BRIDGING THE HIGH-FREQUENCY DATA GAP: A MILLISECOND-RESOLUTION DATASET FOR ADVANCING TIME SERIES FOUNDATION MODELS

Anonymous authorsPaper under double-blind review

ABSTRACT

Time series foundation models (TSFMs) require diverse, real-world datasets to adapt across varying domains and temporal frequencies. However, current largescale datasets predominantly focus on low-frequency time series with sampling intervals, i.e., time resolution, in the range of seconds to years, hindering their ability to capture the nuances of high-frequency time series data. To address this limitation, we introduce a novel dataset that captures millisecond-resolution wireless and traffic conditions from an operational 5G wireless deployment, expanding the scope of TSFMs to incorporate high-frequency data for pre-training. Further, the dataset introduces a new domain, wireless networks, thus complementing existing more general domains like energy and finance. The dataset also provides use cases for short-term forecasting, with prediction horizons spanning from 100 milliseconds (1 step) to 9.6 seconds (96 steps). By benchmarking traditional machine learning models and TSFMs on predictive tasks using this dataset, we demonstrate that TSFMs perform poorly on this new data distribution in both zero-shot and fine-tuned settings. Our work underscores the importance of incorporating highfrequency datasets during pre-training and forecasting to enhance architectures, fine-tuning strategies, generalization, and robustness of TSFMs in real-world applications.

1 Introduction

Foundation models (FMs) have significantly enhanced machine learning (ML) by utilizing large-scale pre-training on diverse datasets, enabling them to generalize across a wide array of tasks and domains (Thakur, 2024). Recently, time series foundation models (TSFMs) have attracted more interest due to their capability to handle complex temporal tasks, with a particular focus on generalizing across varying time scales and domains, including forecasting, anomaly detection, and classification (Liang et al., 2024). However, developing effective TSFMs requires access to datasets that capture diverse real-world scenarios at varying frequencies and across different domains. The blue dots in Fig. 1 demonstrate that the existing benchmark datasets predominantly focus on low-frequency time series with sampling intervals in the range of seconds to years.

Hence, the focus of this paper is to develop and benchmark a high-frequency dataset in the millisecond resolution by comparing the performance of TSFMs with shallow machine learning models to enable new architectures and fine-tuning strategies that can extend to high-frequency data use cases and potentially provide generalizable and diverse characteristics that can improve the accuracy of TSFMs on existing datasets as well.

The main contributions of this paper and dataset are: (1) Extending the scope of pre-training and generalizability for state-of-the-art TSFMs by providing a dataset at millisecond resolution (Fig. 1). (2) Introduction of a new domain, namely, wireless networks, to the existing domains of open datasets (Fig. 2). (3) Applications with short-term forecasting, with prediction horizons spanning from 100 milliseconds (1 step) to 9.6 seconds (96 steps) (Fig. 3).

The rest of the paper is organized as follows. Related work is discussed in Section 2. Section 3 provides a detailed description of the 5G network data, and its characteristics. Section 4 presents the details of models benchmarked, including experimental evaluation and analysis. Section 5 discusses

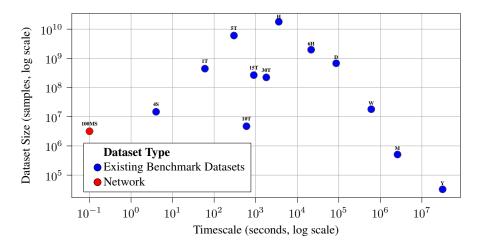


Figure 1: Comparison of timescales and dataset sizes for standard existing datasets used for pretraining (Table 14 in (Aksu et al., 2024)) as compared with the new benchmark. The red dot represents the new dataset that is introduced in this paper.

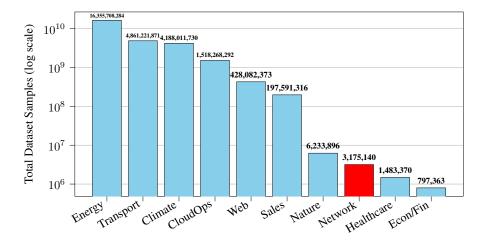


Figure 2: Comparison of existing domains for pre-training (Table 14 in (Aksu et al., 2024)) with the new benchmark. The red bar represents the new dataset that is introduced in this paper.

the limitations of this work. Finally, in Section 6, we conclude and provide directions for future research.

2 Related Work

Time Series Foundation Models (TSFMs) have surged in recent years, with their architectures continually evolving to achieve improved performance in both zero-shot and fine-tuned scenarios. Notably, several TSFMs have garnered widespread attention within the community, including Chronos (Ansari et al., 2024), TTM (Ekambaram et al., 2024), Moirai (Woo et al., 2024), TimesFM (Das et al., 2024), and Time-MOE (Xiaoming et al., 2025). These models can be broadly categorized into two distinct classes: transformer-based and non-transformer-based architectures (Liang et al., 2024). Our work complements these developments by introducing a high-frequency, real-world dataset from a novel domain (wireless networks), which provides an additional and challenging benchmark for evaluating the robustness and adaptability of TSFMs.

Transformer-based TSFMs largely follow established self-supervised (e.g., Moirai) or supervised transformer frameworks (e.g., TimeXer), which have garnered significant recognition within the field. In contrast, non-transformer-based TSFMs leverage alternative machine learning models such

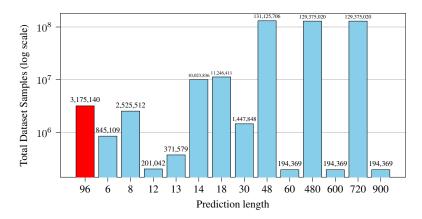


Figure 3: Comparison of prediction lengths of standard test data (Table 2 in (Aksu et al., 2024)) as compared with the new benchmark. The red bar represents the new dataset that is introduced in this paper.

as Multi-Layer Perceptron (MLP) and Convolutional Neural Networks (CNN) (e.g., TTMs). More recent efforts have also focused on enhancing diffusion-based methods (Kollovieh et al., 2023; Su et al., 2025) for modeling and generating data of different characteristics, which is crucial for generative time series forecasting. Furthermore, to address statistical heterogeneity in time series foundation model training and ensure robust generalization, a decentralized cross-domain model fusion approach, as Federated Learning (FL), has been explored in (Chen et al., 2025).

The successful deployment of these TSFMs for accurate zero-shot forecasting relies on the development of pre-trained models that have undergone extensive training on datasets characterized by diverse patterns and resolution properties. This emphasis on data diversity is critical, as it enables TSFMs to exhibit generalizability across a wide range of scenarios and capture complex temporal dynamics with enhanced accuracy. Notably, prior research has underscored the importance of resolution and domain diversity in pre-trained models for optimizing performance (e.g., Section 4 in (Ansari et al., 2024) for Chronos and Section 4.9 and Fig. 3 in (Ekambaram et al., 2024) for TTM.

In practice, a range of open datasets is available for TSFMs, which collectively provide the necessary heterogeneity to ensure that these models generalize effectively to out-of-domain datasets and real-world applications. Specifically, popular datasets such as those from Monash (Godahewa et al., 2021), LIBCITY (Wang et al., 2021), and the UCI Machine Learning archive (Asuncion et al., 2007) have become foundational in pre-training TSFMs and are widely utilized for assessing model performance. These datasets not only serve as data for pre-trained models but also enable out-of-domain testing of pre-trained models when a subset of the datasets are not considered for pre-training. We position our dataset as a complementary resource to these existing open datasets, specifically targeting the gap for millisecond-level time series from communication networks for both training and out-of-domain evaluation of TSFMs. Our dataset directly addresses this need for diversity by introducing a previously underrepresented domain with very fine temporal granularity, thereby contributing to a better understanding of the generalization capabilities of TSFMs when applied to high-frequency wireless data.

This paper provides a benchmark dataset that can fill the critical gap for high-frequency data for TSFMs. In contrast to other high-frequency datasets, our network dataset provides carefully curated use cases for univariate and multivariate forecasting problems ideally suited for TSFMs, along with an initial benchmark study on this dataset.

3 Dataset

3.1 Dataset Overview

We utilize a time series dataset of 5G Radio Access Network (RAN) Performance Measurements (PMs) collected from a real-world deployment of a 5G Open Radio Access Network (O-RAN)

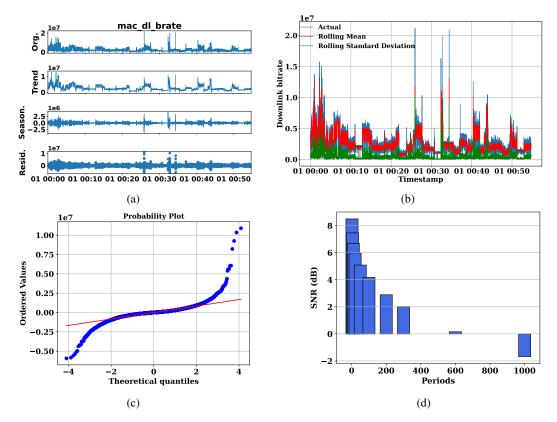


Figure 4: Target variable (Downlink Bitrate; mac_dl_brate): (a) STL decomposition, (b) Rolling mean and standard deviation, (c) Residual Q-Q, (d) Signal-to-Noise Ratio (dB).

within the OpenIreland testbed. O-RAN introduces a modular and open architecture that decomposes the traditional monolithic RAN into standardized, interoperable components (i.e, the Central Unit (CU), Distributed Unit (DU), and Radio Unit (RU)) facilitating multi-vendor deployments and software-driven control. Central to O-RAN's programmability is the near-Real-Time RAN Intelligent Controller (near-RT RIC), which enables rapid, feedback-driven network optimization.

The data was captured using software-defined radios (Ettus USRPs) configured as a base station and multiple user equipments (UEs). To simulate diverse real-world usage, the setup incorporated various mobility profiles (static, pedestrian, car, bus, and train) and generated traffic from both benign applications (web browsing, VoIP, IoT, and video streaming) and malicious activities (DDoS-Ripper, DoS-Hulk, PortScan, Slowloris). PMs were collected at the base station side and span a broad set of physical and medium access control layer features, including the Channel Quality Indicator (CQI), Modulation and Coding Scheme (MCS), Noise ratio interference (SINR), Signal strength (RSSI), buffer occupancy, and packet delivery statistics. In the dataset, each UE is associated with a unique identifier, denoted as *ue_ident*, which serves to distinguish individual UEs across all collected traces. This identifier remains consistent for a given UE, regardless of the mobility pattern or traffic class associated with its traces. The resulting dataset enables temporal modeling of RAN dynamics under realistic operational conditions.

This data context is particularly well-suited for very short-term forecasting, where the goal is to predict network states (e.g., throughput, channel quality, traffic class) over a short horizon ranging from milliseconds to a few seconds. Such forecasting enables predictive control strategies in scenarios characterized by rapid fluctuations in load, mobility, or interference (see Section 3.2 for dataset characteristics). Short-term throughput predictions enhance scheduling efficiency and application-level rate control, especially in latency-sensitive services like cloud gaming or interactive video. Forecasting CQI, for example, allows the network to proactively steer users to cells with better anticipated radio conditions, support load-aware handovers, and preemptively adjust adaptive bitrate algorithms for video streaming. Likewise, anticipating traffic class transitions supports early enforcement of

216217

Table 1: Summary of STL Decomposition of all datasets.

230231232

232233234

235236237

244

245246247248249250251252253

254

255

256

257

263

264

265

266

267

268

269

STL Decomposition Dataset Trend Seasonality Residuals Sharp spikes. Weak short-term Bursts of noise. Network Unstable, step-like shifts. periodic patterns. Lots of unpredictable Hidden by noise. variation. Mostly steady with small Small, regular ETTh1 Tiny random changes. rises and falls. repeating pattern. Electricity Remains steady throughout. Strong repeating pattern. Occasional bursts of noise. Almost flat but interrupted Mostly small, but with rare Weather No seasonality. by sudden sharp spikes. sudden jumps. Slowly increasing Strong, regular **Traffic** Small random changes. trend over time. repeating pattern.

QoS policies, dynamic resource allocation (e.g., in network slicing), and intrusion detection mechanisms capable of identifying malicious activity before it significantly degrades the service.

3.2 Dataset Characteristics

While Section 3.1 provides a broad overview of the 5G network dataset, our analysis and experiments are carried out on a carefully filtered subset of the data. We filter the raw data on the basis of the mobility pattern and benign traffic class. In particular, the *static* mobility pattern for the *video streaming* traffic class. Therefore, the results presented here represent the characteristics of the filtered dataset rather than those of the complete dataset.

The time series of the 5G network demonstrates several important characteristics. Fig. 4a shows the STL (Seasonal and Trend decomposition using Loess) of the time series, which separates the original data (labeled Org. in Fig. 4a) into distinct structural components, i.e., the trend, seasonal and residual components. Here, the trend component reflects the underlying structure of the series; however, it appears unstable, as characterized by step-like shifts rather than a smooth trajectory. The seasonal component captures only weak short-term periodic patterns, which are easily obfuscated by the stronger irregular behavior in the data. The residual component contains the remaining variability, including sharp spikes and bursts of endogenous noise that cannot be explained by trend or seasonality. Similarly, as illustrated in Fig. 4b, both the rolling mean and the standard deviation are observed to change substantially over time, confirming that the process is non-stationary and heteroskedastic. This means that the statistical properties of the data are not constant. The data exhibit extreme outlier events that are more prominent in specific time periods than in random events throughout the series. The autocorrelation analysis (see Section 8) reveals a strong temporal persistence with slow decay, confirming the clustering of extreme events observed in the data. In Fig. 4c, the residuals deviate strongly from the reference line, particularly in the tails, indicating a heavy-tailed distribution. Finally, the signal-to-noise ratio (SNR) analysis in Fig. 4d provides a quantitative view of this instability. The SNR values highlight that the series is dominated by shortterm periodic structures (high SNR in periods 2-20), while medium-term cycles exist but are weaker, and long-term seasonality is essentially absent (SNR nears to zero and even negative beyond period 600). Overall, the time series is mostly influenced by short-term changes, bursts of volatility and clustered anomalies, rather than stable long-term trends.

Next, we provide a summary of the overview on the comparison between our 5G network dataset and other common pre-trained datasets (further experimental details are presented in Appendix A.2). The pre-trained datasets used for comparison are: ETTh1 (Zhou et al., 2021) is an hourly subset of the Electricity Transformer Temperature (ETT) dataset, containing two years of transformer oil temperature and related power load data from two counties in China. Electricity (Wu et al., 2021) dataset contains the hourly electricity consumption(in kWh) from 321 clients, recorded between 2012 and 2014. Weather (Wu et al., 2021) data from 2020 in Germany, recorded every 10 minutes, with 21 indicators such as air temperature, humidity, and wind speed. Traffic (Wu et al., 2021) is a collection of hourly road occupancy rates (0–1) from sensors on San Francisco Bay Area freeways,

collected by the California Department of Transportation between 2015 and 2016. Table 1 summarizes the key differences among the datasets based on their STL decomposition, highlighting that our dataset is notably different due to its unstable trend, weak seasonality, and spiky residuals. Appendix A.2 includes other data characteristics, such as temporal dependencies, and statistical variability.

4 Benchmark

In this section, we provide a comprehensive analysis of the benchmarked models (as explained in 4.1) for the considered target variable *downlink bitrate* (*bitrate*) in the 5G network dataset. In the multivariate setting, all considered models use four input features, with descriptions provided in Table 2. Section 4.3 provides implementation details, including the data processing pipeline, that reflects our consideration of only a subset of data to illustrate the impact of this high frequency dataset.

Table 2: Features used in multivariate setting.

Feature	Description			
CQI	Channel Quality Indicator			
MCS	Modulation and Coding Scheme			
pkt ok/nok	Number of packets sent/dropped			

4.1 Models benchmarked

We selected three state-of-the-art tree-based ensemble models: Random Forest (RF) (Breiman, 2001), implemented using Scikit-learn, eXtreme Gradient Boosting (XGBoost, hereafter XGB) (Chen & Guestrin, 2016), and Adaptive Random Forest (ARF) (Gomes et al., 2018) implemented using the River library. Similarly, we selected a non-parametric baseline, referred to as naive forecast (Naive) (Beck et al., 2025), for a fair evaluation on high-frequency data.

In addition, we evaluated three time series foundation models (TSFMs): TinyTimeMixer (TTM) (Ekambaram et al., 2024), Chronos (Ansari et al., 2024), and Lag-Llama (Rasul et al., 2023), each specifically designed for time series forecasting. TTM is an extremely light-weight pre-trained model, with effective transfer learning capabilities based on the light-weight TSMixer architecture. Likewise, Chronos is a language modeling framework for time series for pre-trained probabilistic time series models. In this work, we specifically adopted the Chronos-bolt-small variant (46M parameters) as the representative Chronos model for our experiments. Lag-Llama is a general-purpose foundation model for univariate probabilistic time series forecasting based on a decoder-only transformer architecture that uses lags as covariates.

Table 3: Parameters used in model training.

Parameter	Univariate	Multivariate	
n_models	10	20	
max_features	None	0.5	
grace_period	50	100	
max_depth	None	5	

⁽a) Hyper-parameters specific to ARF.

Parameter	Value		
Target variable	Downlink bitrate		
No. of features	4		
Mobility patterns	Static		
Past observations	5		
Prediction horizon	96		
Train set:Test set	80:20		

(b) Common parameters for all shallow models.

4.2 System specification

The experiments are carried out on a local machine with the following hardware and software specifications: **Operating System:** Microsoft Windows 10 Enterprise, Version 22H2; **Processor:** 11th

Gen Intel(R) Core(TM) i7-1165G7 CPU @ 2.80 GHz with 4 cores and 8 threads; **Memory:** 32 GB RAM.

4.3 IMPLEMENTATION DETAILS

Pre-processing: During data pre-processing, we changed the time resolution of our dataset by converting the original millisecond-level observations into 100-millisecond intervals. This choice reflects practical constraints in O-RAN networks, where collecting performance measurements at every millisecond would impose excessive overhead. For shallow models, input sequences are constructed using a sliding-window approach, where past observations within a fixed window are used to predict future target values. For TSFMs, we follow the original implementation protocols described in their respective papers. The prediction horizons range from 1 millisecond up to 9.6 seconds. Short-term horizons are often straightforward, as the target variable (i.e., bitrate) tends to remain stationary across very small timescales. In contrast, longer horizons provide more meaningful insights, enabling applications such as video streaming to anticipate changes in bitrate and proactively adjust parameters like encoding level. These long horizon forecasts are valuable both for adapting Quality of Service (QoS) and for estimating the stability of the bitrate, that is, how frequently it is expected to change.

Model parameters: Table 3 summarizes the parameters used during model training. For common parameters shared across all models, offline experiments were conducted to select optimal values based on prediction accuracy, ensuring fair benchmarking conditions. Both RF and XGB models used these optimized common parameters along with their respective default model-specific hyper-parameters without additional tuning. For the ARF model, while using the same optimized common parameters, model-specific hyper-parameter tuning was performed using random search methodology, with parameter ranges detailed in Table 3a. The best performing Root Mean Square Error (RMSE)-based ARF configuration was selected for the final evaluation. Furthermore, the prediction horizon for all models is set at 96 steps, with each step representing a 100-millisecond interval; which corresponds to predicting the next 9.6 seconds (9600 ms).

Model training: For RF and XGB, we utilized Scikit-learn's *MultiOutputRegressor* wrapper to enable direct multi-step forecasting. For TSFMs, we follow the original implementation protocols described in their respective papers.

Post-processing: Both Chronos and Lag-Llama are trained to predict a fixed length horizon **H** from a given context window. By default, these models produce forecasts only for the final prediction window of each series and skip series that do not meet the minimum context length. This default evaluation framework differs from shallow models that generate forecasts for every test sample. To ensure a consistent comparison across models, we implemented a rolling evaluation procedure for both Chronos and Lag-Llama. Specifically, starting with each timestamp **t**, we provide the model with all historical data available up to **t** and generate the next steps **H**. We then slide the starting point forward by one time step and repeat the prediction until the end of the series. This produces overlapping multi-step forecasts aligned with each test timestamp, allowing direct comparison with the shallow models.

Table 4: Performance metrics of benchmarked models.

	Univ	ariate	Multivariate		
Model	RMSE	MAE	RMSE	MAE	
RF	0.0344 ± 0.0000	0.0227 ± 0.0000	0.0342 ± 0.0000	0.0226 ± 0.0000	
XGB	0.0354 ± 0.0000	0.0232 ± 0.0000	0.0356 ± 0.0000	0.0232 ± 0.0000	
ARF	$\textbf{0.0270} \pm \textbf{0.0002}$	$\textbf{0.0189} \pm \textbf{0.0001}$	$\textbf{0.0175} \pm \textbf{0.0007}$	0.0130 ± 0.0005	
Naive	0.0418 ± 0.0000	0.0240 ± 0.0000	0.0418 ± 0.0000	0.0240 ± 0.0000	
TTM (Zero-shot)	0.0359 ± 0.0000	0.0229 ± 0.0000	0.0359 ± 0.0000	0.0230 ± 0.0000	
TTM (Fine-tuning)	0.0360 ± 0.0000	0.0228 ± 0.0000	0.0391 ± 0.0000	0.0249 ± 0.0000	
Chronos (Zero-shot)	0.0313 ± 0.0000	0.0185 ± 0.0000	-	-	
Chronos (Fine-tuning)	0.0313 ± 0.0000	0.0185 ± 0.0000	-	-	
Lag-Llama (Zero-shot)	0.0617 ± 0.0002	0.0384 ± 0.0001	-	-	
Lag-Llama (Fine-tuning)	0.0474 ± 0.0039	0.0268 ± 0.0009	-	-	

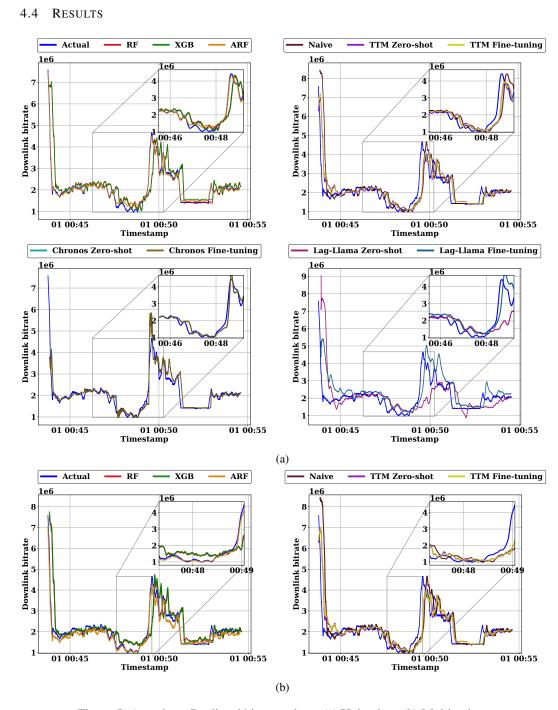


Figure 5: Actual v.s. Predicted bitrate values: (a) Univariate, (b) Multivariate.

In this section, we evaluate the performance of shallow models and TSFMs in both univariate and multivariate settings. Table 4 presents the performance of the benchmarked shallow models and TSFMs, evaluated using RMSE and Mean Absolute Error (MAE). In both settings, ARF consistently outperforms the other shallow models and TSFMs. The performance gain is consistent with the data characteristics observed in Section 3.2; our 5G network dataset is dominated by irregular spikes, step-like changes, and lack of stable seasonality. Static models such as RF or XGB struggle in performance because they assume that the training distribution does not change over time, leading to poor generalization when sudden data shifts occur. Similarly, TSFMs performance degrades due to a shift in data distribution in the zero-shot scenario, as these pre-trained models are trained only

on low-frequency data, limiting their ability to capture high-frequency dynamics with unpredictable spikes and irregular patterns. Even after fine-tuning on our dataset, the performance of TSFMs remains suboptimal, as they fail to generalize effectively. In contrast, ARF is designed to handle concept drift by dynamically updating its ensemble of trees as new patterns appear. This allows it to quickly adapt to data distribution changes and maintain predictive accuracy even in the presence of strong irregularities.

The performance of these models is more clearly reflected in Fig. 5. We observe that ARF follow the curve/trend of the **bitrate** much better than the other shallow models and TSFMs. For the purpose of visualization, we average the actual and predicted values for each test sample. While it is observed that Chronos offers a competitive performance in the univariate setting, it does not operate directly on the multivariate time series, which is the considered scenario for this work. The benefits of including the exogenous features are therefore missed out, limiting the practical applicability of using as it is.

5 LIMITATIONS

Our current study provides valuable insights into the performance of shallow models and TSFMs for millisecond resolution wireless network data, and shows the need to utilize this dataset to enhance the generalizability and applicability of TSFM pre-training and fine-tuning capabilities. However, there are certain limitations in the study that highlight areas for potential improvement in future research. These include:

- The empirical benchmark results for shallow models such as XGBoost and Random Forest
 only had limited hyper-parameter tuning, whereas standard Hyperparameter Optimization
 (HPO) techniques could have been applied to further optimize their performance. Given
 the paper's primary focus on comparing benchmark performance between shallow models
 and TSFMs, any potential marginal improvements through HPO were deemed secondary
 to the main objective.
- Further, default implementations of the TSFMs were considered for the performance on zero-shot models. Feature engineering and data preprocessing strategies can potentially improve the performance of TSFMs but this was not considered. Since shallow models work directly on the raw data and perform reliable forecasting, the same was done for TSFMs to make the comparison fair.
- Default fine-tuning implementations were explored for each TSFM, but novel techniques such as autotuning and Low-Rank Adaptation (LoRA) (Hu et al., 2022) strategies were not considered since the focus was on zero-shot and few-shot learning. Future work on ablation studies is proposed to investigate whether optimizing few-shot learning parameters can significantly enhance the performance of TSFMs.

6 CONCLUSION AND FUTURE WORK

We present a novel high-frequency time series dataset capturing millisecond-resolution measurements from real-world wireless network. This dataset fills a critical gap in existing large-scale resources, which largely lack fine-grained, real-time wireless network data. Our experiments reveal the limitations of current TSFMs and highlight the need to incorporate diverse, high-resolution datasets during pre-training to improve generalization. In the future, we will use this dataset for the use case of anomaly detection and transfer learning across various mobility profiles.

REFERENCES

Taha Aksu, Gerald Woo, Juncheng Liu, Xu Liu, Chenghao Liu, Silvio Savarese, Caiming Xiong, and Doyen Sahoo. Gift-eval: A benchmark for general time series forecasting model evaluation. *arXiv preprint arXiv:2410.10393*, 2024.

Abdul Fatir Ansari, Lorenzo Stella, Caner Turkmen, Xiyuan Zhang, Pedro Mercado, Huibin Shen, Oleksandr Shchur, Syama Syndar Rangapuram, Sebastian Pineda Arango, Shubham Kapoor,

Jasper Zschiegner, Danielle C. Maddix, Michael W. Mahoney, Kari Torkkola, Andrew Gordon Wilson, Michael Bohlke-Schneider, and Yuyang Wang. Chronos: Learning the language of time series. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=gerNCVqqtR.

- Arthur Asuncion, David Newman, et al. Uci machine learning repository. 2007. Published by Irvine, CA, USA.
- Nico Beck, Jonas Dovern, and Stefanie Vogl. Mind the naive forecast! a rigorous evaluation of forecasting models for time series with low predictability: N. beck et al. *Applied Intelligence*, 55 (6):395, 2025.
- Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001. doi: 10.1023/A: 1010933404324. Published by Springer.
- Shengchao Chen, Guodong Long, Jing Jiang, and Chengqi Zhang. Federated foundation models on heterogeneous time series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 15839–15847, 2025.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2016.
- 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.
- Vijay Ekambaram, Arindam Jati, Pankaj Dayama, Sumanta Mukherjee, Nam Nguyen, Wesley M Gifford, Chandra Reddy, and Jayant Kalagnanam. Tiny time mixers (ttms): Fast pre-trained models for enhanced zero/few-shot forecasting of multivariate time series. *Advances in Neural Information Processing Systems*, 37:74147–74181, 2024.
- Rakshitha Godahewa, Christoph Bergmeir, Geoffrey I. Webb, Rob J. Hyndman, and Pablo Montero-Manso. Monash time series forecasting archive. In *Neural Information Processing Systems Track* on *Datasets and Benchmarks*, 2021.
- Heitor Murilo Gomes, Jean Paul Barddal, Luis Eduardo Boiko Ferreira, and Albert Bifet. Adaptive random forests for data stream regression. In *ESANN*, 2018.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *International Conference on Learning Representations, ICLR*, 1(2):3, 2022.
- Marcel Kollovieh, Abdul Fatir Ansari, Michael Bohlke-Schneider, Jasper Zschiegner, Hao Wang, and Yuyang Bernie Wang. Predict, refine, synthesize: Self-guiding diffusion models for probabilistic time series forecasting. *Advances in Neural Information Processing Systems*, 36:28341–28364, 2023.
- Yuxuan Liang, Haomin Wen, Yuqi Nie, Yushan Jiang, Ming Jin, Dongjin Song, Shirui Pan, and Qingsong Wen. Foundation models for time series analysis: A tutorial and survey. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 6555–6565, 2024.
- Kashif Rasul, Arjun Ashok, Andrew Robert Williams, Arian Khorasani, George Adamopoulos, Rishika Bhagwatkar, Marin Biloš, Hena Ghonia, Nadhir Hassen, Anderson Schneider, et al. Laglama: Towards foundation models for time series forecasting. In *R0-FoMo: Robustness of Fewshot and Zero-shot Learning in Large Foundation Models*, 2023.
- Chen Su, Zhengzhou Cai, Yuanhe Tian, Zhuochao Chang, Zihong Zheng, and Yan Song. Diffusion models for time series forecasting: A survey. *arXiv preprint arXiv:2507.14507*, 2025.
- Suresh Chandra Thakur. Foundation models for time series forecasting. *International IT Journal of Research, ISSN: 3007-6706*, 2(4):144–156, 2024.

Jingyuan Wang, Jiawei Jiang, Wenjun Jiang, Chao Li, and Wayne Xin Zhao. Libcity: An open library for traffic prediction. In *Proceedings of the 29th International Conference on Advances in Geographic Information Systems*, SIGSPATIAL '21, pp. 145–148, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450386647.

Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sahoo. Unified training of universal time series forecasting transformers. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.

Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. Advances in Neural Information Processing Systems, NeurIPS, 34:22419–22430, 2021.

Shi Xiaoming, Wang Shiyu, Nie Yuqi, Li Dianqi, Ye Zhou, Wen Qingsong, and Ming Jin. Timemoe: Billion-scale time series foundation models with mixture of experts. In *ICLR 2025: The Thirteenth International Conference on Learning Representations*. International Conference on Learning Representations, 2025.

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

A APPENDIX

A.1 Performance Evaluation Metrics

The Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and Mean Absolute Percentage Error (MAPE) are calculated as follows:

$$RMSE(Y_t, \hat{Y}_t) = \sqrt{\frac{1}{T} \sum_{t=1}^{T} (Y_t - \hat{Y}_t)^2},$$
 (1)

$$MAE(Y_t, \hat{Y}_t) = \frac{1}{T} \sum_{t=1}^{T} |(Y_t - \hat{Y}_t)|,$$
 (2)

where Y_t and \hat{Y}_t are the actual and predicted bitrate values, and T is the total number of samples in the test dataset.

A.2 DATA CHARACTERISTICS COMPARISON

In this section, we compare our 5G network dataset with those used in the pre-training of TSFMs. The comparison focuses on key data characteristics, including statistical distributions, temporal dependencies, and statistical variability, as illustrated in Figs. 6, 7, 8, and 9. We compare the datasets using STL decomposition, rolling mean and standard deviation, autocorrelation (ACF), and residual QQ plots. Our 5G network data is clearly the most different; its trend shifts abruptly in steps, seasonality is weak and mostly hidden by noise, rolling statistics change suddenly, the ACF shows strong temporal persistence with slow decay, and the residual QQ plot departs strongly from normality due to sharp spikes. In contrast, the ETTh1 dataset has a mostly steady trend with mild rises and falls, small but regular seasonal cycles, stable rolling statistics, weak cyclical autocorrelation, and residuals close to normal. The Electricity dataset also remains steady in its trend but shows stronger repeating seasonal patterns, its rolling mean is flat and variance is stable, clear cycles in the ACF, and residuals with occasional deviations. The Weather dataset is mostly flat with rare sharp jumps, no meaningful seasonality, sudden variance spikes in rolling statistics, weak ACF signals, and QQ plots highlighting outliers. Finally, Traffic dataset combines a smooth upward trend with strong, consistent seasonality, gradually increasing rolling mean with stable variance, clear seasonal autocorrelation, and residuals that follow normality fairly well. To conclude, our dataset differs from the others because its persistence comes from clustered extremes and abrupt shifts rather than smooth or cyclical structure, making it the least regular and most unpredictable series.

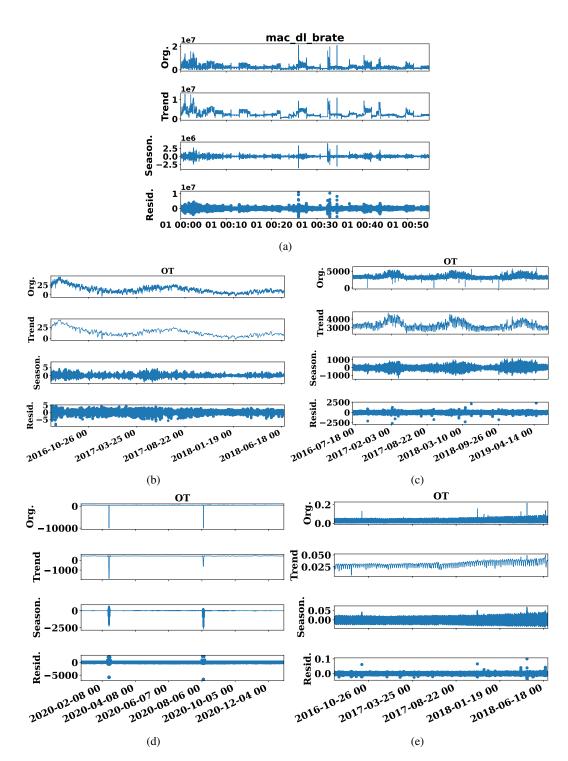


Figure 6: STL decomposition of time series: (a) Network, (b) ETTh1, (c) Electricity, (d) Weather, (e) Traffic.

A.3 ABLATION STUDY

In this section, we evaluate the performance of the benchmarked models on mobility patterns and traffic classes that differ from those presented in Section 3.2. The raw data is filtered based on mobility patterns and traffic generated from malicious activities. In particular, we focus on the *train*

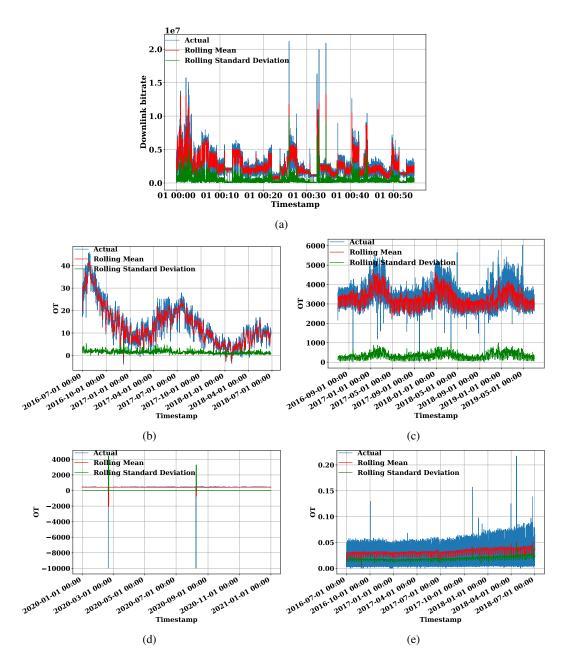


Figure 7: Rolling mean and standard deviation of time series: (a) Network, (b) ETTh1, (c) Electricity, (d) Weather, (e) Traffic.

mobility pattern for the *Dos-Hulk-C* traffic class. This analysis also demonstrates the potential of the dataset for transfer learning use case; by training models on one set of mobility patterns and traffic classes and evaluating them on a different set, we can assess how well knowledge learned in one context generalizes to another.

Table 5 presents the performance of selected benchmark shallow models and TSFMs, evaluated using RMSE and MAE in both univariate and multivariate settings. We specifically include TTM in our analysis because it supports both univariate and multivariate prediction tasks. For the filtered dataset, we observe that TTM outperforms ARF in the univariate setting. However, in the multivariate setting, ARF achieves better performance compared to the other models. Fig. 10 further illustrates how these models follow the trend of the bitrate.

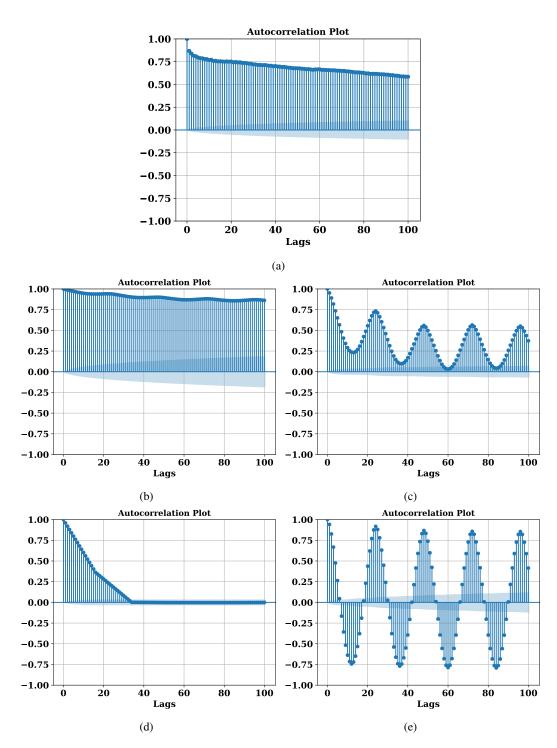


Figure 8: Autocorrelation of time series: (a) Network, (b) ETTh1, (c) Electricity, (d) Weather, (e) Traffic.

A.4 USE OF LARGE LANGUAGE MODELS

Large language models (LLMs) have been used exclusively for the purpose of text editing.

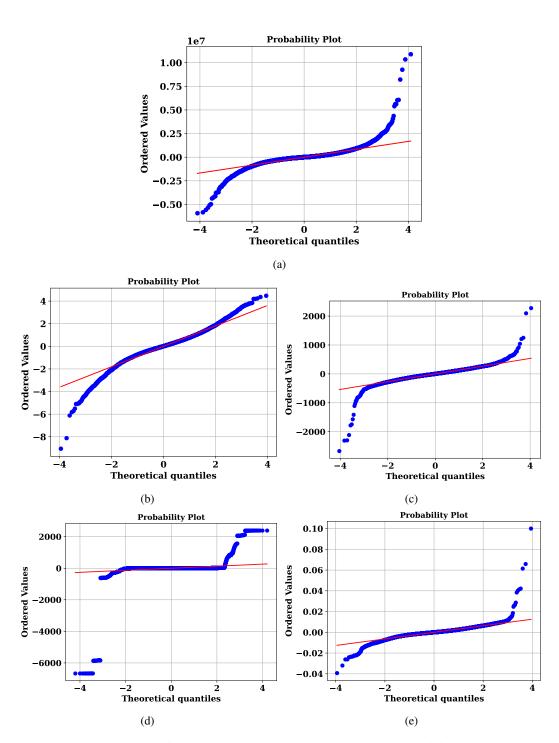


Figure 9: Autocorrelation of time series: (a) Network, (b) ETTh1, (c) Electricity, (d) Weather, (e) Traffic.

Table 5: Performance metrics of benchmarked models.

	Univ	ariate	Multivariate		
Model	RMSE	MAE	RMSE	MAE	
XGB ARF Naive TTM (Zero-shot)	0.1440 0.1728 0.1309 0.1279	0.1087 0.1125 0.0932 0.0922	0.1440 0.0968 0.1309 0.1279	0.1087 0.0634 0.0932 0.0922	

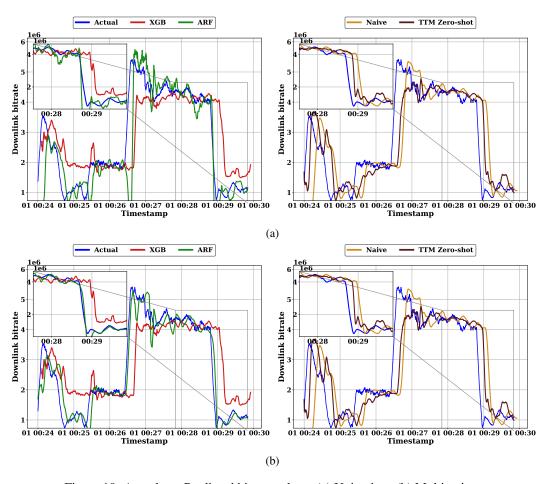


Figure 10: Actual v.s. Predicted bitrate values: (a) Univariate, (b) Multivariate.