Cross-Frequency Time Series Meta-Forecasting

Mike Van Ness* Stanford University Huibin Shen † AWS AI Labs

Hao Wang AWS AI Labs

Xiaoyong Jin AWS AI Labs

Danielle C. Maddix AWS AI Labs **Karthick Gopalswamy** AWS AI Labs

Abstract

Meta-forecasting is a newly emerging field which combines meta-learning and time series forecasting. The goal of meta-forecasting is to train over a collection of source time series and generalize to new time series one-at-a-time. Previous approaches in meta-forecasting achieve competitive performance, but with the restriction of training a separate model for each sampling frequency. In this work, we investigate meta-forecasting over different sampling frequencies, and introduce a new model, the Continuous Frequency Adapter (CFA), specifically designed to learn frequency-invariant representations. We find that CFA greatly improves performance when generalizing to unseen frequencies, providing a first step towards forecasting over larger multi-frequency datasets.

1 Introduction

Time series forecasting is a classical statistical problem with practical applications in several fields, such as finance, business management [16]. Local statistical models such as ARIMA and ETS [9] have long been state-of-the-art for forecasting. In recent years, much effort has been put into matching the performance of local models with deep learning approaches, particularly when modeling several closely-related time series [17, 14, 3].

When less data is available in a target dataset, transfer learning from source to target data is often necessary to compete with local methods. One approach is meta-learning, or domain generalization, where a model is trained to generalize to new target domains after an initial training phase on source data. Recent work has shown that meta-learning for time series forecasting, or meta-forecasting, can achieve competitive performance with only local fine-tuning [7] or even with no fine-tuning [15]. Such approaches are *zero-shot* forecasters, as they can forecast out-of-domain time series one-at-a-time without access to any related time series.

Almost all of the previous transfer learning works use the assumption that source and target data come from the same sampling frequency, e.g. hourly, daily, monthly, etc. We propose a different assumption: all data is seasonal, but not necessarily from the same sampling frequency. The typical seasonality associated with each sampling frequency then creates a correspondence between sampling frequency and signal frequency. As seen in Figure 1, seasonal time series of different signal frequencies appear very similar to humans, but are challenging for typical forecasting models to transfer between. Along with our data assumption, we consider a new task, frequency generalization, in which we task a model to generalize to unseen frequencies during meta-test time. For successful frequency generalization, we propose a new model, the Continuous Frequency Adapter (CFA). As shown in Figure 1, CFA can forecast the correct frequency on data of new unseen frequencies, which other methods cannot. CFA

^{*}Work done while doing an internship at AWS AI Labs.

[†]Correspondence to huibishe@amazon.com

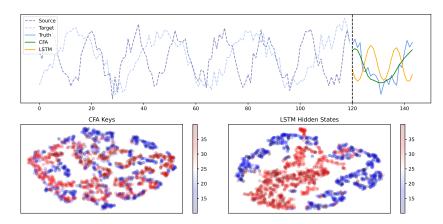


Figure 1: For frequency generalization, CFA outperforms LSTM. Both models are trained on synthetic source data with period length randomly sampled between [10, 20], and are zero-shot applied to synthetic target data with period length randomly sampled from [30, 40]. In the top plot, the LSTM model predicts the wrong seasonality, whereas CFA successfully adapts to the new seasonality. In the bottom plot, we see that CFA keys are invariant to frequency (color), but LSTM hidden states have different distributions for source and target frequencies ranges.

uses continuous domain adaptation [19] to enforce frequency-invariant hidden states, which is vital for frequency generalization. We summarize the novelty and contribution of this paper:

- We explore meta-learning over different sampling frequencies, which is previously unexplored. For this, we introduce a new task, *frequency generalization*, in which we task a model to generalize to time series of new sampling frequencies during test time.
- We develop a new model, CFA, which achieves improved performance for frequency generalization
 than previous meta-learning models. CFA uses continuous domain adaptation [19] to adapt to new
 sampling frequencies, a new technique for the time series literature. Specifically, CFA uses signal
 frequencies obtained from a Fourier transform of the input time series to define continuous domain
 indices, a novel technique for time series domain generalization.

Related Work Transfer learning has been explored previously for time series through several subfields. Time series representation learning [20, 21] does self-supervised pretraining for time series data. Domain adaptation approaches [11, 8] learn models that can adapt from one source dataset to a different target dataset, often utilizing adversarial training. In few-shot learning [10], a model transfers knowledge by learning how to learn from a small support set of related time series.

Recently, a few papers have particularly addressed meta-forecasting. In [15], the popular forecasting model N-BEATS is shown to fit a meta-learning framework, and achieves competitive zero-shot performance. In Meta-GLAR [7], a local closed-form head is used to adapt global representations to new time series. Our work is most similar to these meta-learning approaches in that performance on the target dataset is evaluated in a zero-shot manner. Our model, however, takes inspiration primarily from [11] in its use of adversarial training and self-attention. We emphasize that all of the above cited papers only consider transferring between datasets of the same frequency, with the exception of [11] which considers domain adaptation opposed to the harder task of zero-shot meta-learning.

2 Problem Definition

Time series forecasting is the problem of predicting future observations, i.e. forecasting, given a past context window. That is, for some time series $(z_t)_{t>0}$, data samples are of the form

$$oldsymbol{x} = oldsymbol{z}_{1: au_c}, \quad oldsymbol{y} = oldsymbol{z}_{ au_c+1: au_c+ au_f}$$

where τ_c is the length of the context and τ_f is the length of the forecast. A model f then estimates g from g. If f has parameters g, we aim to find the parameters g that minimize the forecasting loss, i.e.

$$\underset{a}{\operatorname{argmin}} \mathbb{E}[L_f(F(\boldsymbol{x}), \boldsymbol{y}; \theta)] \tag{1}$$

where L_f is a forecasting loss such as MSE.

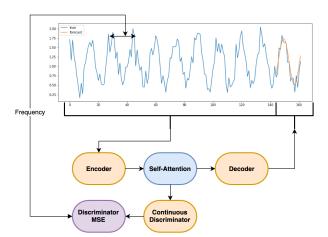


Figure 2: CFA Architecture. The encoder takes in a context window and produces keys, queries, and values, which are used by the self-attention module to produce representations for each timestep in the forecast window, which the decoder uses to make forecasts. Meanwhile, the keys and queries are passed to the continuous discriminator, which predicts the top k frequencies of the context window sorted by FFT amplitude. Adversarial training via Equation 2 is used to learn good forecasts while also making keys and queries frequency-invariant.

What distribution the expectation is taken over, i.e. what distribution z comes from, depends on the nature of the forecasting problem. In this paper, we consider the *zero-shot* framework, in which case $z \sim \mathcal{T}$ comes from some target distribution \mathcal{T} , but we only have training data $(x_1, y_1), \ldots, (x_n, y_n) \sim \mathcal{S}$ from a source distribution \mathcal{S} . This defines the problem of *meta-forecasting*, where the model f must meta-learn on \mathcal{S} with the goal of minimizing 1 on \mathcal{T} with no additional training on \mathcal{T} , where \mathcal{T} could be *any* target distribution unseen during training.

In this paper, we focus on *frequency generalization*, which we define as meta-forecasting with $\mathcal S$ and $\mathcal T$ representing distinct frequencies. This is more challenging than the setting considered in previous meta-forecasting papers, in which $\mathcal T$ only contains frequencies that are already seen in $\mathcal S$. We do explore this second easier setting to compare to previous papers, and present the results in Appendix A.2.

Since minimizing 1 on all potential \mathcal{T} is an ambitious and likely unrealistic desire, we restrict \mathcal{S} and \mathcal{T} to be distributions representing seasonal datasets. Under this assumption, different sampling frequencies correspond to different signal frequencies, and thus generalizing to new sampling frequencies corresponds to generalizing to new signal frequencies. This makes the frequency generalization problem realistic, since seasonal data of different signal frequencies appear quite similar to the human eye but are difficult for machine learning models to generalize between as shown in Figure 1. Relaxing this assumption would be significantly more challenging, and we leave this to future work.

3 Continuous Frequency Adaptation

Most meta-forecasting models struggle to learn frequency-invariant signal vital for meta-forecasting. This challenge is demonstrated in the bottom right portion of Figure 1, where an LSTM model learns hidden states whose distribution depends on the input signal frequency.

To overcome this challenge, we introduce the Continuous Frequency Adapter (CFA), which is specifically designed to learn frequency-invariant representations (see bottom left of Figure 1). CFA is primarily a self-attention network, and utilizes adversarial training to enforce frequency invariance in the attention keys and queries (but not the values). The model architecture is inspired by the Domain Adaptation Forecaster (DAF) [11], but uniquely uses continuous domain indices [19] obtained by a Fourier transform to generalize to unseen signal frequencies.

CFA Architecture The CFA architecture is summarized in Figure 2. The encoder, self-attention, and decoder blocks are similar to the DAF architecture [11]. The encoder consists of a position-wise MLP follows by a series of 1D convolutional layers, and the decoder is a position-wise MLP. The self-attention block is a standard multi-head attention block as in transformer architectures [18]. The CFA discriminator is an MLP, like in DAF, but unlike DAF, our discriminator outputs a continuous response to match the continuous domain index (see next paragraph). Also, unlike DAF, the encoder and decoder blocks are shared for all time series, since a continuous domain index does not allow for separate encoders and decoders for each possible index since there are infinitely many domain indices.

The discriminator takes as input all keys and queries from the self-attention block, in order to enforce domain invariance in the keys and queries via adversarial training (see Equation 2). The motivation for this choice is that the keys and queries are used to generate the attention weights, which tell the model, for any given time point, which other time points are most relevant. For time series data, especially seasonal time series data, this importance weighting can be independent of the signal frequency, as the attention weights only need to capture the phase of the signal. Meanwhile, the values in the self-attention block are left independent of the discriminator, and thus can learn information that is dependent on the given time series, e.g. what the time series typically looks like at each phase.

Continuous Domain Generalization A key component of CFA is the use of continuous domain indices [19]. These continuous domain indices serve as labels for the discriminator, which takes as input the self-attention keys and queries and outputs a continuous domain index prediction. The domain indices are obtained via an FFT on the inputted context window to capture the signal frequencies of the time series. Specifically, we take the absolute value of the FFT outputs to obtain the amplitude corresponding to each frequency bin. We then select the top k frequencies, sorted by their amplitudes, and use their inverses (i.e. the period lengths) as the discriminator labels. We normalize the labels to be between 0 and 1 to stabilize the discriminator loss. For synthetic data we use k=1, since the synthetic data is noist sine waves without subfrequencies, and on real data we use k=2.

Adversarial Loss CFA utilizes adversarial training via a typical minimax loss [6]. Let E be an encoder, F a forecasting decoder, and D a discriminator (for CFA, E generates the self-attention keys, queries, and values, and F produces forecasts using these self-attention inputs). The goal of CFA is to solve the following minimax problem:

$$\min_{E,F} \max_{D} \mathbb{E}[L_f(\boldsymbol{x}; E, F)] - \lambda \mathbb{E}[L_d(\boldsymbol{x}; E, D)]$$
(2)

where L_f is the forecasting loss and L_d is the discriminator loss. In words, D is trained to minimize L_d , thereby training a strong discriminator, while E and F are trained to both produce good forecasts (i.e. minimize L_f) while maintaining hidden states that the strong discriminator cannot predict well (i.e. maximizing L_d). Such adversarial training allows the model to learn frequency-invariant discriminator inputs (keys and queries) while still producing good forecasts. In practice, the forecast/generative parameters (E, F) and the discriminative parameters (D) are updated in an alternating fashion, see Algorithm 1.

Training Procedure Since CFA is designed to learn frequency-invariant keys and queries, it is essential that CFA is trained over source data with varied frequencies. The multi-source training procedure is illustrated in Algorithm 1. The training works by sampling one batch from each of the source datasets, and updating the generative and discriminative parameters from the sum of the batch losses. The forecast/generative parameters (E, F) and the discriminative parameters (D) are updated in an alternating fashion, as typical for adversarial training.

Algorithm 1 CFA Training Algorithm

```
1: Input: source datasets S_1, \ldots S_d, forecast loss L_f, discriminator loss L_d.
 2: for epoch = 1 to E do
        for i=1 to n_batches_per_epoch do
 3:
           Sample \boldsymbol{x}_i, \boldsymbol{y}_i \sim S_i for j = 1, \dots, d
 4:
           Compute generative loss L_j = L_f(\boldsymbol{x}_j, \boldsymbol{y}_j) - \lambda L_d(\boldsymbol{x}_j) for each j
 5:
           Update generative parameters via L = L_1 + \cdots + L_d
 6:
 7:
           Compute discriminative loss L_j = L_d(\boldsymbol{x}_j, \text{FFT}(\boldsymbol{x}_j)) for each j
 8:
           Update discriminative parameters via L = L_1 + \cdots + L_d
 9:
        end for
10: end for
```

4 Experiments

Models For our experiments, we consider the following models.

• Mean: simple baseline that forecasts the mean from the context window.

Source Range	Target Range	Mean	CFA	LSTM	NBEATS
(10, 15)	(15, 20)	0.260 ± 0.002	0.088 ± 0.011	0.312 ± 0.057	0.424 ± 0.017
	(20, 25)	0.259 ± 0.003	0.190 ± 0.025	0.456 ± 0.081	0.351 ± 0.007
	(25, 30)	0.262 ± 0.003	0.223 ± 0.041	0.426 ± 0.028	0.320 ± 0.004
(15, 20)	(10, 15)	0.259 ± 0.003	0.071 ± 0.014	0.388 ± 0.035	0.417 ± 0.011
	(20, 25)	0.259 ± 0.003	0.055 ± 0.004	0.185 ± 0.055	0.469 ± 0.011
	(25, 30)	0.262 ± 0.003	0.082 ± 0.007	0.499 ± 0.072	0.526 ± 0.012
(20, 25)	(10, 15)	0.259 ± 0.003	0.209 ± 0.099	0.497 ± 0.037	0.333 ± 0.005
	(15, 20)	0.260 ± 0.002	0.066 ± 0.013	0.282 ± 0.082	0.425 ± 0.008
	(25, 30)	0.262 ± 0.003	0.064 ± 0.006	0.134 ± 0.052	0.387 ± 0.020
(25, 30)	(10, 15)	0.259 ± 0.003	0.257 ± 0.056	0.442 ± 0.015	0.303 ± 0.006
	(15, 20)	0.260 ± 0.002	0.213 ± 0.081	0.533 ± 0.045	0.417 ± 0.007
	(20, 25)	0.259 ± 0.003	0.069 ± 0.015	0.177 ± 0.027	0.397 ± 0.021

Table 1: Frequency generalization on synthetic data, measured as MSE of forecast. Source and target range indicate the range of uniformly random period lengths in the source and target data, respectively. For each source/target pair, the model is trained on source and applied zero-shot to target. Across all pairs of ranges, CFA has the best performance.

- LSTM: an auto-regressive LSTM model similar to DeepAR [17].
- NBEATS: deep model with mostly linear layers and doubly-residual connections [14].
- CFA: Our model described in Section 3.

Since frequency generalization is a new task, there are some baselines which cannot readily be adapted. For one, any method which requires separate modules for different sampling frequencies cannot be used, e.g. DAF [11], because it is impossible to train a new module for the target sampling frequency in the zero-shot regime. Additionally, we think that is it critical on real data to have different forecasting lengths for different sampling frequencies, and thus require models that can forecast an arbitrary number of time steps ahead during test time. CFA and LSTM can easily do this since they are autoregressive forecasters, i.e. they use the previous forecasts to make each successive forecast. NBEATS, on the other hand, requires a fixed forecast length that cannot be adjusted during test time, and thus we do not use it as a baseline for real data experiments. We still consider NBEATS as a baseline for synthetic data generated with uniform forecast length, though, since NBEATS has been shown to be a strong meta-forecaster [15].

Synthetic Data We generate synthetic time series data using sinusoidal curves with Gaussian noise and uniformly random period length (inverse of frequency), see Appendix A.1 for full details. We designate one period range for source data and one period range for target data. Models are trained on source data and applied zero-shot to target data, evaluated by mean squared error (MSE) in the forecast window. Results are shown in Table 1. Across all combinations of source and target period ranges, CFA is either the best model or within one standard deviation of the best model. As shown in Figure 1, CFA is able to learn good forecasts on the source data while maintaining frequency-invariant keys and queries, allowing CFA to generalize to new frequencies. In comparison, LSTM and NBEATS learn frequency-dependent signal, and thus fail to generalize to new frequencies, even failing to beat the simple mean baseline.

Real Data We focus on real world datasets that exhibit clear seasonality. We use the following datasets, each with a different sampling frequency: elec (hourly), uber (daily), tourism monthly (monthly), and tourism quarterly (quarterly). We load all datasets using GluonTS [1]. More information on each dataset, data preprocessing, and setup can be found in Appendix A.1. For each experiment, we designate one dataset as the target dataset and use the other 3 as the source datasets. We evaluate models by their Normalized Deviation (ND) [11] on the forecasting window. We do not run NBEATS because it requires equal context/forecast lengths across datasets, which we do not enforce for frequency generalization. The results are shown in Table 2. As was the case with synthetic data, CFA outperforms LSTM for frequency generalization.

5 Conclusion

Previous meta-forecasting papers have shown strong performance, but only when training one model per sampling frequency. In this paper, we instead consider frequency generalization, i.e. generalizing to unseen frequencies, for which it is not possible to train one model per sampling frequency. While previous meta-forecasting models are not successful, our CFA model provides much improved

	Mean	CFA	LSTM
elec tourism_monthly tourism_quarterly uber	$\begin{array}{c} 0.404 \pm 0.000 \\ 0.312 \pm 0.000 \\ 0.229 \pm 0.000 \\ 0.203 \pm 0.000 \end{array}$	$\begin{array}{c} 0.281 \pm 0.065 \\ 0.226 \pm 0.019 \\ 0.178 \pm 0.021 \\ 0.166 \pm 0.012 \end{array}$	$\begin{array}{c} 0.376 \pm 0.031 \\ 0.296 \pm 0.033 \\ \textbf{0.19} \pm \textbf{0.018} \\ 0.252 \pm 0.02 \end{array}$

Table 2: Frequency generalization on real data, evaluated by test Normaized Deviation. Each dataset has a different frequency, and each row corresponds to one target dataset, using all other datasets as source. As on synthetic data, CFA is better at frequency generalization than LSTM. We do not evaluate N-BEATS in this setting because each dataset has a different context and forecast length, making N-BEATS incompatible.

performance. This is an important first step towards building forecasting models robust to signal frequency, which could be trained over larger and less constrained datasets.

References

- [1] A. Alexandrov, K. Benidis, M. Bohlke-Schneider, V. Flunkert, J. Gasthaus, T. Januschowski, D. C. Maddix, S. S. Rangapuram, D. Salinas, J. Schulz, et al. Gluonts: Probabilistic and neural time series modeling in python. *J. Mach. Learn. Res.*, 21(116):1–6, 2020.
- [2] G. Athanasopoulos, R. J. Hyndman, H. Song, and D. C. Wu. The tourism forecasting competition. *International Journal of Forecasting*, 27(3):822–844, 2011.
- [3] C. Challu, K. G. Olivares, B. N. Oreshkin, F. Garza, M. Mergenthaler, and A. Dubrawski. N-hits: Neural hierarchical interpolation for time series forecasting. *arXiv preprint arXiv:2201.12886*, 2022.
- [4] D. Dua and C. Graff. UCI machine learning repository, 2017.
- [5] fivethirtyeight. Uber tlc foil response dataset: https://github.com/fivethirtyeight/uber-tlc-foil-response, 2015.
- [6] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [7] R. Grazzi, V. Flunkert, D. Salinas, T. Januschowski, M. Seeger, and C. Archambeau. Metaforecasting by combining global deeprepresentations with local adaptation. *arXiv* preprint arXiv:2111.03418, 2021.
- [8] H. Hu, M. Tang, and C. Bai. Datsing: Data augmented time series forecasting with adversarial domain adaptation. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 2061–2064, 2020.
- [9] R. J. Hyndman and G. Athanasopoulos. Forecasting: principles and practice. OTexts, 2018.
- [10] T. Iwata and A. Kumagai. Few-shot learning for time-series forecasting. *arXiv preprint* arXiv:2009.14379, 2020.
- [11] X. Jin, Y. Park, D. Maddix, H. Wang, and Y. Wang. Domain adaptation for time series forecasting via attention sharing. In *International Conference on Machine Learning*, pages 10280–10297. PMLR, 2022.
- [12] G. Lai, W.-C. Chang, Y. Yang, and H. Liu. Modeling long-and short-term temporal patterns with deep neural networks. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 95–104, 2018.
- [13] Neural Forecasting Competition. Time series forecasting competition for computational intelligence, 2008. http://www.neural-forecasting-competition.com/NN5/, last accessed on 2022-10-22.
- [14] B. N. Oreshkin, D. Carpov, N. Chapados, and Y. Bengio. N-beats: Neural basis expansion analysis for interpretable time series forecasting. *arXiv preprint arXiv:1905.10437*, 2019.

- [15] B. N. Oreshkin, D. Carpov, N. Chapados, and Y. Bengio. Meta-learning framework with applications to zero-shot time-series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9242–9250, 2021.
- [16] F. Petropoulos, D. Apiletti, V. Assimakopoulos, M. Z. Babai, D. K. Barrow, S. B. Taieb, C. Bergmeir, R. J. Bessa, J. Bijak, J. E. Boylan, et al. Forecasting: theory and practice. *International Journal of Forecasting*, 2022.
- [17] D. Salinas, V. Flunkert, J. Gasthaus, and T. Januschowski. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.
- [18] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [19] H. Wang, H. He, and D. Katabi. Continuously indexed domain adaptation. In *International Conference on Machine Learning*, pages 9898–9907. PMLR, 2020.
- [20] Z. Yue, Y. Wang, J. Duan, T. Yang, C. Huang, Y. Tong, and B. Xu. Ts2vec: Towards universal representation of time series. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 8980–8987, 2022.
- [21] G. Zerveas, S. Jayaraman, D. Patel, A. Bhamidipaty, and C. Eickhoff. A transformer-based framework for multivariate time series representation learning. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2114–2124, 2021.

A Appendix

A.1 Additional Dataset Details

Synthetic Data Synthetic data is generated using GluonTS [1]. Each time series is generated as the sum of a sine curve and Gaussian noise. The parameter are chosen according to

- Phase: $\sim \text{Unif}(0, 2\pi)$.
- Period: $\sim \text{Unif}(p_{\min}, p_{\max})$.
- Amplitude: $\sim \text{Unif}(0.5, 2)$.
- Noise: $\sim \text{Normal}(\mu = 0, \sigma = 0.2)$.

The period range (p_{\min}, p_{\max}) for source and target data is indicated in each row of Table 1. Each time series has 120 samples in the context window, and 24 samples in the forecast window. For each synthetic dataset, we generate 5000 time series with exactly one context length plus one forecast length's worth of data. We designate 4000 time series as the training set and the 1000 as the test set.

Real Data Table 3 shows the details of the datasets used in the real data experiments. All datasets are accessed via GluonTS [1], with the GluonTS name for each dataset provided. Following our assumptions, we focus on datasets that exhibit clear seasonality. The hourly, daily, monthly, and quarterly datasets have period lengths of 24, 7, 12, and 4 respectively.

Unlike on synthetic data, where training and test samples come from different time series, training and test sets come from different time windows across all time series. The last forecast length of each time series is always used as the test set. For the training set, we sample batches using an expected number sampler, such that we sample as many unique context/forecast windows as possible from the training window.

Lastly, on real datasets, it is essential to scale each time series, as is typical in deep learning forecasting approaches. We scale each sample by subtracting by the mean and dividing by the standard deviation from the context window. We then rescale our forecasts and compute the forecast loss with the true unscaled forecast labels. Note that we use Normalized Deviation as our forecasting loss on real data, which itself does normalization and thus should be fed unscaled values.

Dataset	GluonTS Name	Frequency	Number of Time Series	Original Source
elec	electricity_nips	hourly	370	UCI [4]
traffic	traffic_nips	hourly	963	UCI [4]
solar	solar_nips	hourly	137	[12]
uber	uber_tlc_daily	daily	262	fivethirtyeight [5]
NN5	nn5_daily_without_missing	daily	111	NN5 challenge [13]
tourism monthly	tourism_monthly	monthly	366	Tourism competition [2]
tourism quarterly	tourism_quarterly	quarterly	366	Tourism competition [2]

Table 3: Description of all real-world datasets used in experiments.

	Multi LSTM	LSTM	Multi NBEATS	NBEATS	CFA
NN5	0.165 ± 0.009	0.158 ± 0.006	0.22 ± 0.011	0.18 ± 0.005	0.193 ± 0.01
Solar	0.969 ± 0.115	0.891 ± 0.111	0.686 ± 0.021	0.602 ± 0.019	1.005 ± 0.037
Traffic	0.325 ± 0.014	0.298 ± 0.007	0.249 ± 0.003	0.238 ± 0.003	0.449 ± 0.03

Table 4: Comparison of zero-shot performance for Question 1 on real data. LSTM and NBEATS are trained on one source dataset of the same frequency as the listed target (elec \rightarrow solar, traffic, uber \rightarrow NN5), while Multi LSTM, Multi NBEATS, and CFA are all trained on 4 source datasets of different frequency (elec, uber, tourism monthly, tourism quarterly). We see that in general, the models trained over uni-frequency source data are superior.

A.2 Uni vs Multi Frequency Experiments

Previous meta-forecasting papers [7, 15] only consider meta-learning across time series of the same sampling frequency. For most of the paper, we focus on frequency generalization, i.e. generalization to unseen frequencies in the target data. Another valid question to ask is: does using source data of multiple frequencies improve target performance even when the target data frequency is already seen? Even if the performance is comparable, multi-frequency source data has the additional benefit of needing to only train one model instead of one model per frequency.

To investigate this question, we run the following experiment. We designate the following datasets as source datasets: Elec, Uber, Tourism Monthly, Tourism Quarterly, and the following datasets as target datasets: NN5, Solar, Traffic. For each target dataset, the multi-source models (Multi-LSTM, Multi-NBEATS, and CFA) are trained jointly on all source datasets, while uni-source models (LSTM, NBEATS) are trained only on the 1 source dataset whose sampling frequency matches the given target dataset. For Multi-NBEATS, in order to train the model over datasets of different context/forecast lengths, we add an encoder to the inputs and a decoder to the outputs, unique to each frequency, to enforce an equal latent backcast/forecast length for NBEATS. We evaluate each model by the zero-shot performance on target after training on source, using Normalized Deviation (ND) [11] as the evaluation metric.

The results are shown in Table 4. We find that the uni-source models always outperform their multi-source counterparts, and CFA has the worst performance for 2 out of 3 target datasets. This supports the conclusion that multi-source training deteriorates performance for zero-shot meta-forecasting. Further, learning frequency-invariant signal, as CFA does, is not suitable when source data is available of the same frequency as the target data. This is an area of potential improvement for future work.