TOTO: TIME SERIES OPTIMIZED TRANSFORMER FOR OBSERVABILITY

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023 024 025

026 027 028

029

031

033 034

037

040

041

042

043 044 045 Paper under double-blind review

ABSTRACT

We introduce the Time Series Optimized Transformer for Observability (Toto), a foundation model designed for time series forecasting with a focus on observability metrics. Toto features a novel proportional factorized attention mechanism and a Student-T mixture model head, enabling it to efficiently handle high-dimensional, sparse, and non-stationary data. Trained on one trillion time series data points, including 75% proprietary observability data, Toto demonstrates state-of-the-art zero-shot performance on standard benchmarks such as electricity and weather forecasting. Furthermore, it significantly outperforms existing models in observability-specific tasks, making it an ideal solution for real-time system monitoring and anomaly detection. Toto's architectural innovations make it a versatile tool for both general-purpose forecasting and domain-specific applications, setting a new benchmark for scalability and accuracy in time series analysis.

1 INTRODUCTION

We present Toto, a time series forecasting foundation model specifically designed to handle the complexities of observability data. It leverages a novel transformer-based architecture to deliver state-of-the-art accuracy and performance. Toto is trained on a massive dataset of diverse time series data, enabling it to excel in zero-shot predictions. Our model is tailored to allow compute and memory-efficient scalability to very large data volumes, thereby providing robust solutions for high-frequency and high-dimensional data commonly encountered in observability metrics.

We detail the following key contributions:

Multivariate Time Series A-OBSERVABILITY DATA SYSTEM SYSTEM SYSTEM A-DECORTION A-DECORTIO

047 048 049



054	• Proportional factorized space-time attention: We introduce an advanced attention mech-
055	anism that allows for efficient grouping of multivariate time series features, reducing com-
056	putational overhead while maintaining high accuracy.
057	• Student-T mixture model head: This novel use of a probabilistic model that robustly gen-
058	eralizes Gaussian mixture models enables Toto to more accurately capture the complex dy-
059	namics of time series data and provides superior performance over traditional approaches.
060	• Domain anagina training data. In addition to consul multi domain time caries data
061	• Domain-specific training data: In addition to general multi-domain time series data, Toto is specifically trained on a large scale detaset of observability metrics, encompossing
062	unique characteristics not present in open-source datasets. This targeted training ensures
063	enhanced performance in observability metric forecasting.
064	······································
065	1.1. Orservarii ity data
066	
067	Observability data encompasses a comprehensive array of metrics collected to monitor and optimize
068	the performance and reliability of modern infrastructure and applications (Li et al., 2020). These
069	metrics are essential for providing insights into the health and performance of systems and include:
070	
071	• Infrastructure metrics: Data related to hardware and system performance, such as mem-
072	ory usage, CPU load, disk I/O, and network throughput.
073	• Application performance indicators: Metrics that capture the performance and behavior
074	of applications, including hit counts, error rates, and latency.
075	
076	Observability data is typically gathered from a variety of sources, including on-premise systems,
077	cloud services, and third-party tools. The integration of these diverse data sources enables a holistic view of system performance, but also introduces several challenges for time series forecesting:
078	view of system performance, but also infroduces several chancinges for time series forecasting.
079	1 High temporal resolution: Observability data often requires high-resolution timestamps
080 081	capturing data at intervals of seconds or minutes to detect rapid changes and anomalies.
082	2. Sparsity and zero-inflation: Many observability metrics are sparse, characterized by nu-
083	merous zero values due to the monitoring of infrequent events, such as system errors or
084	rare performance issues.
085	3. Extreme dynamic range and skewed distributions: Metrics can exhibit wide dynamic
086	ranges and heavy-tailed distributions, especially in latency measurements where occasional
087	extreme values occur.
088	4. Dynamic and non-stationary nature: The monitored systems are dynamic, undergoing
089	frequent changes due to software updates, infrastructure scaling, feature toggles, and vary-
090	ing user behaviors, all of which contribute to non-stationary data patterns.
091	5. High-cardinality multivariate data: Observability data often involves high-dimensional
092	metrics, segmented by various attributes like service type, region, or instance. This results
093	in a large number of time series, each with potentially limited historical data.
094	6 Historical anomalies: Historical data can contain anomalies and outliers resulting from
095	past performance issues or incidents, complicating the forecasting process.
096	
097	Effectively forecasting observability data requires advanced time series models that can manage
098	these complexities. Traditional forecasting methods often fall short due to their inability to scale
100	and adapt to the dynamic, high-dimensional nature of observability data. Therefore, there is a need
100	tor innovative models that can capture intricate patterns and dependencies, ultimately enhancing the
101	ability to proactively detect and mitigate performance issues in real-time systems.
102	
103	1.2 TRADITIONAL MODELS
104	

Historically, time series forecasting has relied on classical models such as ARIMA, exponential
 smoothing, and basic machine learning techniques (Hyndman & Athanasopoulos, 2021). While
 foundational, these models necessitate individual training for each metric, presenting several lim itations (Fildes et al., 1998). The need to develop and maintain separate models for each metric

impedes scalability, especially given the extensive range of metrics in observability data. More over, these models often fail to generalize across different types of metrics, leading to suboptimal
 performance on diverse datasets (Stevenson, 2007; Christodoulos et al., 2010).

112 1.3 RECENT WORK

Neural models, particularly those based on transformer architectures, have shown promise for im-114 proving the accuracy of time series forecasts. These models have demonstrated state-of-the-art 115 performance on benchmark datasets (Nie et al., 2023), frequently surpassing traditional models in 116 both accuracy and robustness. Their capacity to process high-dimensional data efficiently (Lin et al., 117 2021) makes them ideal for applications involving numerous time series metrics with varying char-118 acteristics, such as observability. However, in the full-shot setting, continuous retraining and tuning 119 to adapt to evolving data patterns create a significant operational burden for observability use cases. 120 This scaling limitation has hindered the adoption of deep learning-based methods for time series 121 analysis, even as they show promise in terms of accuracy (Salinas et al., 2020). 122

Even more recently, a number of time series "foundation models" have been released (Das et al., 2024; Ansari et al., 2024; Woo et al., 2024; Garza & Mergenthaler-Canseco, 2023; Rasul et al., 2023; Gruver et al., 2023). By pre-training on extensive, multi-domain datasets, these large models achieve impressive zero-shot prediction capabilities, significantly reducing the need for constant retraining.

127 128 129

133

138

139

140

141

142 143

144

145

146 147

111

1.4 ATTENTION MECHANISMS

To address the unique challenges of time series data, and particularly to adapt transformer architec tures for multivariate time-series forecasting, several works have implemented modifications to the
 attention mechanism. These strategies have included:

- Concatenating variates along the time dimension and computing full self-attention between every space/time location, as in the "any-variate attention" used by Woo et al. (2024). This can capture every possible space and time interaction, but it is costly in terms of computation and memory usage.
 - Assuming channel independence, and computing attention only in the time dimension as in Nie et al. (2023). This is efficient, but throws away all information about space-wise interactions.
 - Computing attention only in the space dimension, and using a feed-forward network in the time dimension (Ilbert et al., 2024; Liu et al., 2024).
 - Computing "factorized attention," where each transformer block contains a separate space and time attention computation (Zhang & Yan, 2023; Rao et al., 2021; Arnab et al., 2021). This allows both space and time mixing, and is more efficient than full cross-attention. However, it doubles the effective depth of the network.
- In Section 2.4, we propose a novel approach that allows for both space and time interactions, while reducing the computational cost and improving overall scalability.
- 150 151
 - 1.5 PROBABILISTIC OUTPUTS

Practitioners who rely on time series forecasting typically prefer probabilistic predictions. A common practice in neural time series models is to use an output layer where the model regresses the parameters of a probability distribution. This allows for prediction intervals to be computed using Monte Carlo sampling (Salinas et al., 2020).

- Common choices for an output layer are Normal (Salinas et al., 2020) and Student-T (Das et al., 2023; Rasul et al., 2023), which can improve robustness to outliers. Moirai (Woo et al., 2024) allows for more flexible residual distributions by proposing a novel mixture model incorporating a weighted combination of Gaussian, Student-T, Log-Normal, and Negative-Binomial outputs.
- 161 However, real-world time series can often have complex distributions that are challenging to fit, with outliers, heavy tails, extreme skew, and multimodality. In order to accommodate these scenarios, we



Figure 2: Toto architecture. Input time series of L steps (univariate example used for simplicity here) are first embedded using the patch embedding layer which produces. They then pass through the transformer stack, which contains F identical segments. Each segment of the transformer consists of one space-wise transformer block followed by N time-wise blocks. The flattened transformer outputs are projected to form the parameters of the Student-T mixture model (SMM) head. The final outputs are the forecasts for the input series, shifted P steps (the patch width) into the future.

introduce an even more flexible output likelihood in Section 2.5 based on a Student-T mixture model (Peel & McLachlan, 2000).

2 MODEL ARCHITECTURE

207 2.1 TRANSFORMER DESIGN

We build upon the ideas discussed above to define a novel architecture that efficiently models multivariate time series data.

Transformer models for time series forecasting have variously used encoder-decoder (Zhou et al., 2020; Wu et al., 2021; Ansari et al., 2024), encoder-only (Nie et al., 2023; Woo et al., 2024; Liu et al., 2024), and decoder-only architectures (Rasul et al., 2023; Das et al., 2024). For Toto, we employ a decoder-only architecture (Fig. 2). Decoder architectures have been shown to scale well (Radford & Narasimhan, 2018; Radford et al., 2019), and allow for arbitrary prediction horizons. The causal next-patch prediction task also simplifies the pre-training process.



Figure 3: The patch embedding takes as input a multivariate time series of M variates by L time steps. It divides each variate along the time dimension into patches of size P and projects these linearly into an embedding space of latent dimension D. This results in an output of size $M \times \frac{L}{D} \times D$ which is fed to the transformer decoder.

244

245

238

239

240

We utilize techniques demonstrated to yield performance and efficiency improvements in contemporary transformer literature, including pre-normalization (Xiong et al., 2020), RMSNorm (Zhang & Sennrich, 2019), and SwiGLU feed-forward layers (Shazeer, 2020).

- 246 247
- 248

253 254

257

260

2.2 INPUT/OUTPUT SCALING

249 As in other time series models, we perform instance normalization on input data before passing it 250 through the patch embedding, in order to make the model generalize better to inputs of different scales (Kim et al., 2022). We scale the inputs to have zero mean and unit standard deviation. The 251 output predictions are then rescaled back to the original units. 252

2.3 INPUT EMBEDDING

Time series transformers in the literature have used various approaches for creating input embed-256 dings. We use non-overlapping patch embeddings (Cordonnier et al., 2020; Dosovitskiy et al., 2021; Nie et al., 2023) (Fig. 3) of size P = 32, to project input time-series containing L = 4096 points to 258 embeddings of size $128 \times D$ per variate, where D = 512 is the embedding dimension. 259

261 2.4 ATTENTION MECHANISM

262 Observability metrics are often high-cardinality, multivariate time series. Therefore, we designed 263 our model to natively handle multivariate forecasting by analyzing relationships both in the time 264 dimension ("time-wise" interactions) and in the channel dimension ("space-wise" interactions). 265

266 In order to model both space and time-wise interactions, we adapt the traditional multi-head attention architecture (Vaswani et al., 2017) from one to two dimensions. We follow the intuition that 267 for many time series, the time relationships are more important or predictive than the space relation-268 ships. As evidence, we observe that even models that completely ignore space-wise relationships 269 (such as PatchTST (Nie et al., 2023) and TimesFM (Das et al., 2024)) can still achieve competitive performance on multivariate datasets. However, other studies (e.g. Woo et al. (2024)) have shown
 clear benefit to including space-wise attention in ablation studies.

We therefore propose a novel variant of factorized attention, which we call "Proportional Factorized Space-Time Attention." We use a mixture of alternating space-wise and time-wise attention blocks. As a configurable hyperparameter, we can change the ratio of time-wise to space-wise blocks, thus allowing us to devote more or less compute budget to each type of attention. For our base model, we selected a configuration with one space-wise attention block for every two time-wise blocks. This method allows for reduced computational complexity when compared to a traditional attention scheme (see Section A.1).

279 280

281

2.5 PROBABILISTIC PREDICTION

282 In order to produce probabilistic forecasts across the wide range of output distributions present in 283 observability data, we employ a method based on Gaussian mixture models (GMMs), which can approximate any density function (Goodfellow et al., 2016). We find that fitting GMMs leads to numer-284 ical instability in training, so we utilize a Student-T mixture model (SMM) of K distributions. This 285 model robustly generalizes GMMs (Peel & McLachlan, 2000), and has previously shown promise 286 for modeling heavy-tailed financial time series (Meitz et al., 2018; WONG et al., 2009). The model 287 predicts K Student-T distributions (where K is a hyperparameter) for each time step, as well as a 288 learned weighting. Formally, the SMM is defined by: 289

292 293 294

295

301

302

303

304

308

310 311 $p(x) = \sum_{k=1}^{K} \pi_k t(x \mid \mu_k, \tau, \nu_k)$ (1)

where $\pi_{k \in K}$ are nonnegative mixing coefficients which sum to 1 for the kth Student's t-distribution t with ν_k degrees of freedom, mean μ_k , and scale σ_k . $t(x \mid \mu, \sigma, \nu)$ is defined as:

$$t(x \mid \mu, \tau, \nu) = \frac{\Gamma\left(\frac{\nu+d}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)(\nu\pi)^{d/2}|\tau|^{1/2}} \left(1 + \frac{1}{\nu}(x-\mu)^T\tau^{-1}(x-\mu)\right)^{-\frac{\nu+d}{2}}$$
(2)

When we perform inference, we draw samples from the mixture distribution at each timestamp, then feed each sample back into the decoder for the next prediction. This allows us to produce prediction intervals at any quantile, limited only by the number of samples; for more precise tails, we can choose to spend more computation on sampling (Fig. 4).

As a decoder-only model, Toto is pre-trained on the next-patch prediction task. We minimize the negative log-likelihood of the next predicted patch with respect to the distribution output of the model, defined by the objective function:

$$\text{NLL} = -\log\left(\sum_{k=1}^{16} \pi_k t(x \mid \mu_k, \tau, \nu_k)\right)$$
(3)

Additionally, we utilize a dual softmax function on output logits for the mixing coefficients (Cheng et al., 2021), which has been demonstrated to improve training stability with highly heterogeneous data.

We train the model using the AdamW optimizer (Loshchilov & Hutter, 2019). The hyperparameters used for Toto are detailed in Table A1, with 103 million total parameters. In Section A.2, we perform an ablation study on the impact of various model components.

319

3 TRAINING DATA

320 321

We trained Toto with a dataset of approximately one trillion time series points. Of these, roughly three-quarters are anonymous observability metrics from an observability platform. The remaining points come from the LOTSA dataset (Woo et al., 2024), a compilation of publicly-available time



Figure 4: Example of Toto's 96-step zero-shot forecasts on the ETTh1 dataset, showing multivariate probabilistic predictions. Solid lines represent ground truth, dashed lines represent median point forecasts, and shaded regions represent 95% prediction intervals.

series datasets across many different domains. Additionally, we include synthetically generated time series data which we found to improve model performance.

3.1 **OBSERVABILITY DATASET**

346 An observability platform ingests more than a hundred trillion events per day. However, much of 347 this data is sparse, noisy, or too granular or high in cardinality to be useful in its raw form. To 348 curate a high-quality dataset for efficient model training, we sample queries based on quality and 349 relevance signals from dashboards, monitor alerts, and notebooks. This provides a strong signal that 350 the data resulting from these queries is of critical importance and sufficient quality for observability 351 of real-world applications. 352

Observability metrics are accessed using a specialized query language supporting filters, group-bys, 353 time aggregation, and various transformations and postprocessing functions. We consider groups 354 returned from the same query to be related variates in a multivariate time series (Fig. A1). After we 355 retrieve the query results, we discard the query strings and group identifiers, keeping only the raw 356 numeric data. 357

Handling this vast amount of data requires several preprocessing steps to ensure consistency and 358 quality. We describe the details of preprocessing and data augmentation in Section A.3.1. 359

360 361 362

363

337

338

339 340 341

342

343 344 345

3.2 SYNTHETIC DATA

We use a synthetic data generation process similar to TimesFM (Das et al., 2024) to supplement our 364 training datasets, improving the diversity of the data and helping to teach the model basic structure. The procedure used to generate synthetic data is detailed in Section A.3.2. 365

366 367

4 RESULTS

368 369

To evaluate predictions, we sequentially divide a time series into context and forecast segments. 370 We input the context segment into Toto and autoregressively generate output patches by sampling 371 from the Student-T mixture model distribution. We forecast a number of steps equal to the nearest 372 multiple of the patch size, then truncate the predictions to the desired length. In order to keep 373 inference time consistent, we vary the number of samples generated based on the cardinality and 374 length of the dataset, with a minimum of 100 samples. We take the median sample at each time 375 step as the final point prediction. This prediction is then compared against the ground-truth forecast 376 segment for evaluation. 377

We report experimental results for a pre-trained Toto model in Section 4.1 and Section 4.2.

				Zero Sho	ot					Full Shot				
Dataset	Metric	Toto	Moirai _{Small}	Moirai _{Base}	Moirai _{Large}	TimesFM*	iTransformer	TimesNet	PatchTST	Crossformer	TiDE	DLinear	SCINet	FEDformer
ETTh1	MAE MSE	0.389 0.363	$\frac{0.424}{0.400}$	0.438 0.434	0.469 0.510	0.426	0.448 0.454	0.450 0.458	0.455 0.469	0.522 0.529	0.507 0.541	0.452 0.456	0.647 0.747	0.460 0.440
ETTh2	MAE MSE	0.261 0.170	0.379 <u>0.341</u>	0.382 0.345	$\frac{0.376}{0.354}$	0.410	0.407 0.383	0.497 0.414	0.407 0.387	0.684 0.942	0.550 0.611	0.515 0.559	0.723 0.954	0.449 0.437
ETTm1	MAE MSE	0.375 0.372	0.409 0.448	0.388 0.381	0.389 0.390	0.388	0.410 0.407	0.406 0.400	0.400 0.387	0.495 0.513	0.419 0.419	0.407 0.403	0.481 0.486	0.452 0.448
ETTm2	MAE MSE	0.319 0.272	0.341 0.300	0.321 0.272	0.320 0.276	0.334	0.332 0.288	0.333 0.291	0.326 0.281	0.611 0.757	0.404 0.358	0.401 0.350	0.537 0.571	0.349 0.305
Electricity	MAE MSE	0.246 0.157	0.320 0.233	0.274 0.188	0.273 0.188	-	0.270 0.178	0.295 0.193	0.304 0.216	0.334 0.244	0.344 0.252	0.300 0.212	0.365 0.268	0.327 0.214
Weather	MAE MSE	0.284 0.256	0.267 0.242	0.261 0.238	0.275 0.259	-	0.278 0.258	0.287 0.259	0.281 0.259	0.315 0.259	0.320 0.271	0.317 0.265	0.363 0.292	0.360 0.309
Mean	MAE MSE	0.312 0.265	0.357 0.328	0.341 0.315	0.350 0.330	-	0.357 0.328	0.378 0.336	0.362 0.333	0.493 0.541	0.424 0.409	0.399 0.374	0.519 0.533	0.400 0.359

Table 1: Comparison of different models with Toto on the LSF benchmark datasets. Results are averaged across prediction lengths of 96, 192, 336, and 720 steps. For Toto, we use a stride of 512 steps and a historical context window of 512 steps. For other models, we use the results reported in Woo et al. (2024) and Das et al. (2024). Metrics for each prediction length are available in Table A3. *TimesFM only reports values for MAE on ETTh1, ETTh2, ETTm1, and ETTm2. Key: **Best results**, <u>Second-best results</u>.

393 394 395

396 397

388

389

390

391

392

4.1 LSF BENCHMARKS

To assess general-purpose time series forecasting performance, we use the Long Sequence Fore-398 casting (LSF) benchmark datasets (ETTh1, ETTh2, ETTm1, ETTm2, Electricity, and Weather) (Wu 399 et al., 2021). For Toto, we used a historical context window of 512 time steps and took the median of 400 200 samples. Following standard practice, we report normalized Mean Absolute Error (MAE) and 401 Mean Squared Error (MSE), fitted on a training split, in order to be able to compare performance 402 across different datasets. We compared Toto's performance with the reported results of other recent 403 zero-shot foundation models (Woo et al., 2024; Das et al., 2024), as well as full-shot time series 404 forecasting models (Liu et al., 2024; Wu et al., 2023; Nie et al., 2023; Zhang & Yan, 2023; Das 405 et al., 2023; Zeng et al., 2023; LIU et al., 2022; Zhou et al., 2022). We evaluate with forecast lengths of 96, 192, 336, and 720 time steps, in sliding windows with stride 512, and average the results. We 406 display these results in Table 1. 407

Toto demonstrates exceptional performance across a variety of benchmark datasets, excelling in zero-shot scenarios. In the LSF datasets, Toto consistently outperforms other models in terms of MAE and MSE. For example, on the ETTh1 dataset, Toto achieves an MAE of 0.389 and an MSE of 0.363, outperforming all zero-shot models, including the previously reported Moirai series and TimesFM. Macro-averaging across the six LSF datasets, Toto achieves an MAE of 0.312 and MSE of 0.265, again exceeding Moirai's reported zero-shot performance as well as the reported performance of the full-shot models.

While Toto generally excels, there are areas where its performance is closely matched by other models. In full-shot scenarios, models like PatchTST, Crossformer, and FEDformer show competitive
results. For example, on the Electricity dataset, while Toto achieves a leading zero-shot MAE of
0.246 and MSE of 0.157, iTransformer and TimesNet also show strong performance, indicating that
these models can catch up when additional training data is available.

420 421

422

4.2 Observability benchmark

We created a benchmark using anonymous observability data to assess performance across various observability metrics. To ensure a representative and realistic sample, we sampled data based on quality and relevance signals from dashboards, monitor alerts, and notebooks. This benchmark comprises 983,994 data points from 82 distinct multivariate time series, encompassing 1,122 variates.

We analyzed summary statistics of the series in our benchmark to identify characteristics that make observability time series challenging to forecast. The categories and their definitions are as follows:

430 431

• **Sparse:** Series with a low density of observations, indicating infrequent recording of data or rare events.

432	Metric	Toto	$Chronos\text{-}T5_{Tiny}$	$Chronos\text{-}T5_{Mini}$	Chronos-T5 _{Small}	$Chronos\text{-}T5_{Base}$	Chronos-T5 _{Large}	Moirai _{Small}	Moirai _{Base}	Moirai _{Large}	TimesFM
433	SMAPE	0.672	0.809	0.788	0.800	0.796	0.805	0.808	0.742	0.736	1.246
434	sMdAPE	0.318	0.406	0.391	0.401	0.393	0.396	0.418	0.370	0.365	0.639

Table 2: Performance of Toto and other zero-shot models on the observability benchmark dataset. Key: **Best results**, <u>Second-best results</u>.

- Extreme right skew: Series with a distribution heavily skewed to the right, characterized by a few very high values and many lower values.
- **Seasonal:** Series exhibiting regular and recurring patterns, often linked to daily, weekly, or yearly cycles.
- Flat: Series with minimal variability, showing little to no change over time.

445 To assess the prediction of other zero-shot models on the observability Benchmark, we follow sam-446 pling procedures delineated in their respective manuscripts. In short, for Chronos models, we gen-447 erate 20 samples and take the median prediction. For Moirai models, we take the median of 100 448 samples and set the patch size to "auto". TimesFM only produces point predictions of the mean, so 449 we use those directly. Since TimesFM and Chronos only support univariate forecasting, we process 450 each variate independently. Moirai, on the other hand, like Toto, makes joint predictions for each 451 group of related variates. For Toto, we utilize the same evaluation procedure we used on the LSF benchmarks. 452

- The relative proportion of these cases are displayed in Table A4. The evaluation results (Table 2) demonstrate that Toto outperforms the other models.
- Because observability data can have extreme variation in both magnitude and dispersion, we select
 symmetric mean absolute percentage error (sMAPE) as a scale-invariant performance metric (Armstrong, 1985). We also report symmetric median absolute percentage error (sMdAPE), a robust
 version of sMAPE (Hyndman & Koehler, 2006) that minimizes the influence of the extreme outliers
 present in observability data. With the lowest sMAPE of 0.672 and sMdAPE of 0.318, Toto proves
 to be the most accurate for forecasting observability time series data.
- These results suggest that current open datasets may not provide sufficient information to extrapolate
 to the specific nuances of observability data, highlighting the importance of training on more relevant
 data as demonstrated by Toto's superior performance.
- 465 466

467

435

436

437 438 439

440 441

442

443

444

5 CONCLUSIONS

Toto demonstrates state-of-the-art performance across both public and proprietary benchmarks. By
 leveraging a proportional factorized attention mechanism and a Student-T mixture model, Toto
 achieves impressive results in both zero-shot and full-shot settings, showcasing its scalability and
 flexibility in handling complex, high-dimensional data.

472 Despite its success, there are areas where further refinement is possible. Future work could involve integrating additional input modalities or exploring new attention mechanisms to enhance scalability and accuracy.
 475

With its demonstrated robustness and ability to manage observability data at scale, Toto not only advances time series forecasting but also opens new pathways for real-time system monitoring and infrastructure optimization, setting the stage for further innovations in the field.

479 480

481

6 IMPACT STATEMENT

In developing Toto, we followed a structured approach to ensure responsible development, focusing
 on identifying, assessing, and mitigating potential risks associated with the use of our model. Given
 that Toto is not intended for mass distribution and specifically generates time series forecasts for
 observability data, the potential harms are considerably lower compared to more general-purpose
 models. Our primary focus was ensuring the accuracy, reliability, and security of the forecasts

generated by Toto, which are crucial for maintaining and optimizing infrastructure and application
 performance.

We carefully analyze the potential for Toto to produce incorrect or misleading forecasts that could impact decision-making processes in critical systems. Additionally, we consider the implications of Toto's performance across diverse datasets, ensuring it can generalize well without introducing significant errors.

7 FUTURE DIRECTIONS

Many exciting areas of exploration remain for further study. Some future research questions that particularly intriguing include:

- Multi-modal inputs: Incorporate additional input modalities such as query metadata and captions to enhance forecast performance.
- Autonomous troubleshooting agents: Creating AI agents for troubleshooting and incident response by integrating modality-specific foundation models like Toto to improve their reasoning and planning capabilities with telemetry data.
- **Conversational interfaces:** Align time series forecasting models with LLMs to develop conversational agents capable of interpreting and reasoning about time series data.
- **Model enhancements and scaling:** Explore ways to improve and scale model performance through optimizations such as new types of input embeddings, attention mechanisms, and examining alternative variate groupings to capture richer interactions.
- 509

526

527

528

493

494 495

496

497 498

499

500

501

504

505

506

507

510 **REFERENCES**

- Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Sundar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor, Jasper Zschiegner, Danielle C. Maddix, Hao Wang, Michael W. Mahoney, Kari Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Yuyang Wang. Chronos: Learning the language of time series, 2024. URL https://arxiv.org/abs/2403.07815.
- 516
 517
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
 518
- Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid.
 Vivit: A video vision transformer. In 2021 IEEE/CVF International Conference on Computer
 Vision (ICCV), pp. 6816–6826, 2021. doi: 10.1109/ICCV48922.2021.00676.
- Xingyi Cheng, Hezheng Lin, Xiangyu Wu, F. Yang, and Dong Shen. Improving video-text retrieval by multi-stream corpus alignment and dual softmax loss. *ArXiv*, abs/2109.04290, 2021. URL https://api.semanticscholar.org/CorpusID:237454570.
 - Charisios Christodoulos, Christos Michalakelis, and Dimitris Varoutas. Forecasting with limited data: Combining arima and diffusion models. *Technological Forecasting and Social Change*, 77: 558–565, 5 2010. ISSN 00401625. doi: 10.1016/j.techfore.2010.01.009.
- Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between selfattention and convolutional layers. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020. URL https: //openreview.net/forum?id=HJlnClrKPB.
- Abhimanyu Das, Weihao Kong, Andrew Leach, Shaan K Mathur, Rajat Sen, and Rose Yu.
 Long-term forecasting with tiDE: Time-series dense encoder. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=
 pCbC3aQB5W.
- Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
 time-series forecasting. In *Forty-first International Conference on Machine Learning*, 2024. URL
 https://openreview.net/forum?id=jn2iTJas6h.

540 541 542 543	Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In <i>International Conference on Learning Representations</i> , 2021. URL https:
544	//openreview.net/forum?id=YicbFdNTTy.
545	Pohert Fildes Michèle Hibon Spyros Makridakis and Nigel Meade, Generalising about univariate
546	forecasting methods: further empirical evidence International Journal of Forecasting 14:339-
547	358, 9 1998, ISSN 01692070, doi: 10.1016/S0169-2070(98)00009-0.
548	
549	Azul Garza and Max Mergenthaler-Canseco. Timegpt-1, 2023.
550	Ian Goodfellow, Voshua Bengio, and Aaron Courville. Deen Learning, MIT Press, 2016, http://
551	//www_deeplearningbook_org
552	, ,
555 555 555	Nate Gruver, Marc Anton Finzi, Shikai Qiu, and Andrew Gordon Wilson. Large language models are zero-shot time series forecasters. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023. URL https://openreview.net/forum?id=md68e8iZK1.
557 558	R. J Hyndman and A. B. Koehler. Another look at measures of forecast accuracy. <i>International Journal of Forecasting</i> , 22, 2006.
559	Deb I Hundman and Coorse Athonesenergies Foreservices Duinsinkes and Duratics OTarts 2nd
560	edition 2021 URL https://otexts.com/fpp3/
561	Curron, 2021. OKE https://ocexcs.com/ipps/.
562	Romain Ilbert, Ambroise Odonnat, Vasilii Feofanov, Aladin Virmaux, Giuseppe Paolo, Themis Pal-
563	panas, and Ievgen Redko. SAMformer: Unlocking the potential of transformers in time series
564	forecasting with sharpness-aware minimization and channel-wise attention. In <i>Forty-first Interna-</i>
565	id-9kL zI 50Pb2
566	IU-OKUZUJQDHZ.
567	Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Re-
568	versible instance normalization for accurate time-series forecasting against distribution shift. In
509	International Conference on Learning Representations, 2022. URL https://openreview.
570	net/forum?id=CGDAkQofCUp.
572	Ze Li, Qian Cheng, Ken Hsieh, Yingnong Dang, Peng Huang, Pankaj Singh, Xinsheng Yang, Qing-
573	wei Lin, Youjiang Wu, Sebastien Levy, and Murali Chintalapati. Gandalf: an intelligent, end-to-
574	end analytics service for safe deployment in cloud-scale infrastructure. In <i>Proceedings of the 17th</i>
575	Usenix Conference on Networked Systems Design and Implementation, NSDI 20, pp. 389–402, USA 2020 USENIX Association ISBN 0781020122127
576	USA, 2020. USEINIA Association. ISBN 9781939133137.
577	Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. A survey of transformers. CoRR,
578	abs/2106.04554,2021. URL https://arxiv.org/abs/2106.04554.
579	Minhao IIII Ailing Zang Muyi Chan Zhijian Yu Guyia IAI Lingna Ma and Ojang Yu
580	SCINet: Time series modeling and forecasting with sample convolution and interaction. In Al-
581	ice H. Oh. Alekh Agarwal. Danielle Belgrave. and Kvunghvun Cho (eds.). Advances in Neu-
582	ral Information Processing Systems, 2022. URL https://openreview.net/forum?id=
583	AyajSjTAzmg.
584	Vara Lin Tarana II. Hanna Zhana Haim We Oli - West Litte Mersel Million I
585	itransformer: Inverted transformers are effective for time series for acousting 2024 LIDI https://
200 507	//openreview.net/forum?id=JePfAI8fah.
589	
580	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In International Confer-
590	ence on Learning Representations, 2019. URL https://openreview.net/forum?id=
591	вкдекісді /.
592	Mika Meitz, Daniel P. A. Preve, and Pentti Saikkonen. A mixture autoregressive model based on
593	student's t-distribution. <i>Communications in Statistics - Theory and Methods</i> , 52:499 – 515, 2018. URL https://api.semanticscholar.org/CorpusID:73615847.

617

634

641 642

643

644

645

- Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. 2023. URL https://openreview.net/forum?id=Jbdc0vTOcol.
- 598 D. Peel and G.J. McLachlan. Robust mixture modelling using the t distribution. *Statistics and Computing*, 10(4):339–348, 2000.
- Alec Radford and Karthik Narasimhan. Improving language understanding by generative pretraining. 2018. URL https://api.semanticscholar.org/CorpusID:49313245.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.
 Language models are unsupervised multitask learners. 2019. URL https://api.
 semanticscholar.org/CorpusID:160025533.
- Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and Alexander Rives. Msa transformer. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8844–8856. PMLR, 18–24 Jul 2021. URL https://proceedings.
 mlr.press/v139/rao21a.html.
- Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Arian Khorasani, George Adamopoulos, Rishika Bhagwatkar, Marin Biloš, Hena Ghonia, Nadhir Hassen, Anderson Schneider, Sahil Garg, Alexandre Drouin, Nicolas Chapados, Yuriy Nevmyvaka, and Irina Rish. Lag-Ilama: Towards foundation models for time series forecasting. In *R0-FoMo:Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*, 2023. URL https://openreview.net/forum? id=jYluzCLFDM.
- David Salinas, Valentin Flunkert, Jan Gasthaus, and Tim Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. International Journal of Forecasting, 36:1181–1191, 2020. ISSN 0169-2070. doi: https://doi.org/10.1016/j.ijforecast.
 2019.07.001. URL https://www.sciencedirect.com/science/article/pii/ S0169207019301888.
- Noam Shazeer. Glu variants improve transformer, 2020. URL https://arxiv.org/abs/
 2002.05202.
- Simon Stevenson. A comparison of the forecasting ability of arima models. *Journal of Property Investment & Finance*, 25:223–240, 5 2007. ISSN 1463-578X. doi: 10.1108/14635780710746902.
- Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding, 2021.
- Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. A length-extrapolatable transformer. In ACL 2023, December 2022. URL https://www.microsoft.com/en-us/research/publication/a-length-extrapolatable-transformer/.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. volume 30. Curran Associates, Inc., 2017. URL https://papers.nips.cc/paper_files/paper/2017/ hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.
- C. S. WONG, W. S. CHAN, and P. L. KAM. A student t -mixture autoregressive model with applications to heavy-tailed financial data. *Biometrika*, 96(3):751–760, 2009. ISSN 00063444, 14643510. URL http://www.jstor.org/stable/27798861.
 - Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo. Unified training of universal time series forecasting transformers. 2024. URL https: //openreview.net/forum?id=Yd8eHMY1wz.
- Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. 2021. URL https: //openreview.net/forum?id=J4gRj6d5Qm.

648 649 650	Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Temporal 2d-variation modeling for general time series analysis. In <i>International Conference on Learning Representations</i> , 2023.
652 653 654	Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. On layer normalization in the transformer architecture, 2020. URL https://openreview.net/forum?id=B1x8anVFPr.
655 656 657 658	Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series forecasting? <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 37(9):11121–11128, Jun. 2023. doi: 10.1609/aaai.v37i9.26317. URL https://ojs.aaai.org/index.php/ AAAI/article/view/26317.
659 660 661 662	Biao Zhang and Rico Sennrich. Root Mean Square Layer Normalization. In <i>Advances in Neural</i> <i>Information Processing Systems 32</i> , Vancouver, Canada, 2019. URL https://openreview. net/references/pdf?id=S1qBAf6rr.
663 664 665	Yunhao Zhang and Junchi Yan. Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=vSVLM2j9eie.
666 667 668 669	Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wan Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. 2020. URL https://api.semanticscholar.org/CorpusID:229156802.
670 671 672 673 674	Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting. In <i>Proc. 39th International</i> <i>Conference on Machine Learning (ICML 2022)</i> , 2022.
675 676 677	
678 679 680	
681 682 683	
684 685 686	
687 688 689	
690 691 692	
693 694 695	
696 697	
698 699 700	

APPENDIX А

A.1 MODEL ARCHITECTURE

After the patchwise embedding layer, we have inputs of shape $\mathbf{X} \in \mathbb{R}^{B \times M \times \frac{L}{P} \times D}$, where B is the batch dimension, M is the number of variates per batch item, $\frac{L}{D}$ is time steps divided by patch width, and D is the model embedding dimension.

Time-wise attention. We parallelize along the time dimension by reshaping the input tensor:

 $\mathbf{X} \in \mathbb{R}^{B \times M \times \frac{L}{P} \times D} \to \mathbf{X}_{\text{time}} \in \mathbb{R}^{(B \times M) \times \frac{L}{P} \times D}$

This allows for attention to be calculated independently in parallel per variate, giving a complexity of:

$$\mathcal{O}(M \times (\frac{L}{P})^2 \times D)$$

In the time-wise attention blocks, we use causal masking and rotary positional embeddings (Su et al., 2021) with XPOS (Sun et al., 2022) in order to autoregressively model time-dependent features.

Space-wise attention. We similarly parallelize along the time dimension by reshaping the input tensor:

$$\mathbf{X} \in \mathbb{R}^{B \times M \times \frac{L}{P} \times D} \to \mathbf{X}_{\text{space}} \in \mathbb{R}^{(B \times \frac{L}{P}) \times M \times D}$$

We calculate attention in parallel for each time step, with complexity:

$$\mathcal{O}(\frac{L}{P} \times M^2 \times D)$$

In the space-wise blocks, we use full bidirectional attention (without causal masking) in order to preserve permutation invariance of the covariates, with a block-diagonal ID mask to ensure that only related variates attend to each other. This masking allows us to pack multiple independent multivariate time series into the same batch, in order to improve training efficiency and reduce the amount of padding.

Computational complexity. Each transformer block in our model contains N timewise attention layers and 1 spacewise layer. The complexity for full self-attention over N + 1 layers, where interactions can occur across all variates and sequence positions, would be of complexity:

$$\mathcal{O}\left((N+1) \times M^2 \times \left(\frac{L}{P}\right)^2 \times D\right)$$
 (A1)

(A2)

This reflects the quadratic dependence on both the sequence length $\frac{L}{P}$ and the variate dimension M, with linear dependence on the embedding dimension D. However, by utilizing factorized attention, we can reduce the computational complexity of the attention calculation to:

- $\mathcal{O}\left(N \times M \times \left(\frac{L}{P}\right)^2 \times D + \frac{L}{P} \times M^2 \times D\right) =$

$$\mathcal{O}\left(D \times \frac{L}{P} \times M \times \left(N \times \frac{L}{P} + M\right)\right)$$

We demonstrate that factorized space-wise attention is asymptotically smaller in computational complexity than full self-attention (see Equation A1 and Equation A2). When comparing a model with full self-attention, we can assume N and D are fixed. Therefore:

$$\mathcal{O}\left(M\times \left(\frac{L}{P}\right)^2 + \frac{L}{P}\times M^2\right) < \mathcal{O}\left(M^2\times \left(\frac{L}{P}\right)^2\right)$$

which reduces to:

765 766

767 768 769

770

771 772 773

774

Thus, by factorizing attention into time-wise and space-wise components, the computational complexity is reduced, especially for large numbers of variates M or long sequences $\frac{L}{D}$, making it more

 $\mathcal{O}\left(M + \frac{L}{P}\right) < \mathcal{O}\left(M \times \frac{L}{P}\right)$

A.1.1 HYPERPARAMETERS

scalable than full self-attention.

775 Hyperparameter Value 776 **Embedding Dimension** 512 777 MLP Dimension 2048 778 # Layers 24 779 8 # Heads 780 32 # Variates 781 (0.9, 0.95) (β_1, β_2) 782 Weight Decay 0.01 Spacewise Layer Cadence 3 783 32 Patch Size 784 # Student-T Mixture Model Components 16 785 0.001 Initial Learning Rate 786 Annealing Schedule Cosine 787 Batch Size 192 788 Warmup Steps 5000 789 **Total Train Steps** 193000

Table A1: Hyperparameters for Toto

A.2 ABLATIONS

In this ablation study, we compare several versions of the Toto model using Negative Log Likelihood
 (NLL) loss on the validation set of our observability dataset. In addition to the full Toto model, we
 train separate variants with:

799 800

801

802 803

804

791

792 793 794

- 1. No space-wise attention (Time-wise Attention layers only)
- 2. No Student-T mixture model (instead, we replace the output with a single Student-T distribution)
 - 3. No observability data (instead, we train only on the full LOTSA dataset with synthetic data)

All models (except the "no observability data" model) were trained on a scaled down dataset with 620B points, with the number of training steps proportionally reduced to 117,000 steps. For each model, we report the NLL at its minimum during training and present the relative performance as a percentage decrease in comparison to the full Toto model. Table A2 presents the performance of each model variant, showing the percentage increase in NLL relative to the full Toto model (lower percentages indicate worse performance).

810	Model	NLL (% Increase Relative to Toto)
812	Toto	(baseline) 0%
813	No Space-wise Attention	4.37%
814	Single Student-T	11.48%
815	No Observability Data	14.21%

Table A2: Percentage increase in NLL relative to the full Toto model.

We observe that the full Toto model achieves the lowest NLL at its best validation point, serving as the baseline. The "No Space-wise Attention" variant shows a 4.37% increase in NLL, while the "Single Student-T" and "No observability Data" variants show larger decreases in performance, with NLL increases of 11.48% and 14.21%, respectively. These results indicate that space-wise attention, the Student-T mixture model, and the inclusion of observability -specific data are essential for optimal model performance. The percentage differences highlight the impact of these components on the model's ability to effectively model the underlying data distribution.

826 827

828

- A.3 TRAINING DATA PREPROCESSING
- A.3.1 OBSERVABILITY DATASET

Initially, we apply padding and masking techniques to align the series lengths, making them divisible
by the patch stride. This involves adding necessary left-padding to both the time series data and the
ID mask, ensuring compatibility with the model's requirements.

834 Various data augmentations are employed to enhance the dataset's robustness. We introduce random 835 time offsets to prevent memorization caused by having series always align the same way with the 836 patch grid. After concatenating the observability and LOTSA datasets for training, we also imple-837 ment a variate shuffling strategy to maintain diversity and representation. Specifically, we randomly 838 combine variates from either observability, LOTSA, and/or synthetic data with a probability of 839 10%, thus creating new, diverse combinations of data points. To sample the indices when mixing in this fashion, we employ a normal distribution with a standard deviation of 1000, favoring data points 840 that were closer together in the original datasets. This Gaussian sampling ensures that, while there is 841 a preference for adjacent data points, significant randomness is introduced to enhance the diversity 842 of the training data. This approach improves the model's ability to generalize across different types 843 of data effectively. 844



Figure A1: Example metric query in the observability platform. The metric name (1) determines which metric is being queried. The filter clause (2) limits which contexts are queried, in this case restricting the query to the prod environment. The space aggregation (3) indicates that the average metric value should be returned for each unique combination of the group-by keys. The time aggregation (4) indicates that metric values should be aggregated to the average for each 60-second interval. The query results will be a multivariate time series with 1-minute time steps, and with separate individual variates for each unique service, datacenter tuple.

A.3.2 SYNTHETIC DATA

We simulate time series data through the composition of components such as piecewise linear trends, ARMA processes, sinusoidal seasonal patterns, and various residual distributions. We randomly combine five of these processes per variate, introducing patterns not always present in our real-world datasets. The generation process involves creating base series with random transformations, clipping extreme values, and rescaling to a specified range. By making synthetic data approximately 5% of our training dataset, we ensure a wide range of time-series behaviors are captured. This diversity exposes our models to various scenarios during training, improving their ability to generalize and effectively handle real-world data.

A.4 RESULTS

A.4.1 LSF BENCHMARKS

Dataset Prediction Lenge Matrix Model, and Data Tomosfield Tom	878						Zero Sho	t					Full Shot				
87.9 96 Mac 0.362 0.402 0.398 0.436 0.442 0.445 0.445 0.448 0.449 0.448 0.444 0.448 0.444 0.448 0.444 0.448 0.444 0.448 0.444 0.445 0.448 0.444 0.448 0.444 0.448 0.444 0.448 0.444 0.448 0.444 0.448 0.444	070	Dataset	Prediction Length	Metric	Toto	Moirai _{Small}	Moirai _{Base}	Moirai _{Large}	TimesFM	iTransformer	TimesNet	PatchTST	Crossformer	TIDE	DLinear	SCINet	FEDformer
980 pg NAE 0.386 0.436 0.436 0.436 0.445 0.474 0.472 0.472 0.472 0.473 0.474 0.488 0.421 0.471 0.474 0.488 0.421 0.471 0.473 0.488 0.421<	879		96	MAE	0.366	0.402	0.402	0.398	0.398	0.405	0.402	0.419	0.448	0.464	0.400	0.599	0.419
FTh: MSE 0.239 0.399 0.423 0.440 0.441 0.456 0.466 0.456 0.437 0.437 0.439 0.435 881 36 MSE 0.399 0.412 0.436 0.436 0.448 0.466 0.516 0.535 0.437 0.439 0.449 882 720 MSE 0.399 0.212 0.336 0.481 0.470 0.536 0.5	880		192	MAE	0.368	0.375	0.384	0.380	0.424	0.386	0.384	0.414	0.423	0.479	0.380	0.634	0.376
881	000	ETTh1		MSE	0.329	0.399	0.425	0.440	-	0.441	0.436	0.460	0.471	0.525	0.437	0.719	0.420
720 MAE 0.23 0.243 0.473 0.568 0.491 0.500 0.888 0.621 0.558 0.516 0.699 0.507 883 96 MAE 0.197 0.327 0.225 0.356 0.591 0.531 0.591 0.538 0.561 0.591 0.538 0.507 0.338 0.577 0.338 0.577 0.338 0.577 0.338 0.577 0.338 0.377 0.338 0.377 0.338 0.377 0.338 0.377 0.338 0.377 0.338 0.377 0.338 0.377 0.388 0.440 0.407 0.441 0.400 0.565 0.599 0.476 0.668 0.439 885 0.366 0.382 0.371 <t< th=""><th>881</th><th></th><th>336</th><th>MAE</th><th>0.399</th><th>0.429</th><th>0.450</th><th>0.474</th><th>0.436</th><th>0.458</th><th>0.469</th><th>0.466</th><th>0.546</th><th>0.515</th><th>0.459</th><th>0.659</th><th>0.465</th></t<>	881		336	MAE	0.399	0.429	0.450	0.474	0.436	0.458	0.469	0.466	0.546	0.515	0.459	0.659	0.465
MSE 0.419 0.410 0.470 0.705 - 0.503 0.521 0.500 0.653 0.594 0.519 0.836 0.507 883 96 MAE 0.093 0.231 0.237 0.232 0.336 0.344 0.344 0.344 0.344 0.334 0.333 0.707 0.338 884 ETTh2 336 MAE 0.137 0.137 0.137 0.137 0.137 0.137 0.137 0.138 0.414 0.400 0.414 0.400 0.414 0.400 0.656 0.509 0.477 0.889 0.437 885 70 MAE 0.131 0.237 0.127 0.428 0.442 0.438 0.877 0.528 0.417 0.488 0.451 0.446 0.763 0.679 0.878 0.888 0.446 886 0.385 0.412 0.435 0.334 0.338 0.327 0.438 0.437 0.445 0.436 0.447 0.441 <	000		720	MAE	0.424	0.444	0.473	0.568	0.445	0.491	0.500	0.488	0.621	0.558	0.516	0.699	0.507
883 96 MAE 0.197 0.334 0.327 0.232 0.356 0.349 0.348 0.548 0.440 0.387 0.621 0.377 884 PTD2 MAE 0.231 0.237 0.374 0.348 0.440 0.400 0.445 0.400 0.435 0.475 0.400 0.436 0.400 0.445 0.400 0.445 0.400 0.445 0.433 0.475 0.476 0.489 0.448 885 20.46 0.456 0.421 0.437 0.428 0.445 0.451 0.433 0.431 0.447 0.448 0.451 0.431 0.140 0.745 0.467 0.448 0.439 0.431 0.148 0.438 0.441 0.433 0.446 0.338 0.329 0.445 0.451 0.436 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 0.337 <t< th=""><th>882</th><th></th><th></th><th>MSE</th><th><u>0.419</u></th><th>0.413</th><th>0.470</th><th>0.705</th><th>-</th><th>0.503</th><th>0.521</th><th>0.500</th><th>0.653</th><th>0.594</th><th>0.519</th><th>0.836</th><th>0.506</th></t<>	882			MSE	<u>0.419</u>	0.413	0.470	0.705	-	0.503	0.521	0.500	0.653	0.594	0.519	0.836	0.506
bit of the second sec	883		96	MAE	0.197	0.334	0.327	0.325	0.356	0.349	0.374	0.348	0.584	0.440	0.387	0.621	0.397
884 ETh2 MSE 0.340 0.340 0.347 - 0.380 0.402 0.388 0.877 0.528 0.477 0.860 0.429 885 MAE 0.260 0.335 0.441 0.377 - 0.428 0.452 0.425 0.426 0.433 0.571 <	005		192	MAE	0.093	0.281	0.374	0.287	0.400	0.400	0.414	0.302	0.656	0.509	0.333	0.689	0.338
336 MAE 0.200 0.333 0.401 0.332 0.428 0.432 0.431 0.433 0.571 0.571 0.541 0.744 0.487 885 720 MAE 0.355 0.416 0.422 0.421 0.452 0.426 0.431 0.140 0.871 0.571 0.571 0.531 0.249 0.480 886 70 MAE 0.232 0.334 0.441 0.457 0.442 0.431 1.140 0.871 0.331 0.420 887 MSE 0.336 0.402 0.335 0.427 0.337 0.367 0.426 0.337 0.377 0.428 0.441 0.344 0.439 0.445 0.448 0.389 0.445 0.448 0.389 0.445 0.448 0.389 0.445 0.449 0.347 0.387 0.387 0.387 0.437 0.449 0.449 0.449 0.449 0.449 0.449 0.449 0.449 0.449 0.449 0.449 <th>884</th> <th>ETTh2</th> <th></th> <th>MSE</th> <th>0.135</th> <th>0.340</th> <th>0.340</th> <th>0.347</th> <th>-</th> <th>0.380</th> <th>0.402</th> <th>0.388</th> <th>0.877</th> <th>0.528</th> <th>0.477</th> <th>0.860</th> <th>0.429</th>	884	ETTh2		MSE	0.135	0.340	0.340	0.347	-	0.380	0.402	0.388	0.877	0.528	0.477	0.860	0.429
885 720 MAE 0.350 0.280 0.031 0.427 0.447 0.457 0.442 0.446 0.1043 0.446 0.043 0.446 0.043 0.446 0.043 0.446 0.043 0.446 0.043 0.446 0.043 0.446 0.043 0.446 0.043 0.446 0.043 0.0467 0.043 0.0462 0.446 0.043 0.0462 0.043 0.0461 0.043 0.0462 0.043 0.0446 0.043 0.038 0.328 0.038 0.038 0.0460 0.037 0.037 0.037 0.038 0.038 0.0450 0.0426 0.038 0.0460 0.0426 0.038 0.0461 0.0440 0.038 0.0461 0.0440 0.038 0.0426 0.038 0.0426 0.038 0.0461 0.0438 0.0446 0.038 0.0461 0.0485 0.0461 0.0485 0.0461 0.0485 0.0461 0.0485 0.0461 0.0485 0.0461 0.0485 0.0461 0.0485 0.0486 0.0485 0.0461 0.0397 0.420 0.411 0.410 0.358 0.4461 0.438 0.0485 0.0480 0.0461 0.0485 0.0480 0.0461 0.0485 0.0480 0.0461 0.0438			336	MAE	0.260	0.393	0.401	0.393	0.428	0.432	0.541	0.433	0.731	0.571	0.541	0.744	0.487
Nise 0.294 0.380 0.394 0.404 - 0.427 0.462 0.431 1.104 0.874 0.831 1.249 0.463 886 96 MAE 0.328 0.363 0.345 0.368 0.375 0.367 0.426 0.887 0.372 0.448 0.418 0.418 0.419 887 192 MAE 0.353 0.127 0.391 0.387 0.385 0.441 0.440 0.389 0.480 0.426 888 36 MAE 0.389 0.416 0.399 0.426 0.411 0.410 0.515 0.422 0.411 0.410 0.515 0.422 0.411 0.410 0.515 0.423 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.439 0.449 0.435 0.449 0.441 0.441 0.444 0.444 0.444	885		720	MAE	0.355	0.416	0.426	0.421	0.457	0.445	0.452	0.446	0.763	0.679	0.657	0.838	0.490
96 MAE 0.323 0.360 0.433 0.345 0.367 0.426 0.387 0.372 0.438 0.419 887 MSE 0.306 0.404 0.353 0.375 0.367 0.426 0.387 0.347 0.446 0.354 0.438 0.418 0.418 0.419 888 36 MAE 0.335 0.402 0.377 0.374 0.387 0.385 0.411 0.440 0.388 0.449 0.445 0.449 0.445 0.449 0.445 0.449 0.445 0.449 0.440 0.451 0.440 0.431 0.445 0.449 0.439 0.589 0.461 0.449 0.433 0.550 0.443 0.449 0.439 0.589 0.461 0.444 0.449 0.459 0.450 0.439 0.589 0.461 0.447 0.454 0.447 0.454 0.447 0.454 0.447 0.595 0.543 0.448 0.444 0.595 0.543 0.561				MSE	0.294	0.380	0.394	0.404	-	0.427	0.462	0.431	1.104	0.874	0.831	1.249	0.463
887 PTM1 N3E 0.306 0.404 0.335 0.335 0.337 0.338 0.329 0.404 0.364 0.345 0.418 0.379 888 AMAE 0.335 0.435 0.435 0.435 0.426 0.377 0.378 0.377 0.428 0.411 0.410 0.377 0.428 0.433 0.463 0.463 0.463 0.459 0.436 0.459 0.436 0.451 0.413 0.413 0.413	886		96	MAE	0.328	0.383	0.360	0.363	0.345	0.368	0.375	0.367	0.426	0.387	0.372	0.438	0.419
887 ETml 92 MAE 0.325 0.432 0.379 0.379 0.387 0.385 0.485 0.441 0.449 0.389 0.480 0.441 888 36 MAE 0.389 0.430 0.411 0.410 0.515 0.425 0.430 0.441 0.410 0.515 0.422 0.413 0.448 0.449 889 70 MAE 0.429 0.399 - 0.426 0.410 0.595 0.422 0.413 0.448 0.449 889 70 MAE 0.429 0.431 0.443 0.459 0.450 0.439 0.585 0.461 0.443 0.445 890 70 MAE 0.220 0.232 0.426 0.411 0.416 0.459 0.450 0.439 0.585 0.461 0.447 0.591 0.543 890 MAE 0.270 0.284 0.269 0.261 0.247 0.291 0.246 0.203 0.264	007			MSE	0.306	0.404	0.335	0.353		0.334	0.338	0.329	0.404	0.364	0.345	0.418	0.379
Bit Init 336 Mice 0.389 0.437 0.399 0.397 0.317 0.310 0.339 0.393 0.439 <th< th=""><th>887</th><th>ETTm1</th><th>192</th><th>MAE</th><th>0.353</th><th>0.402</th><th>0.379</th><th>0.380</th><th>0.374</th><th>0.391</th><th>0.387</th><th>0.385</th><th>0.451</th><th>0.404</th><th>0.389</th><th>0.450</th><th>0.441</th></th<>	887	ETTm1	192	MAE	0.353	0.402	0.379	0.380	0.374	0.391	0.387	0.385	0.451	0.404	0.389	0.450	0.441
888 mse 0.390 0.462 0.197 0.476 0.416 0.199 0.538 0.428 0.413 0.490 0.443 889 720 MAE 0.429 0.437 0.419 0.432 0.436 0.491 0.478 0.444 0.666 0.487 0.411 0.990 0.538 0.461 0.453 0.550 0.543 890 96 MAE 0.270 0.282 0.266 0.266 0.261 0.275 0.287 0.289 0.366 0.305 0.292 0.377 0.287 891 192 MAE 0.315 0.318 0.303 0.300 0.309 0.302 0.492 0.364 0.390 0.322 0.472 0.591 0.289 892 70 MAE 0.319 0.333 0.334 0.344 0.343 0.542 0.421 0.399 0.325 0.414 0.399 0.325 0.415 0.318 0.316 0.314 0.354 0.415 <td< th=""><th>000</th><th>E11m1</th><th>336</th><th>MAE</th><th>0.328</th><th>0.435</th><th>0.394</th><th>0.395</th><th>0.397</th><th>0.420</th><th>0.374</th><th>0.307</th><th>0.515</th><th>0.398</th><th>0.380</th><th>0.439</th><th>0.459</th></td<>	000	E11m1	336	MAE	0.328	0.435	0.394	0.395	0.397	0.420	0.374	0.307	0.515	0.398	0.380	0.439	0.459
889 720 MAE 0.429 0.417 0.436 0.459 0.450 0.447 0.474 0.595 0.543 890 96 MAE 0.270 0.282 0.269 0.261 0.263 0.291 0.377 0.287 0.270 0.187 0.270 0.284 0.202 0.291 0.371 0.286 0.261 0.277 0.287 0.287 0.278 0.284 0.399 0.264 0.390 0.284 0.399 0.269 0.261 0.399 0.269 0.261 0.399 0.269 0.261 0.399 0.291 0.355 0.333 0.341 0.341 0.343 0.542 0.422 0.421 892 70 MAE 0.374 0.419 <th0.327< th=""> <th0.399< th=""> 0.343<</th0.399<></th0.327<>	888			MSE	0.390	0.462	0.391	0.399	-	0.426	0.410	0.399	0.532	0.428	0.413	0.490	0.445
889 Mike 0.463 0.493 0.493 0.0478 0.478 0.474 0.0474	000		720	MAE	0.429	0.437	0.419	0.417	0.436	0.459	0.450	0.439	0.589	0.461	0.453	0.550	0.490
890 96 MAE 0.270 0.282 0.263 0.263 0.264 0.267 0.289 0.366 0.305 0.292 0.377 0.287 891 H2 MAE 0.315 0.318 0.303 0.309 0.309 0.302 0.492 0.364 0.392 0.484 0.392 0.482 0.364 0.390 0.302 0.492 0.364 0.390 0.302 0.492 0.364 0.390 0.302 0.492 0.364 0.390 0.302 0.492 0.364 0.390 0.269 0.261 0.247 0.291 0.284 0.399 0.232 0.433 0.542 0.421 0.290 0.284 0.399 0.235 0.318 0.332 0.343 0.542 0.421 0.355 0.521 0.775 0.415 0.415 0.415 0.415 0.403 0.400 1.042 0.522 0.775 0.412 893 MAE 0.374 0.411 0.335 0.315 0.315 0.	889			MSE	0.463	0.490	0.434	0.432	-	0.491	0.478	0.454	0.666	0.487	0.474	0.595	0.543
bit bit<	800		96	MAE	0.270	0.282	0.269	0.260	0.263	0.264	0.267	0.259	0.366	0.305	0.292	0.377	0.287
891 ETm2 336 MSE 0.247 0.247 0.247 0.249 0.241 0.141 0.290 0.284 0.399 0.299 0.268 0.364 0.341 0.349 0.349 0.349 0.349 0.349 0.349 0.341 0.349 0.344 0.351 0.353 0.542 0.422 0.422 0.591 0.356 0.356 0.331 0.331 0.335 0.542 0.422 0.571 0.366 0.366 892 720 MAE 0.354 0.415 0.355 0.332 0.377 0.366 0.415 0.372 0.412 0.408 0.402 1.730 0.558 0.554 0.928 0.345 0.345 893 96 MAE 0.212 0.239 0.242 - 0.412 0.408 0.402 1.730 0.558 0.554 0.308 0.308 0.308 0.327 0.237 0.377 0.237 0.377 0.237 0.377 0.237 0.377 0.237 <th>090</th> <th></th> <th>192</th> <th>MAE</th> <th>0.200</th> <th>0.205</th> <th>0.195</th> <th>0.189</th> <th>0 309</th> <th>0.180</th> <th>0.187</th> <th>0.175</th> <th>0.287</th> <th>0.207</th> <th>0.195</th> <th>0.286</th> <th>0.205</th>	090		192	MAE	0.200	0.205	0.195	0.189	0 309	0.180	0.187	0.175	0.287	0.207	0.195	0.286	0.205
336 MAE 0.319 0.355 0.333 0.334 0.349 0.348 0.351 0.343 0.452 0.422 0.472 0.591 0.356 892 720 MAE 0.374 0.419 0.377 0.386 0.415 0.401 0.403 0.400 1.042 0.524 0.529 0.735 0.415 893 96 MAE 0.374 0.419 0.376 0.415 0.401 0.403 0.400 1.042 0.524 0.525 0.534 0.415 893 96 MAE 0.212 0.239 0.325 0.415 0.418 0.400 1.042 0.374 0.525 0.541 0.355 0.415 0.418 894 Fleetricity 96 MAE 0.212 0.228 0.235 0.314 0.325 0.413 0.305 0.375 0.377 0.360 0.237 0.279 0.237 0.279 0.237 0.279 0.237 0.279 0.285 0.335	801	ETTm2		MSE	0.269	0.261	0.247	0.247	-	0.250	0.249	0.241	0.414	0.290	0.284	0.399	0.269
892 70 MAE 0.264 0.319 0.291 0.291 0.393 - 0.311 0.312 0.305 0.577 0.377 0.369 0.637 0.322 893 MAE 0.374 0.410 0.373 0.345 0.415 0.417 0.448 0.402 1.730 0.554 0.522 0.735 0.415 0.440 893 MAE 0.321 0.235 0.372 - 0.412 0.408 0.402 1.730 0.554 0.522 0.735 0.415 0.416 894 MAE 0.412 0.202 0.227 0.285 0.314 0.039 0.422 0.345 0.346 0.429 0.322 0.345 0.368 0.359 0.329 0.282 0.345 0.368 0.359 0.329 0.282 0.345 0.368 0.359 0.329 0.282 0.345 0.368 0.359 0.352 0.359 0.359 0.359 0.359 0.359 0.359 0.359	001		336	MAE	0.319	0.355	0.333	0.334	0.349	0.348	0.351	0.343	0.542	0.422	0.427	0.591	0.366
Visc 0.374 0.410 0.212 0.303 0.413 0.440 0.400 0.400 1.022 0.224 0.102 0.421 893 MSE 0.354 0.415 0.352 - 0.413 0.405 0.400 1.040 1.730 0.525 0.524 0.421 893 96 MAE 0.212 0.299 0.248 0.242 - 0.248 0.416 0.122 0.285 0.314 0.299 0.248 0.402 - 0.248 0.416 0.122 0.277 0.276 0.277 0.277 0.276 0.277 0.277 0.276 0.300 0.305 0.337 0.344 0.301 0.369 0.226 895 MSE 0.155 0.236	892		720	MSE	0.264	0.319	0.291	0.295	0.415	0.311	0.321	0.305	0.597	0.377	0.369	0.637	0.325
893 96 MAE 0.212 0.299	002		720	MSE	0.374	0.415	0.355	0.372	-	0.407	0.403	0.400	1.730	0.558	0.554	0.960	0.413
Nise 0.124 0.205 0.158 0.152 - 0.148 0.168 0.105 0.219 0.237 0.197 0.247 0.108 894 192 MAE 0.232 0.310 0.263 0.263 0.259 - 0.249 0.289 0.237 0.197 0.247 0.108 894 192 MAE 0.232 0.310 0.263 0.253 0.289 0.289 0.239 0.237 0.197 0.235 0.315 895 336 MAE 0.249 0.278 0.278 0.278 0.200 0.030 0.305 0.337 0.344 0.301 0.369 0.329 896 70 MAE 0.211 0.377 0.273 0.272 0.230 0.337 0.343 0.330 0.330 0.333 0.330 0.330 0.333 0.330 0.330 0.333 0.333 0.330 0.330 0.337 0.344 0.245 0.290 0.246 0.246 0.2	893		96	MAE	0.212	0.200	0.248	0.242		0.240	0.272	0.285	0.314	0.320	0.282	0.345	0.308
894 P12 MAE 0.233 0.210 0.263 0.263 0.255 0.289 0.289 0.322 0.300 0.285 0.315 0.210 895 MSE 0.138 0.220 0.174 0.171 - 0.162 0.184 0.199 0.231 0.236 0.928 0.235 0.315 0.210 895 MKE 0.249 0.323 0.278 0.278 - 0.129 0.305 0.337 0.344 0.300 0.369 0.329 0.246 0.399 0.226 0.290 0.214 0.279 0.210 0.305 0.337 0.344 0.300 0.309 0.231 0.246 0.249 0.290 0.266 0.249 0.290 0.266 0.249 0.290 0.266 0.246 0.240 0.290 0.266 0.250 0.241 0.270 0.229 0.226 0.226 0.226 0.226 0.226 0.226 0.226 0.226 0.226 0.226 0.226 0.226<			20	MSE	0.124	0.205	0.158	0.152	-	0.148	0.168	0.195	0.219	0.237	0.197	0.247	0.193
Electricity MSE 0.138 0.220 0.174 0.171 - 0.162 0.184 0.199 0.231 0.236 0.196 0.237 0.201 895 36 MAE 0.129 0.232 0.278 - 0.260 0.303 0.331 0.236 0.199 0.237 0.239 896 700 MAE 0.155 0.246 0.191 0.192 - 0.178 0.198 0.215 0.246 0.249 0.239 0.237 896 700 MAE 0.211 0.270 0.313 - 0.177 0.120 0.337 0.333 0.390 0.355 897 MAE 0.211 0.270 0.228 - 0.224 0.226 0.284 0.201 0.245 0.290 0.246 897 MAE 0.167 0.173 0.167 0.177 - 0.174 0.177 0.178 0.230 0.266 0.246 0.250 0.241 0.210 <	894		192	MAE	0.232	0.310	0.263	0.259	-	0.253	0.289	0.289	0.322	0.330	0.285	0.355	0.315
895		Electricity	226	MSE	0.138	0.220	0.174	0.171	-	0.162	0.184	0.199	0.231	0.236	0.196	0.257	0.201
896 70 MAE 0.291 0.347 0.397 0.313 - 0.317 0.320 0.337 0.363 0.373 0.338 0.339 0.358 896 MSE 0.211 0.270 0.225 0.220 0.235 0.220 0.236 0.280 0.284 0.248 0.249 0.246 897 MSE 0.180 0.173 0.167 0.177 - 0.214 0.220 0.216 0.250 0.244 0.240 0.244 0.240 0.246 0.241 0.270 0.238 0.246 0.241 0.270 0.217 0.174 0.174 0.177 0.158 0.200 0.214 0.220 0.216 0.276 0.230 0.214 0.210 0.217 0.173 0.167 0.174 0.177 0.158 0.202 0.216 0.236 0.340 0.336 0.336 0.336 0.336 0.336 0.336 0.336 0.336 0.336 0.336 0.336 0.336 0.336 </th <th>895</th> <th></th> <th>330</th> <th>MSE</th> <th>0.155</th> <th>0.236</th> <th>0.191</th> <th>0.192</th> <th></th> <th>0.178</th> <th>0.198</th> <th>0.215</th> <th>0.246</th> <th>0.249</th> <th>0.209</th> <th>0.369</th> <th>0.214</th>	895		330	MSE	0.155	0.236	0.191	0.192		0.178	0.198	0.215	0.246	0.249	0.209	0.369	0.214
896 MSE 0.211 0.270 0.229 0.236 - 0.225 0.220 0.266 0.280 0.284 0.249 0.246 897 96 MAE 0.223 0.212 0.203 0.203 0.214 0.220 0.218 0.230 0.261 0.255 0.306 0.299 192 MAE 0.267 0.250 0.217 0.174 0.172 0.17 0.188 0.230 0.210 0.210 0.218 0.230 0.202 0.96 0.321 0.217 0.217 0.178 0.167 0.174 0.172 0.17 0.188 0.230 0.211 0.214 0.219 - 0.254 0.259 0.277 0.298 0.241 0.249 - 0.254 0.259 0.277 0.298 0.242 0.237 0.216 0.299 0.276 0.296 0.241 0.249 0.276 0.296 0.296 0.242 0.237 0.242 0.237 0.242 0.237 0.242			720	MAE	0.291	0.347	0.307	0.313	-	0.317	0.320	0.337	0.363	0.373	0.333	0.390	0.355
96 MAE 0.223 0.212 0.203 0.203 0.214 0.220 0.218 0.230 0.261 0.255 0.306 0.296 897 MSE 0.180 0.173 0.167 0.177 - 0.174 0.127 0.18 0.231 0.221 0.218 0.230 0.261 0.255 0.306 0.296 192 MAE 0.267 0.256 0.276 0.271 0.289 0.277 0.28 0.296 0.221 0.219 0.255 0.216 0.276 898 Weather MSE 0.257 0.256 0.276 0.296 0.306 0.297 0.326 0.242 0.237 0.261 0.276 336 MAE 0.251 0.256 0.277 - 0.296 0.306 0.297 0.335 0.335 0.378 0.380 MAE 0.256 0.277 - 0.278 0.278 0.272 0.287 0.283 0.335 0.378 0.380 <th>896</th> <th></th> <th></th> <th>MSE</th> <th>0.211</th> <th>0.270</th> <th>0.229</th> <th>0.236</th> <th>-</th> <th>0.225</th> <th><u>0.220</u></th> <th>0.256</th> <th>0.280</th> <th>0.284</th> <th>0.245</th> <th>0.299</th> <th>0.246</th>	896			MSE	0.211	0.270	0.229	0.236	-	0.225	<u>0.220</u>	0.256	0.280	0.284	0.245	0.299	0.246
MSE 0.180 0.173 0.167 0.172 0.174 0.172 0.177 0.188 0.202 0.216 0.217 0.174 0.172 0.177 0.188 0.202 0.196 0.221 0.217 898 MAE 0.267 0.250 0.224 0.214 0.249 0.254 0.261 0.259 0.298 0.296 0.340 0.336 898 MAE 0.235 0.216 0.209 0.219 - 0.224 0.210 0.225 0.206 0.242 0.237 0.261 0.376 809 MAE 0.252 0.226 0.226 0.208 0.297 0.238 0.340 0.380 809 700 MAE 0.256 0.277 - 0.278 0.278 0.277 0.281 0.418 0.438 0.380 0.380 809 700 MAE 0.356 0.322 0.323 0.335 0.335 0.335 0.336 0.349 0.380	007		96	MAE	0.223	0.212	0.203	0.208	-	0.214	0.220	0.218	0.230	0.261	0.255	0.306	0.296
152 INTE 0.207 0.204 0.201 0.204 0.201 0.205 0.217 0.298 0.296 0.390 0.300 898 MSE 0.235 0.216 0.209 0.219 - 0.221 0.219 0.225 0.216 0.200 0.276 0.291 0.219 0.225 0.216 0.276 0.292 0.297 0.305 0.335 0.337 0.236 0.376 899 MSE 0.252 0.252 0.256 0.276 0.297 - 0.298 0.297 0.335 0.335 0.337 0.339 0.339 899 720 MAE 0.325 0.322 0.322 0.323 0.350 - 0.449 0.359 0.418 0.386 0.381 0.427 0.428	091		102	MSE	0.180	0.173	0.167	0.177	-	0.174	0.172	0.177	0.158	0.202	0.196	0.221	0.217
336 MAE 0.291 0.226 0.292 - 0.296 0.306 0.297 0.335 0.335 0.335 0.335 0.336 0.389 MSE 0.252 0.256 0.277 - 0.296 0.306 0.278 0.237 0.283 0.309 0.389 899 720 MAE 0.356 0.322 0.323 0.350 - 0.439 0.339 0.386 0.381 0.428 0.428 0.428 0.418 0.436 0.438 0.428 <th>000</th> <th>Weather</th> <th>192</th> <th>MSE</th> <th>0.267</th> <th>0.230</th> <th>0.241</th> <th>0.249</th> <th></th> <th>0.234</th> <th>0.201</th> <th>0.239</th> <th>0.206</th> <th>0.298</th> <th>0.296</th> <th>0.261</th> <th>0.336</th>	000	Weather	192	MSE	0.267	0.230	0.241	0.249		0.234	0.201	0.239	0.206	0.298	0.296	0.261	0.336
MSE 0.252 0.260 0.256 0.277 - 0.278 0.280 0.272 0.287 0.283 0.309 0.339 720 MAE 0.356 0.322 0.323 0.350 - 0.349 0.359 0.348 0.418 0.386 0.381 0.427 0.428	030		336	MAE	0.291	0.282	0.276	0.292	-	0.296	0.306	0.297	0.335	0.335	0.335	0.378	0.380
U33 720 MAE 0.356 0.322 0.323 0.350 - 0.349 0.359 0.348 0.418 0.386 0.381 0.427 0.428	200			MSE	0.252	0.260	0.256	0.277	-	0.278	0.280	0.278	0.272	0.287	0.283	0.309	0.339
MSE 0356 0370 0371 0365 - 0358 0365 0354 0398 0351 0345 0377 0403	033		720	MAE	0.356	0.322	0.323	0.350	-	0.349	0.359	0.348	0.418	0.386	0.381	0.42/	0.428
900	900			111012	0.550	0.520	0.321	0.505	-	0.550	0.505	0.554	0.570	0.001	0.545	0.577	0.400

Table A3: Performance metrics for various models. Key: Best results, Second-best results.

A.4.2 **OBSERVABILITY BENCHMARK**

Case	%
Sparse	12.20
Extreme Right Skew	17.07
Seasonal	80.49
Flat	1.22

> Table A4: Breakdown of observability dataset based on case, computed based on the average char-acteristics of variates in each multivariate series. Note that these do not add to 100% because time series may fall into multiple categories.