	Transport	5T	883	1	0	24,921,792
PEMS08	LibCity (Wang et al., 2023a)	Transport	5T	170	3	0	3,035,520
PEMS Bay	LibCity (Wang et al., 2023a)	Transport	5T	325	1	0	16,937,700
Los-Loop	LibCity (Wang et al., 2023a)	Transport	5T	207	1	0	7,094,304
Beijing Subway	LibCity (Wang et al., 2023a)	Transport	30T	276	2	11	248,400
SHMetro	LibCity (Wang et al., 2023a)	Transport	15T	288	2	0	1,934,208
HZMetro	LibCity (Wang et al., 2023a)	Transport	15T	80	2	0	146,000
Q-Traffic	LibCity (Wang et al., 2023a)	Transport	15T	45,148	1	0	264,386,688
Subseasonal	SubseasonalClimateUSA (Mouatadid et al., 2023)	Climate	D	862	4	0	14,097,148
Subseasonal Precipitation	SubseasonalClimateUSA (Mouatadid et al., 2023)	Climate	D	862	1	0	9,760,426
Covid19 Energy	ProEnFo (Wang et al., 2023b)	Energy	H	1	1	6	31,912
GEF12	ProEnFo (Wang et al., 2023b)	Energy	H	20	1	1	788,280
GEF14	ProEnFo (Wang et al., 2023b)	Energy	H	1	1	1	17,520
GEF17	ProEnFo (Wang et al., 2023b)	Energy	H	8	1	1	140,352
PDB	ProEnFo (Wang et al., 2023b)	Energy	H	1	1	1	17,520
Spanish	ProEnFo (Wang et al., 2023b)	Energy	H	1	1	1	35,064
BDG-2 Hog	ProEnFo (Wang et al., 2023b)	Energy	H	24	1	5	421,056
BDG-2 Bull	ProEnFo (Wang et al., 2023b)	Energy	H	41	1	3	719,304
BDG-2 Cockatoo	ProEnFo (Wang et al., 2023b)	Energy	H	1	1	5	17,544
ELF	ProEnFo (Wang et al., 2023b)	Energy	H	1	1	0	21,792
London Smart Meters	Monash (Godaheva et al., 2021)	Energy	30T	5,520	1	0	166,238,880
Wind Farms	Monash (Godaheva et al., 2021)	Energy	T	337	1	0	172,165,370
Wind Power	Monash (Godaheva et al., 2021)	Energy	4S	1	1	0	7,397,147
Solar Power	Monash (Godaheva et al., 2021)	Energy	4S	1	1	0	7,397,222
Oikolab Weather	Monash (Godaheva et al., 2021)	Climate	H	8	1	0	800,456
Elecdemand	Monash (Godaheva et al., 2021)	Energy	30T	1	1	0	17,520
Covid Mobility	Monash (Godaheva et al., 2021)	Transport	D	362	1	0	148,602
Kaggle Web Traffic Weekly	Monash (Godaheva et al., 2021)	Web	W	145,063	1	0	16,537,182
Extended Web Traffic	Monash (Godaheva et al., 2021)	Web	D	145,063	1	0	370,926,091
M1 Yearly	Monash (Godaheva et al., 2021)	Econ/Fin	Y	106	1	0	3,136
M1 Quarterly	Monash (Godaheva et al., 2021)	Econ/Fin	Q	198	1	0	9,854
M1 Monthly	Monash (Godaheva et al., 2021)	Econ/Fin	M	617	1	0	44,892
M3 Yearly	Monash (Godaheva et al., 2021)	Econ/Fin	Y	645	1	0	18,319
M3 Quarterly	Monash (Godaheva et al., 2021)	Econ/Fin	Q	756	1	0	37,004
M3 Monthly	Monash (Godaheva et al., 2021)	Econ/Fin	M	1,428	1	0	141,858
M3 Other	Monash (Godaheva et al., 2021)	Econ/Fin	Q	174	1	0	11,933
NN5 Daily	Monash (Godaheva et al., 2021)	Econ/Fin	D	111	1	0	81,585
NN5 Weekly	Monash (Godaheva et al., 2021)	Econ/Fin	W	111	1	0	11,655
Tourism Yearly	Monash (Godaheva et al., 2021)	Econ/Fin	Y	419	1	0	11,198
Tourism Quarterly	Monash (Godaheva et al., 2021)	Econ/Fin	Q	427	1	0	39,128
Tourism Monthly	Monash (Godaheva et al., 2021)	Econ/Fin	M	366	1	0	100,496
CIF 2016	Monash (Godaheva et al., 2021)	Econ/Fin	M	72	1	0	6,334
Traffic Weekly	Monash (Godaheva et al., 2021)	Transport	W	862	1	0	82,752
Traffic Hourly	Monash (Godaheva et al., 2021)	Transport	H	862	1	0	14,978,112
Australian Electricity Demand	Monash (Godaheva et al., 2021)	Energy	30T	5	1	0	1,153,584
Rideshare	Monash (Godaheva et al., 2021)	Transport	H	2,304	1	0	859,392
Sunspot	Monash (Godaheva et al., 2021)	Nature	D	1	1	0	73,894
Vehicle Trips	Monash (Godaheva et al., 2021)	Transport	D	329	1	0	32,512
Weather	Monash (Godaheva et al., 2021)	Climate	D	3,010	1	0	42,941,700
FRED MD	Monash (Godaheva et al., 2021)	Econ/Fin	M	107	1	0	76,612
Pedestrian Counts	Monash (Godaheva et al., 2021)	Transport	H	66	1	0	3,130,762
Bitcoin	Monash (Godaheva et al., 2021)	Econ/Fin	D	18	1	0	74,824
KDD Cup 2022	LOTSAs_Others (Woo et al., 2024)	Energy	10T	134	1	9	4,727,519
GoDaddy	LOTSAs_Others (Woo et al., 2024)	Econ/Fin	M	3,135	2	0	128,535
Favorita Sales	LOTSAs_Others (Woo et al., 2024)	Sales	D	111,840	1	0	139,179,538
Favorita Transactions	LOTSAs_Others (Woo et al., 2024)	Sales	D	54	1	0	84,408
China Air Quality	LOTSAs_Others (Woo et al., 2024)	Nature	H	437	6	0	5,739,234
Beijing Air Quality	LOTSAs_Others (Woo et al., 2024)	Nature	H	12	11	0	420,768
Residential Load Power	LOTSAs_Others (Woo et al., 2024)	Energy	T	271	3	0	145,994,559
Residential PV Power	LOTSAs_Others (Woo et al., 2024)	Energy	T	233	3	0	125,338,950
CDC Fluview ILLNet	LOTSAs_Others (Woo et al., 2024)	Healthcare	W	75	5	0	63,903
CDC Fluview WHO NREVS	LOTSAs_Others (Woo et al., 2024)	Healthcare	W	74	4	0	41,760
Project Tycho	LOTSAs_Others (Woo et al., 2024)	Healthcare	W	1,258	1	0	1,377,707

R: [Entropy and Hurst: High entropy (suggesting greater forecasting difficulty) and low Hurst (indicating stochasticity thus again high difficulty) appear to favor Transformer-based models (except for Moirai variants). On the other hand, for lower entropy or higher Hurst values, foundation models, achieve better results on average, suggesting they may handle simpler temporal structures more effectively.]







1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

Table 19: Results on GIFT-Eval with all models aggregated by frequency. The best results across each row are **bolded**, while second best results are underlined.

Model	Nv.	S.Nv.	A.Ar.	A.Th.	A.ETS	D.AR	FTT	TIDE	N-B	P.FST	lTr.	DLin.	C.former	Timer	T.M	L-Llama	T.FM	V.TS	Chr.s	Chr.s	Chr.l	Moi.s	Moi.s	Moi.l	Best			
10S	MAPE	1.27	1.00	1.00	4.49e-1	1.19	9.12e-1	1.46	8.48e-1	7.45e-1	6.96e-1	6.47e-1	1.08	1.53	1.34	1.67	1.34	1.80	7.21e-1	1.21	1.19	1.16	1.34	1.78	1.24	A.Th.		
	MASE	1.98	1.00	1.00	1.59e+1	5.51e+1	3.76e+1	3.37e+1	3.23e+1	2.71e+1	2.21e+1	2.35e+1	3.08e+1	6.13e+1	5.17e+1	7.29e+1	5.98e+1	7.97e+1	2.16e+1	5.23e+1	5.23e+1	5.06e+1	7.95e+1	8.44e+1	5.72e+1	A.Th.		
	MAE	1.24	1.00	1.00	3.51e-1	1.01	8.02e-1	6.73e-1	7.36e-1	5.62e-1	5.58e-1	5.44e-1	7.35e-1	6.86e-1	9.00e-1	1.44	1.17	1.35	6.50e-1	8.03e-1	8.72e-1	8.27e-1	1.34	1.20	1.12	A.Th.		
	MSE	1.29	1.00	1.00	1.03e+1	1.06	5.01e+1	4.60e+1	4.31e+1	3.02e+1	3.30e+1	2.02e+1	4.13e+1	3.96e+1	6.00e+1	1.08	1.15	1.17	4.44e+1	5.85e+1	6.02e+1	5.63e+1	1.46	1.06	1.07	A.Th.		
	MAE	1.24	1.00	1.00	3.54e-1	1.01	8.02e-1	6.73e-1	7.36e-1	5.61e-1	5.58e-1	5.44e-1	7.35e-1	6.86e-1	9.00e+1	1.44	1.17	1.35	6.50e-1	8.03e-1	8.72e-1	8.27e-1	1.34	1.20	1.12	A.Th.		
	CRPS	1.44	1.00	1.00	3.15e-1	1.93	7.54e-1	6.72e-1	7.05e-1	5.98e-1	5.36e-1	5.01e-1	7.82e-1	7.29e-1	9.58e-1	1.53	1.12	1.30	6.91e-1	7.93e-1	8.59e-1	8.18e-1	1.24	1.06	1.02	A.Th.		
	Rank	1.95e+1	1.13e+1	1.03e+1	1.00	2.58e+1	1.23e+1	8.83	1.12e+1	7.17	5.00	2.50	1.28e+1	1.22e+1	1.65e+1	2.55e+1	1.92e+1	1.63e+1	1.08e+1	1.12e+1	1.35e+1	1.23e+1	2.26e+1	3.95e+1	1.78e+1	A.Th.		
10T	MAPE	9.42e-1	1.00	1.00	1.32	0.52e-1	1.56	1.15	1.28	1.17	1.11	1.13	1.30	1.40	1.27	1.35	1.34	1.18	1.16	1.24	1.24	1.23	1.10	1.01	9.35e-1	A.Th.		
	MASE	9.42e-1	1.00	1.00	9.84e-1	9.10e-1	1.40	8.36e-1	9.61e-1	8.84e-1	7.87e-1	7.73e-1	1.16	1.45	9.59e-1	1.22	1.02	2.38	8.19e-1	8.72e-1	8.62e-1	8.69e-1	1.89e-1	6.89e-1	6.09e-1	A.Th.		
	MAE	7.99e-1	1.00	1.00	8.85e-1	9.13e-1	8.83e-1	6.95e-1	7.93e-1	7.29e-1	6.54e-1	6.63e-1	8.29e-1	7.30e-1	8.76e-1	8.59e-1	9.84e-1	8.03e-1	7.43e-1	7.72e-1	7.68e-1	7.75e-1	6.32e-1	6.48e-1	6.16e-1	A.Th.		
	MSE	7.23e-1	1.00	1.00	7.41e+1	8.79e+1	7.85e+1	5.12e+1	5.76e+1	5.65e+1	4.81e+1	4.56e+1	5.55e+1	5.49e+1	6.30e+1	6.12e+1	9.02e+1	6.30e+1	5.36e+1	6.09e+1	7.03e+1	5.90e+1	5.09e+1	4.69e+1	4.69e+1	A.Th.		
	MAE	7.99e-1	1.00	1.00	8.84e-1	9.13e-1	8.83e-1	6.96e-1	7.93e-1	7.29e-1	6.54e-1	6.63e-1	8.29e-1	7.30e-1	8.75e-1	8.59e-1	9.84e-1	8.03e-1	7.43e-1	7.71e-1	7.69e-1	7.75e-1	6.32e-1	6.47e-1	6.16e-1	A.Th.		
	CRPS	1.19	1.00	1.00	9.48e-1	9.58e-1	7.43e-1	5.36e-1	6.31e-1	6.09e-1	5.22e-1	5.22e-1	7.81e-1	6.96e-1	8.28e-1	8.12e-1	7.95e-1	6.73e-1	7.02e-1	6.82e-1	6.83e-1	6.87e-1	4.96e-1	4.84e-1	4.61e-1	A.Th.		
	Rank	2.34e+1	2.78e+1	2.34e+1	2.26e+1	9.05e+1	1.37e+1	6.58	1.33e+1	1.45e+1	6.75	7.75	9.92e+1	1.62e+1	2.18e+1	2.02e+1	1.96e+1	1.52e+1	1.63e+1	1.48e+1	1.51e+1	1.58e+1	7.47e+1	8.47e+1	4.58	A.Th.		
10V	MAPE	1.68	1.00	1.00	1.89	1.97	2.33	1.86	1.92	2.06	1.91	1.84	2.39	1.93	1.80	2.18	1.78	1.81	1.98	1.93	1.94	2.01	2.11	2.08	2.11	2.08	S.Nv.	
	MASE	1.28	1.00	1.00	1.62	1.77	1.55	9.42e-1	1.27	1.21	1.19	1.09	1.45	2.22	1.44	9.47e-1	1.48	1.27	9.12e-1	1.20	1.09	1.07	1.10	1.15	1.13	V.TS		
	MSE	1.51	1.00	1.00	1.88	2.11	1.87	1.36e+1	1.28	1.55	1.76e+1	1.21	1.35	1.29	1.75	1.03	1.61	1.09	9.09e-1	1.41	1.22	1.21	1.20	1.31	1.34	V.TS		
	MAE	1.57	1.00	1.00	1.67	2.99	1.36	7.85e-1	1.55	1.55	9.35e-1	1.17	8.25e-1	1.23	1.44	6.88e-1	1.77	1.06	8.02e-1	1.50	1.19	1.18	1.29	1.51	1.48	T.M		
	MSE	1.51	1.00	1.00	2.11	1.87	1.39	1.01	1.28	1.55	1.16	1.21	1.35	1.84	1.75	1.03	1.61	1.17	9.04e-1	1.41	1.22	1.21	1.20	1.31	1.34	V.TS		
	CRPS	2.08	1.00	1.00	2.33	1.42	5.37e+1	3.44e+1	5.69e+1	6.98e+1	4.13e+1	4.43e+1	5.99e+1	5.73e+1	7.78e+1	4.59e+1	6.31e+1	4.59e+1	4.42e+1	5.47e+1	4.75e+1	4.71e+1	4.91e+1	5.04e+1	5.14e+1	FTT		
	Rank	2.67e+1	2.22e+1	2.12e+1	2.80e+1	2.42e+1	1.47e+1	5.67	1.65e+1	1.82e+1	9.50	8.00	1.60e+1	1.62e+1	2.20e+1	1.03e+1	1.77e+1	1.00e+1	9.33	1.55e+1	1.05e+1	9.67	1.10e+1	1.28e+1	1.30e+1	FTT		
10Y	MAPE	1.41	1.00	1.00	9.78e-1	1.11	1.27	1.66	1.01	1.06	1.06	9.15e-1	1.02	1.81	1.23	1.07	1.23	9.96e-1	9.92e-1	1.02	9.68e-1	9.62e-1	9.87e-1	9.72e-1	1.00	1.00	1.00	1.00
	MASE	1.52	1.00	1.00	1.03	1.61	1.76	9.66e-1	1.02	1.03	9.77e-1	9.26e-1	9.92e-1	1.76	1.24	1.02	1.27	9.66e-1	9.92e-1	9.26e-1	9.87e-1	8.87e-1	9.49e-1	9.25e-1	9.77e-1	1.00	1.00	
	MAE	1.40	1.00	1.00	9.78e-1	1.10	1.39	1.56	1.03	1.02	8.74e-1	8.67e-1	1.43	1.30	1.50	1.05	1.36	9.54e-1	9.09e-1	9.13e-1	8.82e-1	8.78e-1	9.22e-1	9.22e-1	9.63e-1	1.00	1.00	
	MSE	1.99	1.00	1.00	9.59e-1	1.10	1.39	2.04	7.99e+1	9.00e+1	9.99e+1	2.20e+1	7.29e+1	8.71e+1	2.95	1.56	9.79e+1	1.89	8.61e+1	8.13e+1	8.51e+1	8.05e+1	8.06e+1	8.67e+1	8.96e+1	9.09e+1	1.00	
	MAE	1.49	1.00	1.00	9.78e-1	1.06	1.56	1.64	9.34e-1	1.00	1.02	8.74e-1	8.68e-1	9.83e-1	2.07	1.30	1.05	1.86	9.54e-1	9.09e-1	9.13e-1	8.82e-1	8.79e-1	9.42e-1	9.29e-1	9.64e-1	1.00	
	CRPS	2.20	1.00	1.00	9.52e-1	1.51	1.19e+1	1.26	7.08e-1	7.92e-1	9.63e-1	6.55e-1	9.26e-1	1.34	1.22	9.89e-1	1.02	7.68e-1	8.56e-1	7.73e-1	7.49e-1	7.46e-1	8.07e-1	6.91e-1	7.30e-1	7.47e-1	1.00	
	Rank	2.74e+1	2.03e+1	1.91e+1	2.38e+1	2.47e+1	1.07e+1	8.67	1.37e+1	2.00e+1	5.00	4.67	1.90e+1	1.74e+1	2.96e+1	1.96e+1	1.78e+1	1.07e+1	1.73e+1	1.29e+1	1.08e+1	1.06e+1	1.10e+1	1.10e+1	1.00	6.17	9.58	1.00
11	MAPE	1.46	1.00	1.00	1.21	1.40	1.37	1.34	1.06	1.13	1.08	1.03	1.15	1.61	1.13	1.06	1.13	9.88e-1	1.01	9.71e-1	9.65e-1	9.66e-1	1.09	9.91e-1	9.84e-1	9.84e-1	1.00	
	MASE	1.46	1.00	1.00	1.28	1.27	1.31	8.25e-1	9.50e-1	8.72e-1	7.74e-1	7.40e-1	9.43e-1	1.73	9.41e-1	8.52e-1	8.95e-1	8.27e-1	7.70e-1	7.73e-1	7.63e-1	7.63e-1	9.92e-1	7.86e-1	7.70e-1	7.70e-1	1.00	
	MAE	1.43	1.00	1.00	1.04	1.39	1.21	1.14	7.95e-1	8.99e-1	8.99e-1	7.57e-1	7.90e-1	9.07e-1	8.06e-1	9.47e-1	8.52e-1	8.92e-1	8.27e-1	7.86e-1	7.87e-1	7.77e-1	7.80e-1	8.62e-1	8.66e-1	7.86e-1	P.FST	
	MSE	1.85	1.00	1.00	1.04	1.48	1.35	1.15	6.12e-1	6.96e-1	7.61e-1	6.79e-1	6.26e-1	6.96e-1	6.96e-1	6.18e-1	7.70e-1	6.52e-1	6.96e-1	6.96e-1	6.81e-1	7.48e-1	6.17e-1	7.48e-1	6.17e-1	6.17e-1	P.FST	
	MAE	1.43	1.00	1.00	1.04	1.39	1.21	1.14	7.95e-1	8.99e-1	8.99e-1	7.57e-1	7.90e-1	9.06e-1	1.08	9.47e-1	8.51e-1	8.92e-1	8.27e-1	7.86e-1	7.87e-1	7.77e-1	7.79e-1	8.62e-1	7.66e-1	7.38e-1	P.FST	
	CRPS	1.67	1.00	1.00	7.43e-1	1.57	5.01e+1	6.23e+1	4.28e+1	5.11e+1	6.00e+1	4.07e+1	4.21e+1	6.96e+1	5.34e+1	6.33e+1	5.69e+1	4.89e+1	4.66e+1	5.25e+1	4.68e+1	4.62e+1	4.64e+1	5.13e+1	4.13e+1	4.07e+1	P.FST	
	Rank	2.75e+1	2.48e+1	2.29e+1	2.66e+1	2.64e+1	1.02e+1	9.27e+1	7.77e+1	1.44e+1	1.55e+1	8.32e+1	1.96e+1	1.59e+1	1.64e+1	1.19e+1	1.16e+1	1.19e+1	1.16e+1	1.19e+1	1.16e+1	1.19e+1	1.16e+1	1.19e+1	1.16e+1	1.19e+1	5.23	1.00
D	MAPE	1.00	1.00	1.00	9.37e-1	9.89e-1	8.95e-1	1.03	8.14e-1	1.11	8.60e-1	8.49e-1	9.24e-1	8.95e-1	2.00	1.03	1.00	1.30	8.26e-1	8.95e-1	8.39e-1	8.13e-1	8.13e-1	8.94e-1	8.55e-1	8.75e-1	1.00	
	MASE	1.00	1.00	1.00	8.82e-1	9.36e-1	9.01e-1	1.06	9.77e-1	1.15	7.57e-1	7.49e-1	8.31e-1	8.97e-1	4.83	9.92e-1	9.44e-1	1.19	7.46e-1	8.22e-1								

Table 22: Results on all dataset configs for GIFT-Eval | Table 1/3. The best results across each row are **bolded**, while second best results are underlined.

Table with 30 columns: Dataset, term, frequency, Metric, S.W., A.F., A.Th., A.Sts, D.AR, TPT, FLOP, S.B., P.FST, LFR, DLin, C.former, Timer, TTM, L-Llama, T.FM, V.FTS, Chr.1, Chr.2, Chr.3, Chr.4, Chr.5, Chr.6, Chr.7, Chr.8, Chr.9, Chr.10, Chr.11, Chr.12, Chr.13, Chr.14, Chr.15, Chr.16, Chr.17, Chr.18, Chr.19, Chr.20, Chr.21, Chr.22, Chr.23, Chr.24, Chr.25, Chr.26, Chr.27, Chr.28, Chr.29, Chr.30. Rows list various datasets like 'Dataset, term, frequency', 'hibrains\_fast\_storage\_long\_ST', etc., with corresponding metric values.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

Table 23: Results on all dataset configs for GIFT-Eval | Table 2/3. The best results across each row are **bolded**, while second best results are underlined.

Table with columns: Dataset, term, frequency, Metric, Nv., S.Nv., A.Ar., A.Th., A.FTS, D.AR, FTF, TIDE, N.B., P.TST, I.Tr., D.in., C.Fomer, Timer, TM, L-Liana, T.FM, V.TS, Chr.s, Chr.r, Chr.l, Mo.s, Mo.g, Mo.l, Best. The table contains a dense grid of numerical values for various datasets and metrics.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

Table 24: Results on all dataset configs for GIFT-Eval | Table 3/3. The best results across each row are **bolded**, while second best results are underlined.

Table with columns: Dataset, term, frequency, Metric, Nv., S.Siv., A.Ar., A.Th., A.ETS, D.AR, TPT, T10E, N-B., P.TST, L7E., D1a:n, C.Fomer, T1mer, T1M, L-Llama, T.F.M, V.TS, Chr.s, Chr.g, Chr.l., Mo1.s, Mo1.g, Mo1.l, Best. The table contains a dense grid of numerical values for various datasets and metrics.

Table 25: Moirai vs Moirai-Leakage results on datasets that LOTSA collection and our GIFT-Eval has in common.

Dataset	Model	Short		Medium		Long	
		MAPE	CRPS	MAPE	CRPS	MAPE	CRPS
hierarchical_sales, D	<b>Moi Leak.B</b>	0.51	0.24	NA	NA	NA	NA
hierarchical_sales, D	<b>Moi.B</b>	0.49	0.25	NA	NA	NA	NA
hierarchical_sales, D	<b>Moi Leak.L</b>	0.53	0.25	NA	NA	NA	NA
hierarchical_sales, D	<b>Moi.L</b>	0.52	0.24	NA	NA	NA	NA
hierarchical_sale, D	<b>Moi Leak.S</b>	0.50	0.25	NA	NA	NA	NA
hierarchical_sales, D	<b>Moi.S</b>	0.51	0.25	NA	NA	NA	NA
loop_seattle, 5T	<b>Moi Leak.B</b>	0.67	0.57	0.42	0.37	0.50	0.38
loop_seattle, 5T	<b>Moi.B</b>	0.84	0.64	0.83	0.62	0.78	0.57
loop_seattle, 5T	<b>Moi Leak.L</b>	0.66	0.51	0.33	0.31	0.46	0.36
loop_seattle, 5T	<b>Moi.L</b>	0.83	0.65	0.85	0.65	0.81	0.59
loop_seattle, 5T	<b>Moi Leak.S</b>	0.84	0.69	0.75	0.57	0.70	0.53
loop_seattle, 5T	<b>Moi.S</b>	0.87	0.65	0.77	0.61	0.75	0.57
loop_seattle, D	<b>Moi Leak.B</b>	0.50	0.34	NA	NA	NA	NA
loop_seattle, D	<b>Moi.B</b>	0.52	0.35	NA	NA	NA	NA
loop_seattle, D	<b>Moi Leak.L</b>	0.51	0.35	NA	NA	NA	NA
loop_seattle, D	<b>Moi.L</b>	0.49	0.33	NA	NA	NA	NA
loop_seattle, D	<b>Moi Leak.S</b>	0.53	0.35	NA	NA	NA	NA
loop_seattle, D	<b>Moi.S</b>	0.54	0.35	NA	NA	NA	NA
loop_seattle, H	<b>Moi Leak.B</b>	0.96	0.68	0.54	0.50	0.49	0.26
loop_seattle, H	<b>Moi.B</b>	1.08	0.72	0.65	0.55	0.59	0.30
loop_seattle, H	<b>Moi Leak.L</b>	0.84	0.61	0.53	0.45	0.47	0.23
loop_seattle, H	<b>Moi.L</b>	0.89	0.65	0.71	0.59	1.18	0.45
loop_seattle, H	<b>Moi Leak.S</b>	1.22	0.80	0.73	0.64	0.70	0.35
loop_seattle, H	<b>Moi.S</b>	1.19	0.78	0.70	0.60	0.71	0.33
m_dense, D	<b>Moi Leak.B</b>	0.78	0.35	NA	NA	NA	NA
m_dense, D	<b>Moi.B</b>	0.55	0.27	NA	NA	NA	NA
m_dense, D	<b>Moi Leak.L</b>	0.67	0.32	NA	NA	NA	NA
m_dense, D	<b>Moi.L</b>	0.63	0.31	NA	NA	NA	NA
m_dense, D	<b>Moi Leak.S</b>	0.58	0.28	NA	NA	NA	NA
m_dense, D	<b>Moi.S</b>	0.53	0.26	NA	NA	NA	NA
m_dense, H	<b>Moi Leak.B</b>	0.54	0.50	0.49	0.26	0.53	0.22
m_dense, H	<b>Moi.B</b>	0.65	0.55	0.59	0.30	0.61	0.26
m_dense, H	<b>Moi Leak.L</b>	0.53	0.45	0.47	0.23	0.49	0.21
m_dense, H	<b>Moi.L</b>	0.71	0.59	1.18	0.45	1.69	0.45
m_dense, H	<b>Moi Leak.S</b>	0.73	0.64	0.70	0.35	0.71	0.31
m_dense, H	<b>Moi.S</b>	0.70	0.60	0.71	0.33	0.91	0.33
restaurant	<b>Moi Leak.B</b>	0.70	0.29	NA	NA	NA	NA
restaurant	<b>Moi.B</b>	0.72	0.31	NA	NA	NA	NA
restaurant	<b>Moi Leak.L</b>	0.75	0.30	NA	NA	NA	NA
restaurant	<b>Moi.L</b>	0.76	0.30	NA	NA	NA	NA
restaurant	<b>Moi Leak.S</b>	0.74	0.31	NA	NA	NA	NA
restaurant	<b>Moi.S</b>	0.74	0.31	NA	NA	NA	NA
sz_taxi, 15T	<b>Moi Leak.B</b>	0.90	0.69	0.71	0.33	2.16	0.38
sz_taxi, 15T	<b>Moi.B</b>	0.84	0.69	0.64	0.46	2.42	0.38
sz_taxi, 15T	<b>Moi Leak.L</b>	0.78	0.69	0.71	0.46	2.14	0.38
sz_taxi, 15T	<b>Moi.L</b>	0.82	0.69	0.60	0.47	2.24	0.38
sz_taxi, 15T	<b>Moi Leak.S</b>	0.95	0.69	0.65	0.47	2.12	0.38
sz_taxi, 15T	<b>Moi.S</b>	1.11	0.70	0.60	0.47	2.30	0.39
sz_taxi, H	<b>Moi Leak.B</b>	0.64	0.62	NA	NA	NA	NA
sz_taxi, H	<b>Moi.B</b>	0.72	0.64	NA	NA	NA	NA
sz_taxi, H	<b>Moi Leak.L</b>	0.63	0.64	NA	NA	NA	NA
sz_taxi, H	<b>Moi.L</b>	0.70	0.64	NA	NA	NA	NA
sz_taxi, H	<b>Moi Leak.S</b>	0.65	0.66	NA	NA	NA	NA
sz_taxi, H	<b>Moi.S</b>	0.77	0.65	NA	NA	NA	NA

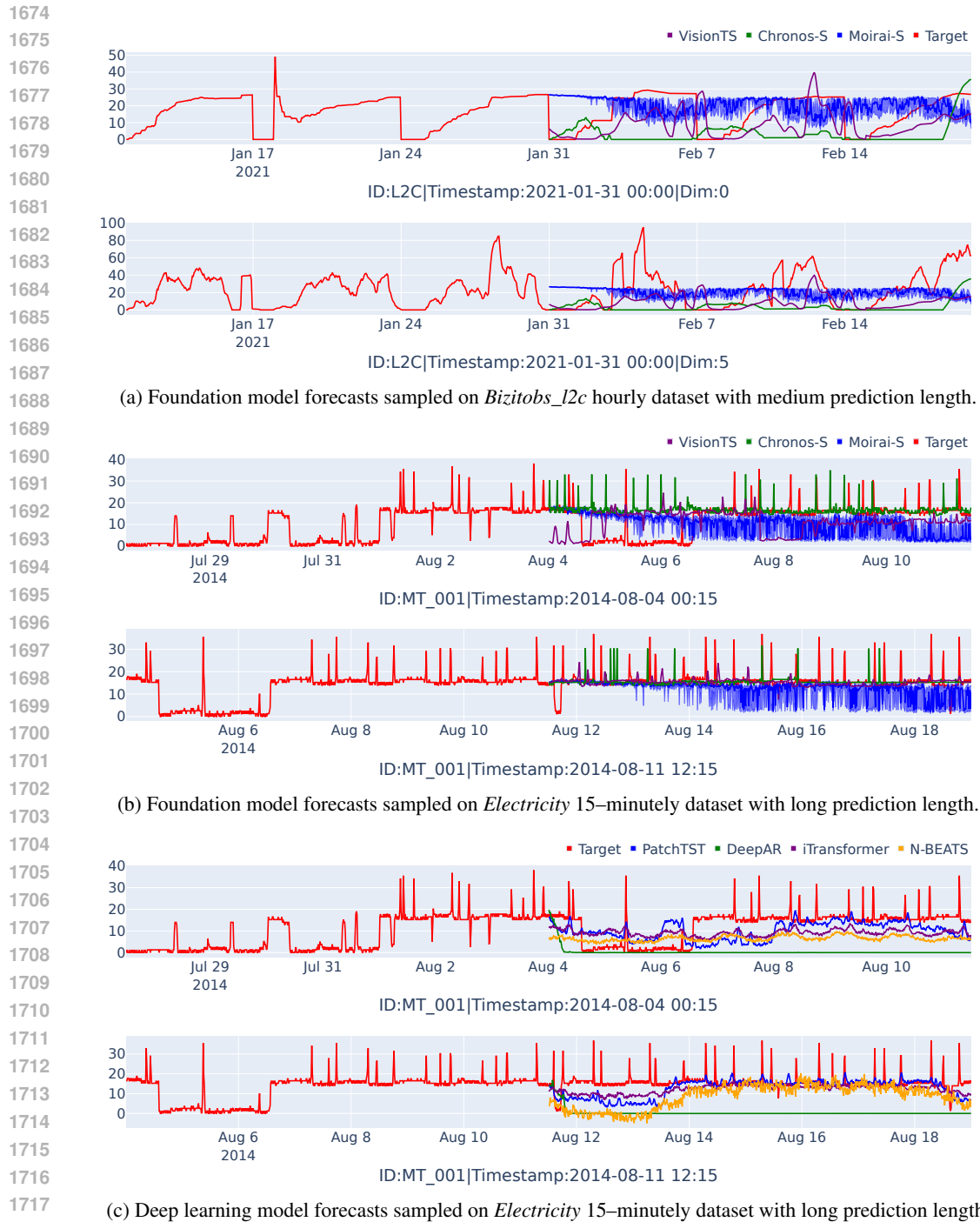


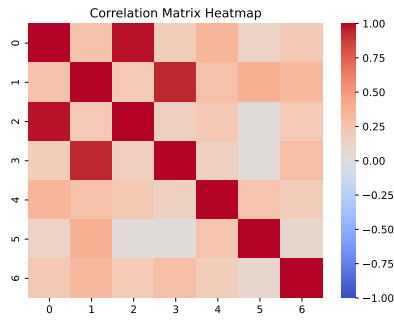
Figure 3: Qualitative plots showing forecasts from various deep learning and foundation models on several time series forecasting datasets.

datasets had entirely uncorrelated variates, there would be little justification for using multivariate models, as univariate models could predict each variate independently with comparable effectiveness.]

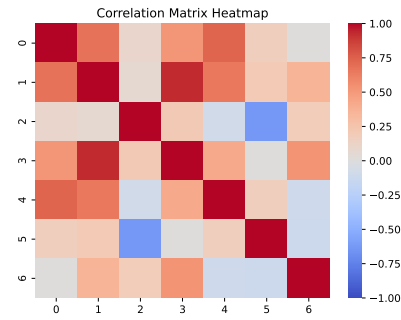
R: [ Figure 4 illustrates the correlation matrices for each multivariate dataset, highlighting the degree of inter-variate correlation. Specific statistics for each dataset are provided in the respective figure captions to offer additional insights into their characteristics.]



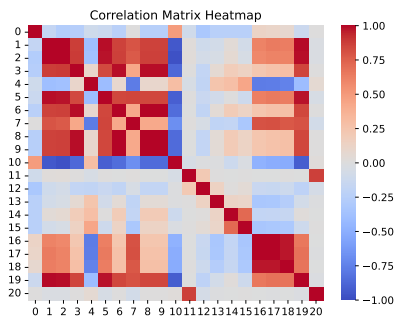
1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781



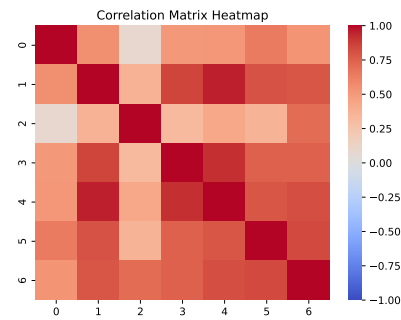
(a) Correlation matrix for Dataset ETT1. Mean: 0.38, Median: 0.25, Std: 0.33.



(b) Correlation matrix for Dataset ETT2. Mean: 0.34, Median: 0.21, Std: 0.41.



(c) Correlation matrix for Dataset Jena Weather. Mean: 0.18, Median: 0.03, Std: 0.50.



(d) Correlation matrix for Dataset Bizitobs\_l2c. Mean: 0.68, Median: 0.74, Std: 0.24.

Figure 4: R: [Inter-variate correlation matrices for selected multivariate datasets in GIFT-Eval. Each heatmap visualizes the correlation across variates for a specific dataset, highlighting the strength and distribution of inter-variate dependencies. Descriptive statistics (mean, median, standard deviation) are provided in the subcaptions for further insight.]