

A decoder-only foundation model for time-series forecasting

Abhimanyu Das¹ Weihao Kong¹ Rajat Sen¹ Yichen Zhou¹

Abstract

Motivated by recent advances in large language models for Natural Language Processing (NLP), we design a time-series foundation model for forecasting whose out-of-the-box zero-shot performance on a variety of public datasets comes close to the accuracy of state-of-the-art supervised forecasting models for each individual dataset. Our model is based on pretraining a decoder style attention model with input patching, using a large time-series corpus comprising both real-world and synthetic datasets. Experiments on a diverse set of previously unseen forecasting datasets suggests that the model can yield accurate zero-shot forecasts across different domains, forecasting horizons and temporal granularities.

1. Introduction

Time-series data is ubiquitous in various domains such as retail, finance, manufacturing, healthcare and natural sciences. In many of these domains, one of the most important use-cases of time-series data is forecasting. Time-series forecasting is critical to several scientific and industrial applications, like retail supply chain optimization, energy and traffic prediction, and weather forecasting. In recent times, Deep learning models (Salinas et al., 2020; Oreshkin et al., 2019) have emerged as a popular approach for forecasting rich, multivariate, time-series data, often outperforming classical statistical approaches such as ARIMA or GARCH (Box & Jenkins, 1968). In several forecasting competitions such as the M5 competition (Makridakis et al., 2022) and IARAI Traffic4cast contest (Kopp et al., 2021), deep neural network based solutions showed good performance.

At the same time, we are witnessing a rapid progress in the Natural Language Processing (NLP) domain on large foundation models for downstream NLP tasks. Large language models (LLMs) are growing in popularity because they can be used to generate text, translate languages, write different kinds of creative content, and answer your questions in

an informative way (Radford et al., 2019). They are trained on massive amounts of data, which allows them to learn the patterns of human language. This makes them very powerful tools that can be used for a variety of downstream tasks, often in a zero-shot learning mode.

This motivates the question: “Can large pretrained models trained on massive amounts of time-series data learn temporal patterns that can be useful for time-series forecasting on previously unseen datasets?” In particular, can we design a time-series foundation model that obtains good zero-shot out-of-the-box forecasting performance? Such a pretrained time-series foundation model, if possible, would bring significant benefits for downstream forecasting users in terms of no additional training burden and significantly reduced compute requirements. It is not immediately obvious that such a foundation model for time-series forecasting is possible. Unlike in NLP, there is no well defined vocabulary or grammar for time-series. Additionally, such a model would need to support forecasting with varying history lengths (context), prediction lengths (horizon) and time granularities. Furthermore, unlike the huge volume of public text data for pretraining language models, vast amounts of time-series data is not readily available. In spite of these issues, we provide evidence to answer the above question in the affirmative.

In particular, we design *TimesFM*, a single foundation model for time-series forecasting that, when applied to a variety of previously-unseen forecasting datasets across different domains, obtains close to state-of-the-art zero-shot accuracy (compared to the best supervised models trained individually for these datasets). Our model can work well across different forecasting history lengths, prediction lengths and time granularities at inference time. The key elements of our foundation model are twofold: 1) a large-scale time-series corpus built using both real-world (mostly time-series data from web search queries¹ and Wikipedia page visits²) and synthetic data, which meets the volume and diversity of data needed for training our foundation model, and 2) a decoder style attention architecture with input patching, that can be efficiently pre-trained on this time-series corpus.

¹Google Research. Correspondence to: Rajat Sen <senrajat@google.com>.

Author names are listed in Alphabetical order.

¹<https://trends.google.com>

²https://wikimedia.org/api/rest_v1/

Compared to the latest large language models, our time-series foundation model is much smaller in both parameter size (200M parameters) and pretraining data size (100B timepoints); yet we show that even at such scales, it is possible to pretrain a practical foundation model for forecasting whose zero-shot performance comes close to the accuracy of fully-supervised approaches on a diverse set of time-series data. Our work also suggests that unlike recent work (Gruver et al., 2023) that recommends Large Language Models such as GPT-3 and LLaMA-2 as out-of-the-box zero-shot forecasters, foundation models trained from scratch exclusively on time-series data can obtain much better zero-shot performance at a tiny fraction of its costs.

2. Related Work

In the last decade, deep learning models (Salinas et al., 2020; Oreshkin et al., 2019) have emerged as powerful contenders in forecasting time-series in the presence of large training datasets and have been shown to outperform traditional statistical methods such as ARIMA and Exponential smoothing (McKenzie, 1984). Forecasting models can be categorized broadly into: (i) Local univariate models that include traditional methods like ARIMA, exponential smoothing (McKenzie, 1984) and non-autoregressive models like Prophet (Taylor & Letham, 2018). These models are trained individually for each time-series in a dataset in order to predict the corresponding time-series’s future. (ii) Global univariate models like DeepAR (Salinas et al., 2020), Temporal Convolutions (Borovykh et al., 2017), N-BEATS (Oreshkin et al., 2019) and long-term forecasting models such as (Nie et al., 2022; Das et al., 2023) that are trained globally on many time-series but during inference they predict the future of a time-series as a function of its own past and other related covariates. (iii) Global multivariate models that take in the past of all time-series in the dataset to predict the future of all the time-series. Such models include the classical VAR model (Zivot & Wang, 2006) as well as deep learning models like (Sen et al., 2019; Zhou et al., 2022; 2021) to name a few.

All the works cited above have primarily been applied in the supervised setting with the notable exception of PatchTST (Nie et al., 2022) and N-BEATS (Oreshkin et al., 2019). PatchTST has a section on dataset-to-dataset transfer learning in the semi-supervised setting. (Oreshkin et al., 2021) also show that the N-BEATS architecture lends itself to transfer learn between various source-target dataset pairs. However, none of these works aim to train a single foundation model that can work on a plethora of datasets. For an in-depth discussion about transfer learning in time-series we refer the reader to the survey in (Ma et al., 2023).

There has been some very recent work on re-using or fine-tuning large language models for time-series forecasting.

In particular, (Gruver et al., 2023) benchmarks pretrained LLMs like GPT-3 and LLaMA-2 for zero-shot forecasting performance. As we show later, our model obtains much superior zero-shot performance at a tiny fraction of these model sizes. Zhou et al. (2023) and (Chang et al., 2023) show how to fine-tune a GPT-2 (Radford et al., 2019) backbone model for time-series forecasting tasks. With the exception of a transfer-learning study (forecasting on a target dataset after having trained on a source dataset), these papers mostly focus on fine-tuning a pretrained model on target datasets, and not on pretraining a single foundation model with good out-of-the box zero-shot performance on a variety of datasets. To the best of our knowledge, the very recent work in TimeGPT-1 (Garza & Mergenthaler-Canseco, 2023) is the only other parallel work on a zero-shot foundation model for time-series forecasting. However the model is not public access, and several model details and the benchmark dataset have not been revealed.

3. Problem Definition

The task at hand is to build a general purpose zero-shot forecaster that takes in the past C time-points of a time-series as context and predicts the future H time-points. Let the context be denoted by $\mathbf{y}_{1:L} := \{y_1, \dots, y_L\}$ where we follow a numpy-like notation for indices. Similarly the actual values in the horizon are denoted by $\mathbf{y}_{L+1:L+H}$. Note that since we are building a single pre-trained model, we cannot have dataset specific dynamic or static covariates during training time. The task is then to learn a foundation model that can map any time-series context to horizon,

$$f : (\mathbf{y}_{1:L}) \longrightarrow \hat{\mathbf{y}}_{L+1:L+H}. \quad (1)$$

The accuracy of the prediction can be measured by a metric that quantifies their closeness to the actual values, for instance, Mean Absolute Error (MAE) defined in Equation 6.

4. Model Architecture

A foundation model for time-series forecasting should be able to adapt to variable context and horizon lengths, while having enough capacity to encode all patterns from a large pretraining datasets. Transformers have been shown to be able to adapt to different context lengths in NLP (Radford et al., 2019). However, there are several time-series specific design choices. The main guiding principles for our architecture are the following:

Patching. Inspired by the success of patch based modeling in the recent long horizon forecasting work (Nie et al., 2022) we also choose to break down the time-series into patches during training. A patch of a time-series is a natural analogue for a token in language models and patching has been shown to improve performance. Moreover this improves inference speed as the number of tokens being

fed into the transformer is reduced by a factor of the patch length. On the other hand, increasing the patch length all the way to the context length moves us away from decoder-only training and the efficiencies that come with it. We delve into this further in Section 6.2.

Decoder-only model. A key difference between our architecture and PatchTST (Nie et al., 2022) is that our model is trained in decoder-only mode (Liu et al., 2018). In other words, given a sequence of input patches, the model is optimized to predict the next patch as a function of all past patches. Similar to LLMs this can be done in parallel over the entire context window, and automatically enables the model to predict the future after having seen varying number of input patches.

Longer output patches. In LLMs the output is always generated in an auto-regressive fashion one token at a time. However, in long-horizon forecasting it has been observed that directly predicting the full horizon yields better accuracy than multi-step auto-regressive decoding (Zeng et al., 2023). But this is not possible when the horizon length is not known apriori, as in the case of zero-shot forecasting which is our primary goal.

We propose a middle ground by allowing our output patches for prediction to be longer than the input patches. As an example, suppose the input patch length is 32 and output patch length is 128. During training, the model is simultaneously trained to use the first 32 time-points to forecast the next 128 time-steps, the first 64 time-points to forecast time-steps 65 to 192, the first 96 time-points to forecast time-steps 97 to 224 and so on. During inference, suppose the model is given a new time-series of length 256 and tasked with forecasting the next 256 time-steps into the future. The model will first generate the future predictions for time-steps 257 to 384, then condition on the initial 256 length input plus the generated output to generate time-steps 385 to 512. On the other hand, if in a model the output patch length was fixed to the input patch length of 32, then for the same task we would have to go through 8 auto-regressive generation steps instead of just the 2 above. However, there is a trade-off. If the output patch length is too long, then it is difficult to handle time-series whose lengths are less than the output patch length for instance monthly, yearly time-series in our pretraining data.

Patch Masking. If we use patches naively, the model might only learn to predict well for context lengths that are multiples of the input patch length. Therefore we make a careful use of masking during training. Parts of patches as well as entire patches from the beginning of the context window can be masked in a data batch. We employ a specific random masking strategy (described later) during training that helps the model see all possible context lengths starting from 1 to a maximum context length.

Now that we have mentioned the guiding principles, we next formally describe each component of our model architecture (illustrated in Figure 1), which we name as TimesFM (Time-series Foundation Model).

Input Layers. The job of the input layers is to preprocess the time-series into input tokens to the transformer layers. We first break the input into contiguous non-overlapping patches. Then each patch is processed by a Residual Block into a vector of size `model_dim`. Along with the input, we also supply a binary padding mask $\mathbf{m}_{1:L}$ where 1 denotes that the corresponding input in $\mathbf{y}_{1:L}$ should be ignored and vice-versa. The Residual Block is essentially a Multi-layer Perceptron (MLP) block with one hidden layer with a skip connection, similar to that defined in (Das et al., 2023).

In other words, the inputs $\mathbf{y}_{1:L}$ are broken down into patches of size `input_patch_len` (p). The j -th patch can be denoted as $\tilde{\mathbf{y}}_j = \mathbf{y}_{p(j-1)+1:pj}$. Similarly the mask can also be patched as $\tilde{\mathbf{m}}_j = \mathbf{m}_{p(j-1)+1:pj}$. Then the j -th input token to the subsequent transformer layers can be denoted as,

$$\mathbf{t}_j = \text{InputResidualBlock}(\tilde{\mathbf{y}}_j \odot (1 - \tilde{\mathbf{m}}_j)) + \text{PE}_j \quad (2)$$

where PE_j denotes the j -th positional encoding as defined in the original transformer paper (Vaswani et al., 2017). There will be $N = \lfloor L/p \rfloor$ such input tokens.

Stacked Transformer. The bulk of the parameters in our model are in `num_layers` (n_l) transformer layers stacked on top of each other. Each of these layers have the standard multi-head self-attention (SA) followed by a feed-forward network (FFN). The main hyperparameters are `model_dim` which is equal to the dimension of the input tokens \mathbf{t}_j 's and number of heads (`num_heads`). We set the hidden size of the FFNs to be equal to `model_dim` as well. We use causal attention that is each output token can only attend to input tokens that come before it in the sequence (including the corresponding input token). This can be described by the equation

$$\mathbf{o}_j = \text{StackedTransformer}((\mathbf{t}_1, \dot{m}_1), \dots, (\mathbf{t}_j, \dot{m}_j)), \quad (3)$$

for all $j \in [N]$. \dot{m}_j is the masking indicator for the j -th token defined as $\min\{\mathbf{m}_{p(j-1)+1:pj}\}$ i.e if a patch has any non-masked time-point the corresponding token marked as not being masked. All patches that are masked out completely are not attended to by the causal self attention.

Output Layers. The remaining task is to map the output tokens into predictions. We train in decoder only mode i.e each output token should be able to be predictive of the part of the time-series that follows the last input patch corresponding to it. This is common for popular large

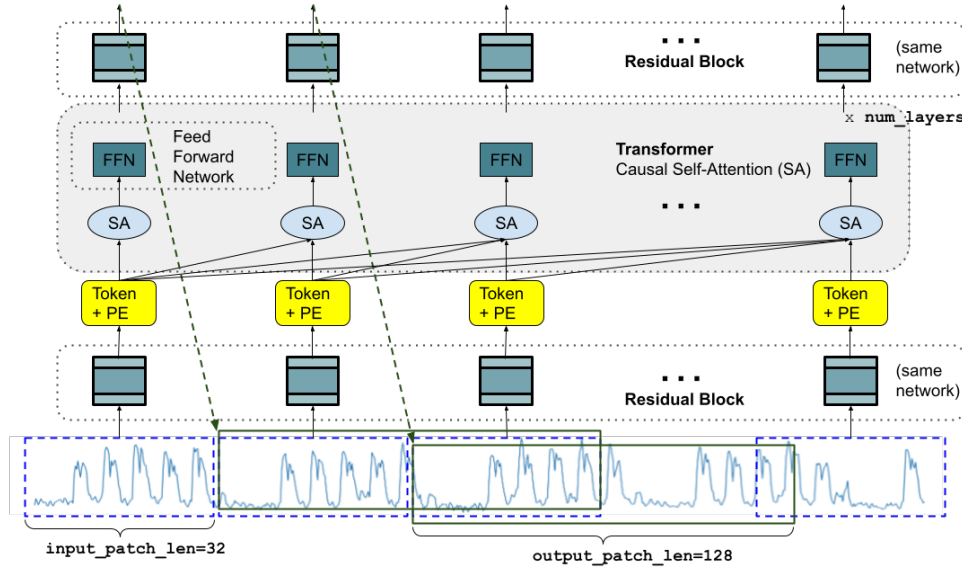


Figure 1. We provide an illustration of the TimesFM model architecture during training, where we show a input time-series of a specific length that can be broken down into input patches. Each patch along is processed into a vector by a residual block (as defined in the model definition) to the model dimension of the transformer layers. The vector is then added to positional encodings and fed into n_l stacked transformer layers. SA refers to self-attention (note that we use multi-head causal attention) and FFN is the fully connected layer in the transformer. The output tokens are then mapped through a residual block to an output of size `output_patch_len`, which is the forecast for the time window following the last input patch seen by the model so far.

language models like (Radford et al., 2019). However, one key difference in our time-series foundation model is that input patch length need not be equal to output patch length i.e we should be able to predict a larger chunk of the time-series based on the encoded information from the input patches seen so far. Let the output patch length be `output_patch_len` (h). We use another **Residual Block** to map the output tokens to the predictions. This can be described as,

$$\hat{\mathbf{y}}_{pj+1:pj+h} = \text{OutputResidualBlock}(\mathbf{o}_j). \quad (4)$$

Thus we encode all the data in $\mathbf{y}_{1:pj}$ into \mathbf{o}_j and use that to predict the subsequent h time-points $\mathbf{y}_{pj+1:pj+h}$. This is done for all patches in one training mini-batch.

Loss Function. In this work, we focus on point forecasting. Therefore we can use a point forecasting loss during training like Mean Squared Error (MSE). The loss that is minimized during training can be expressed as,

$$\text{TrainLoss} = \frac{1}{N} \sum_{j=1}^N \text{MSE}(\hat{\mathbf{y}}_{pj+1:pj+h}, \mathbf{y}_{pj+1:pj+h}). \quad (5)$$

Note that if one is interested in probabilistic forecasting, then it is easy to have multiple output heads for each output patch, each head minimizing a separate quantile loss

as in (Wen et al., 2017). Another approach can be to output the logits of a probability distribution family and minimize the maximum likelihood loss for probabilistic forecasting (Awasthi et al., 2021; Salinas et al., 2020).

Training. We train the model with standard mini-batch gradient descent in decoder-only fashion, that goes through all windows for a time-series and across time-series. The only non-standard part is the way we sample the mask during training. For each time-series in the batch, we sample a random number r between 0 and $p - 1$. Then we set the $\mathbf{m}_{1:r} = 1$ and the rest as zero i.e we mask our a fraction of the first input patch. However, this is sufficient to cover all input context lengths from 1 to the maximum training context length. We explain this using an example below:

Suppose the maximum context length is 512 and $p = 32$. Then if $r = 4$, the output prediction after seeing the first patch (from \mathbf{o}_1) is optimized to predict after seeing $28 = 32 - 4$ time-points, the output of the next patch (from \mathbf{o}_2) is optimized to predict after seeing $28 + 32$ time-points, and so on. When this argument is repeated for all such r 's, the model has seen all possible context lengths till 512.

Inference. The trained network can be used to produce forecasts for *any* horizon using auto-regressive decoding similar to large language models. Given an input $\mathbf{y}_{1:L}$ (assume L is a multiple of p for simplicity) it can first pre-

dict $\hat{y}_{L+1:L+h}$. Then, we can use the concatenated vector $\tilde{y}_{1:L+h} = [y_{1:L}; \hat{y}_{L+1:L+h}]$ as an input to the network to generate the next output patch prediction $\hat{y}_{L+h+1:L+2h}$ and so on. If L is not a multiple of p , we simply append zeros to make it a multiple of p and mark the corresponding entries in the mask as 1.

5. Pretraining Details

We would like our pretraining corpus to include large volumes of temporal data representing a variety of domains, trend and seasonality patterns and time granularities that ideally capture the forecasting use-cases which we are interested in serving by the deployed model. It is challenging to find a large time-series dataset that meets the volume and diversity of data needed for training our foundation model. We address this problem by sourcing the bulk of data used to train our models from three major sources: Google trends, Wiki Pageview statistics and synthetic time-series. In summary the main data sources are:

Google Trends. Google Trends ³ captures search interest over time for millions of queries. We choose around 22k head queries based on their search interest over 15 years from 2007 to 2022. Beyond these head queries the time-series become more than 50% sparse. We download the search interest over time for these queries in hourly, daily, weekly and monthly granularities to form our dataset. The date ranges are Jan. 2018 to Dec. 2019 for hourly and Jan. 2007 to Dec. 2021 for the other granularities. The trends datasets amounts to roughly 1B time-points.

Wiki Pageviews. Wiki Pageviews ⁴ captures the hourly views of all Wikimedia pages. We download all pageview data from Jan. 2012 to Nov. 2023, clean and aggregate the views by page into hourly, daily, weekly and monthly granularities, and filter out pageview time-series with excessive zeros. The final corpus contains roughly 100B time-points.

Synthetic Data. Another major component of our pretraining data is of synthetic origin. We create generators for ARMA (McKenzie, 1984) processes, seasonal patterns (mixture of sines and cosines of different frequencies), trends (linear, exponential with a few change-points) and step functions. A synthetic time-series can be an additive combination of one or more of these processes. We create 3M synthetic time-series each of length 2048 time-points. More details about our synthetic data generation are presented in Appendix A.5.

Other real-world data sources. Along with the wiki and trends data, we also add time-series from several other publicly available datasets to our pretraining corpus. We add

³<https://trends.google.com>

⁴https://en.wikipedia.org/wiki/Wikipedia:Pageview_statistics

all the granularities of the M4 dataset (Makridakis et al., 2022), the hourly and 15 minute Electricity and the hourly Traffic datasets (see (Zhou et al., 2021)). We also add the 10-minute granularity Weather dataset used for evaluations in (Zhou et al., 2021). M4 has a good mix of granularities with around 100k time-series in total. Traffic and Electricity are large long-term forecasting datasets with > 800 and > 300 time-series each having tens of thousands of time-points. In addition, we add all the 15 min granularity traffic time-series from (Wang et al., 2023).

Dataset Mixing and Training. We train on a mixture distribution over these datasets that aims to give sufficient weight to all granularities. The training loader samples 40% real data and 60% synthetic, with the real data mixture providing equal weights to all hourly + sub-hourly, daily, weekly, and monthly datasets. We train with a maximum context length of 512 whenever the length of the time-series allows that. For weekly granularity we do not have sufficiently long time-series; therefore a maximum context length of 256 is used. For the same reason, a maximum context length of 64 is used while training on \geq monthly granularity data. We also use only the standard normalization part of reversible instance normalization (Kim et al., 2021) – i.e, the context of each time-series is scaled by the context mean and standard deviation of the first input patch in the context.

6. Empirical Results

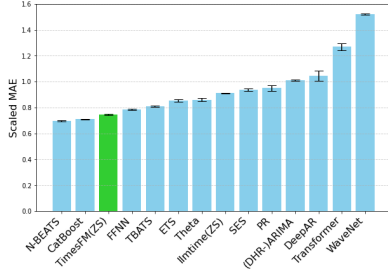
We evaluate our model in zero-shot settings on three groups of well known public datasets against the best performing baselines for each group. These datasets have been intentionally held out from our pretraining data. We show that a *single* pretrained model can come close or surpass the performance of baselines models on the benchmarks even when the baselines are specially trained or tuned for each specific task. Subsequently, we perform ablation studies that justify different choices made in our architecture.

6.1. Zero-shot Evaluation

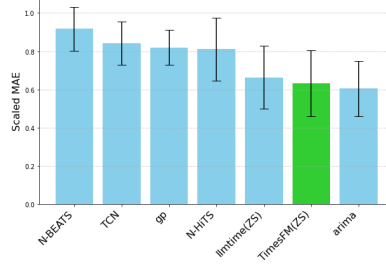
To benchmark our model’s performance, we choose three groups of commonly used forecasting datasets that cover various domains, sizes, granularities, and horizon lengths: Darts (Herzen et al., 2022), Monash (Godaheva et al., 2021) and Informer datasets (Zhou et al., 2021), to test the generalization power of our foundation model against other baselines.

In all cases, we report performance on the official metrics and scalings of the datasets. We present a summary of the results below - more detailed experimental results can be found in Appendix A.2. We provide the hyper-parameters and other details about our model in Appendix A.3.

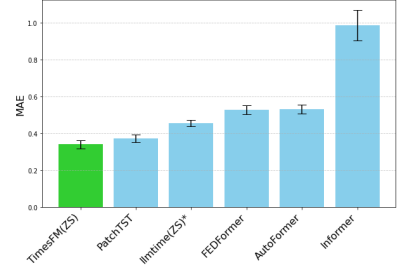
Monash (Godaheva et al., 2021). Monash archive is a



(a) Monash Archive (Godaheva et al., 2021)



(b) Darts (Herzen et al., 2022)



(c) ETT (Horizons 96 and 192) (Zhou et al., 2021)

Figure 2. We report average performance in three groups of datasets. In all figures, the lower the metric the better and the error bars represent one standard error. Note that among the baselines only TimesFM and Ilmtime are zero-shot. In (a) we report results on the Monash datasets. Since the datasets have different scales, we average the MAE’s scaled by the MAE of a naive baseline. We can see that TimesFM is among the top-3 models. In (b), we report the similarly scaled MAE on the Darts benchmarks. TimesFM is within significance of the top of the best performing method which is ARIMA in this case. Note that these datasets have one time-series each and therefore statistical methods are competitive with deep learning ones. Finally, in (c) we report the average MAE for 96 and 192 horizon prediction tasks on 4 ETT datasets i.e 8 tasks in total. TimesFM and PatchTST are the best performing models in this case.

collection of 30 datasets of different training and prediction lengths that covers granularities ranging from minutes to years and domains including finance, demand forecasting, weather and traffic. The archive reports four official metrics for several statistical baselines such as Exponential Smoothing(ETS) and ARIMA, as well as supervised ML baselines like CatBoost (Prokhorenkova et al., 2018), DeepAR (Salinas et al., 2020) and WaveNet (Oord et al., 2016). Following Ilmtime (Gruver et al., 2023) we start from the Monash Huggingface repository⁵ and filter out the datasets that contain missing values. This leaves us with 18 datasets which we specify in Appendix A.2.2.

Out of the four official metrics, following prior work (Gruver et al., 2023), we report our performance in terms of mean MAE (see Appendix A.1). As the datasets have massively different scales, for each dataset we normalize the metric by the metric achieved by a naive baseline that just constantly predicts the last value in the context for each time-series. Then the scaled MAE’s are averaged across all datasets. The scaled aggregation was also used in (Gruver et al., 2023).

The mean scaled MAE across all datasets is plotted in Figure 2a along with standard error bars. We compare the performance of TimesFM with the baseline models implemented in Monash, and the zero-shot Ilmtime (Gruver et al., 2023) model that uses GPT-3 (Radford et al., 2019) with a specific prompting technique. Note that the zero-shot models are marked as (Zero-Shot). TimesFM is among the top 3 models even though we never trained on these datasets. It outperforms deep supervised models like

DeepAR and FFNN (Godaheva et al., 2021), and improves on Ilmtime’s performance by more than 10%.

Darts (Herzen et al., 2022). This is a collection of 8 univariate datasets which include interesting seasonalities and additive+multiplicative trends. We report performance of several baselines implemented in the Darts package like TCN (Lea et al., 2016), N-HITS (Challu et al., 2023) and N-BEATS (Oreshkin et al., 2019). All these baselines are supervised. As before, we also report zero-shot forecasting results from Ilmtime (Gruver et al., 2023) using GPT-3 (Radford et al., 2019). Other supervised baselines in (Gruver et al., 2023) like SM-GP (Wilson & Adams, 2013) and ARIMA (McKenzie, 1984) are also added.

We report the official metric for this dataset group that is MAE for each individual dataset in Appendix A.2. In Figure 2b, we present the average scaled MAE across all 8 datasets, as we did for the Monash datasets. TimesFM is within statistical significance of the best model that is seasonal ARIMA in this case. Note that since there are only 8 individual time-series in this dataset group, the standard errors are not sharp and therefore does not provide a clear ordering among the models. Also, note that for ARIMA, the seasonality needs to be encoded correctly in the parameters for the best results, which needed manual tuning.

Informer (Zhou et al., 2021). The Informer datasets have been widely used for benchmarking various supervised long-horizon forecasting methods. A few of these datasets are used in pretraining, so we focus on the other datasets in this collection (ETTM1, ETTM2, ETTh1 and ETTh2) related to electricity transformer temperatures over a two year period in 1 hour and 15 minutes granularities. Note that the long horizon baselines usually report rolling validation re-

⁵https://huggingface.co/datasets/monash_tsf

sults on the test set which would amount to millions of tokens for evaluating lltime (Gruver et al., 2023) and would be too expensive. Therefore, following lltime, we compare all methods on the last test window. Also, it is reasonable to directly average the MAE for these datasets since the results are reported on standard normalized dataset (using the statistics of the training portion).

We consider the task of predicting horizon length 96 and 192, given a context length of 512 for all methods. The MAE averaged over all 8 tasks (4 datasets with two horizons each) is presented in Figure 2b. TimesFM performs the best and the supervised PatchTST (Nie et al., 2022) baseline (which is a state-of-the-art long horizon deep forecasting method) is within significance of it. The other long horizon methods are quite a bit worse even though they have been trained these datasets. lltime is better than FEDFormer but worse than PatchTST in a statistically significant way.

Visualizing Forecasts. Next, we conduct a visual inspection of the forecasts generated by TimesFM, first on some synthetic examples and then on the benchmark datasets.

In Figure 3 we show 4 different synthetic curves: (1) sum of 5 sine curves of different periods, (2) a sine curve linearly scaled, (3) a sine curve with a linear trend, and (4) minimum of two sine curves with a linear trend. Our results suggests that TimesFM picks up the trend and seasonal components readily interpretable by humans, while ARIMA and (to a lesser extent) lltime fail in some of the instances.

As illustrated in Figure 4, TimesFM also effectively captures these subtle characteristics within both the trend and seasonal patterns of the depicted real world time-series. For instance, in the Air Passenger dataset, TimesFM correctly captures the amplitude increase with trend –this is also reflected by the fact that it attains the best MAE on this dataset (see Table 1). In the traffic hourly example on the left, it can be seen that TimesFM can correctly identify the seasonal peaks even in the presence of outliers in the context, while lltime is thrown off.

More examples are presented in Appendix A.2.

6.2. Ablation

Next, we perform several ablation studies that inform the design decisions we made for our model architecture.

Scaling. Performance curves with respect to number of parameters in a model have been a keenly studied area in the context of LLMs. Kaplan et al. (2020) established a power law like relationship between the number of parameters in a language model and its downstream performance i.e the more the number of parameters the better the performance. However, Hoffmann et al. (2022) established a more nu-

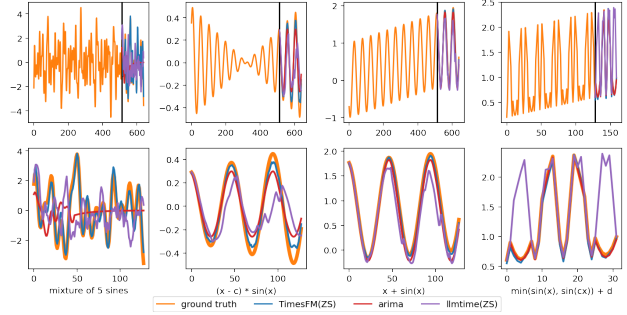


Figure 3. Forecasts visualized on synthetic curves. The bottom row plots zoom in on the prediction horizon for the sake of clarity.

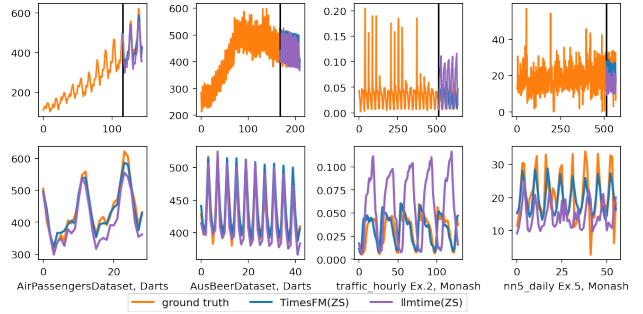


Figure 4. Forecasts visualized on Darts and Monash. The bottom row plots zoom in on the prediction horizon for the sake of clarity.

anced scaling law that lays down methods to train compute optimal models based on the number of tokens available in a training dataset. We perform a preliminary scaling study where we train three TimesFM models of sizes 17M, 70M and 200M parameters, using the same pretraining dataset and till the similar number of iterations.

We provide the mean scaled MAE results on the Monash datasets in Figure 5. It can be clearly seen that the errors decrease monotonically with the number of parameters, thus pointing to a similar observation as (Kaplan et al., 2020), but in the context of time-series.

Autoregressive Decoding. In recent long-term forecasting works (Zeng et al., 2023; Nie et al., 2022; Das et al., 2023) it has been observed that directly predicting the entire forecasting horizon in one shot from a decoder can yield better results than auto-regressive decoding on long horizon benchmarks. For a foundation model, the horizon length of the task is not known before inference time, therefore one-shot decoding might not be possible for very long horizons. However, as mentioned earlier, by keeping the output_patch.len longer than input_patch.len one can ensure fewer autoregressive steps. This was one of the key decisions in the design of TimesFM, that is quite different from LLMs. In order to showcase this we choose the task of predicting 512 time-steps into the fu-

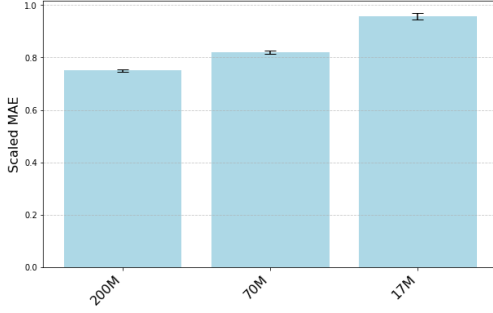


Figure 5. Average scaled MAE on Monash datasets for three different TimesFM model sizes. We can see that the performance improves with increasing model size.

ture for the ETT datasets. In Figure 6, we present results from a model with `output_patch_len=32` vs our original model that uses `output_patch_len=128`. The former has to perform 16 autoregressive steps while the latter has to do only 4. It can be clearly seen that having a larger `output_patch_len` improves performance in this case.

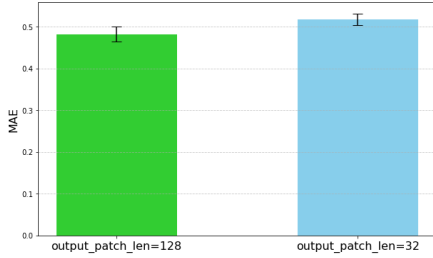


Figure 6. Ablation with respect to output patch length for the task of predicting 512 steps into the future, in the same setting as that of Section 6.1. We report the average across all ETT 4 datasets.

Input Patch Length. The size of `input_patch_len` represents an important trade-off. We have typically seen that increasing its value from 8 to 32 increases performance but having too high a `input_patch_len` is impractical since that makes the model shift from decoder only training more towards encoder-decoder style training. Note that in the "Training" paragraph of Section 4, we describe the mask sampling strategy to support any context length. If in the extreme case p is set the maximum context length we have to individually sample all possible context windows from 1 to maximum context length, which would be required for encoder-decoder style of training.

In Figure 7 we show the mean scaled MAE TimesFM(ZS) - 70M model with `input_patch_len=8` on Monash datasets, which is clearly worse than our original model that uses `input_patch_len=32`.

Importance of Dataset Size. Next we showcase the importance of our large pretraining datasets for the perfor-

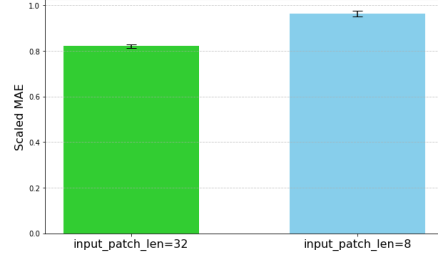


Figure 7. Average scaled MAE for our 70M models on Monash datasets for two different input patch lengths. We can see that the performance improves with increasing patch size.

mance of our model. We train two models with the same model size and hyperparameters, but with different pre-training data. The first model we trained include all the datasets mentioned in Section 5. For the second model, we trained only the smaller real-world datasets (M4, Electricity, Traffic and Weather). As seen by the results in Figure 8 on the Monash benchmark, pretraining on all the datasets clearly results in a significant improvement in model quality, compared to using the smaller real-world datasets.

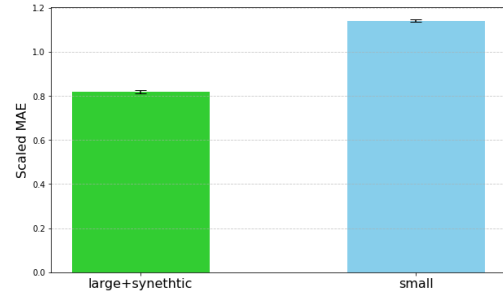


Figure 8. Average scaled MAE on Monash datasets for our 70M models with two pretraining settings. Large+synthetic model is pretrained on all the datasets mentioned in Section 5; Small is pretrained on M4, Electricity, Traffic and Weather only.

7. Conclusion

In this paper, we presented TimesFM, a practical foundation model for forecasting whose zero-shot performance comes close to the accuracy of fully-supervised forecasting models on a diverse set of time-series data. This model is pretrained on real-world and synthetic datasets comprising around 100B timepoints using a patched-decoder style attention architecture with around 200M parameters.

In future work, we plan to delve into a more theoretical understanding of how such a time-series foundation model can obtain good performance for out-of-distribution data, and also investigate the fine-tuning/few-shot performance of this model.

8. Impact Statement

This paper presents work whose goal is to advance the field of Time-Series Forecasting using Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Awasthi, P., Das, A., Sen, R., and Suresh, A. T. On the benefits of maximum likelihood estimation for regression and forecasting. *arXiv preprint arXiv:2106.10370*, 2021.
- Borovykh, A., Bohte, S., and Oosterlee, C. W. Conditional time series forecasting with convolutional neural networks. *arXiv preprint arXiv:1703.04691*, 2017.
- Box, G. E. and Jenkins, G. M. Some recent advances in forecasting and control. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 17(2):91–109, 1968.
- Challu, C., Olivares, K. G., Oreshkin, B. N., Garza, F., Mergenthaler, M., and Dubrawski, A. NHITS: Neural Hierarchical Interpolation for Time Series forecasting. In *The Association for the Advancement of Artificial Intelligence Conference 2023 (AAAI 2023)*, 2023. URL <https://arxiv.org/abs/2201.12886>.
- Chang, C., Peng, W.-C., and Chen, T.-F. Llm4ts: Two-stage fine-tuning for time-series forecasting with pre-trained llms. *arXiv preprint arXiv:2308.08469*, 2023.
- Das, A., Kong, W., Leach, A., Mathur, S. K., Sen, R., and Yu, R. 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>.
- Garza, A. and Mergenthaler-Canseco, M. Timegpt-1. *arXiv preprint arXiv:2310.03589*, 2023.
- Godahehwa, R., Bergmeir, C., Webb, G. I., Hyndman, R. J., and Montero-Manso, P. Monash time series forecasting archive. *arXiv preprint arXiv:2105.06643*, 2021.
- Gruver, N., Finzi, M., Qiu, S., and Wilson, A. G. Large language models are zero-shot time series forecasters. *arXiv preprint arXiv:2310.07820*, 2023.
- Herzen, J., Lässig, F., Piazzetta, S. G., Neuer, T., Tafti, L., Raille, G., Van Pottelbergh, T., Pasieka, M., Skrodzki, A., Huguenin, N., et al. Darts: User-friendly modern machine learning for time series. *The Journal of Machine Learning Research*, 23(1):5442–5447, 2022.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D. d. L., Hendricks, L. A., Welbl, J., Clark, A., et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Kim, T., Kim, J., Tae, Y., Park, C., Choi, J.-H., and Choo, J. Reversible instance normalization for accurate time-series forecasting against distribution shift. In *International Conference on Learning Representations*, 2021.
- Kopp, M., Kreil, D., Neun, M., Jonietz, D., Martin, H., Herruzo, P., Gruca, A., Soleymani, A., Wu, F., Liu, Y., Xu, J., Zhang, J., Santokhi, J., Bojesomo, A., Marzouqi, H. A., Liatsis, P., Kwok, P. H., Qi, Q., and Hochreiter, S. Traffic4cast at neurips 2020 - yet more on the unreasonable effectiveness of gridded geospatial processes. In Escalante, H. J. and Hofmann, K. (eds.), *Proceedings of the NeurIPS 2020 Competition and Demonstration Track*, volume 133 of *Proceedings of Machine Learning Research*, pp. 325–343. PMLR, 06–12 Dec 2021. URL <https://proceedings.mlr.press/v133/kopp21a.html>.
- Lea, C., Vidal, R., Reiter, A., and Hager, G. D. Temporal convolutional networks: A unified approach to action segmentation. In *Computer Vision—ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part III 14*, pp. 47–54. Springer, 2016.
- Liu, P. J., Saleh, M., Pot, E., Goodrich, B., Sepassi, R., Kaiser, L., and Shazeer, N. Generating wikipedia by summarizing long sequences. *arXiv preprint arXiv:1801.10198*, 2018.
- Ma, Q., Liu, Z., Zheng, Z., Huang, Z., Zhu, S., Yu, Z., and Kwok, J. T. A survey on time-series pre-trained models. *arXiv preprint arXiv:2305.10716*, 2023.
- Makridakis, S., Spiliotis, E., and Assimakopoulos, V. M5 accuracy competition: Results, findings, and conclusions. *International Journal of Forecasting*, 38(4):1346–1364, 2022.
- McKenzie, E. General exponential smoothing and the equivalent arma process. *Journal of Forecasting*, 3(3): 333–344, 1984.
- Nie, Y., Nguyen, N. H., Sinthong, P., and Kalagnanam, J. A time series is worth 64 words: Long-term forecasting with transformers. *International conference on learning representations*, 2022.

- Oord, A. v. d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- Oreshkin, B. N., Carpo, D., Chapados, N., and Bengio, Y. N-beats: Neural basis expansion analysis for interpretable time series forecasting. In *International Conference on Learning Representations*, 2019.
- Oreshkin, B. N., Carpo, D., Chapados, N., and Bengio, Y. Meta-learning framework with applications to zero-shot time-series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 9242–9250, 2021.
- Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., and Gulin, A. Catboost: unbiased boosting with categorical features. *Advances in neural information processing systems*, 31, 2018.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Salinas, D., Flunkert, V., Gasthaus, J., and Januschowski, T. Deepar: Probabilistic forecasting with autoregressive recurrent networks. *International Journal of Forecasting*, 36(3):1181–1191, 2020.
- Sen, R., Yu, H.-F., and Dhillon, I. S. Think globally, act locally: A deep neural network approach to high-dimensional time series forecasting. *Advances in neural information processing systems*, 32, 2019.
- Taylor, S. J. and Letham, B. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Wang, J., Jiang, J., Jiang, W., Han, C., and Zhao, W. X. Towards efficient and comprehensive urban spatial-temporal prediction: A unified library and performance benchmark. *arXiv preprint arXiv:2304.14343*, 2023.
- Wen, R., Torkkola, K., Narayanaswamy, B., and Madeka, D. A multi-horizon quantile recurrent forecaster. *arXiv preprint arXiv:1711.11053*, 2017.
- Wilson, A. and Adams, R. Gaussian process kernels for pattern discovery and extrapolation. In *International conference on machine learning*, pp. 1067–1075. PMLR, 2013.
- Wu, H., Xu, J., Wang, J., and Long, M. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in Neural Information Processing Systems*, 34:22419–22430, 2021.
- Zeng, A., Chen, M., Zhang, L., and Xu, Q. Are transformers effective for time series forecasting? *Proceedings of the AAAI conference on artificial intelligence*, 2023.
- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, 2021.
- Zhou, T., Ma, Z., Wen, Q., Wang, X., Sun, L., and Jin, R. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International Conference on Machine Learning*, pp. 27268–27286. PMLR, 2022.
- Zhou, T., Niu, P., Wang, X., Sun, L., and Jin, R. One fits all: Power general time series analysis by pretrained lm. *arXiv preprint arXiv:2302.11939*, 2023.
- Zivot, E. and Wang, J. Vector autoregressive models for multivariate time series. *Modeling financial time series with S-PLUS®*, pp. 385–429, 2006.

A. Appendix

A.1. Metrics

The metrics that are used for reporting results in this paper are:

- MAE (Godaheewa et al., 2021)

$$\text{MAE}(\mathbf{y}_{L+1:L+H}, \hat{\mathbf{y}}_{L+1:L+H}) = \frac{1}{H} \|\mathbf{y}_{L+1:L+H} - \hat{\mathbf{y}}_{L+1:L+H}\|_1. \quad (6)$$

- msMAPE (Godaheewa et al., 2021)

$$\text{msMAPE}(\mathbf{y}_{L+1:L+H}, \hat{\mathbf{y}}_{L+1:L+H}) = \frac{1}{H} \sum_{i=1}^H \frac{2|y_{L+i} - \hat{y}_{L+i}|}{\max\{|y_{L+i}| + |\hat{y}_{L+i}| + \epsilon, 0.5 + \epsilon\}}. \quad (7)$$

In Monash benchmarks (Godaheewa et al., 2021) $\epsilon = 0.1$ was used. This metric is used in order to avoid undefined values in other normalized metrics like MAPE. In multivariate datasets the metrics are calculated for each time-series and then we take the mean or the median. In this paper we only use the mean versions.

Aggregating across datasets. Since the datasets have wildly different scales averaging unnormalized metrics like MAE is not kosher. Therefore following (Gruver et al., 2023) we scale the metric of each baseline for a dataset by the same metric achieved by a naive baseline on that dataset. The naive baseline just makes the constant prediction y_L repeated across the prediction length. We did not need to do that for the Informer datasets since on these datasets metrics are usually reported on standard normalized data (Nie et al., 2022).

A.2. Additional Empirical Results

In this section, we provide more detailed tables for our zero-shot datasets and experiments described in Section 6.1

A.2.1. DARTS

We present the MAE results individually from all 8 datasets in Table 1. It can be seen that TimesFM performs well for all datasets with clear seasonal patterns. We do not perform the among the top-3 only in Sunspot and Woolly datasets. On an average we are within significant level of the best model and a close second on the mean scaled MAE metric.

In Figure 9 we present visual comparisons of our forecasts vs some of the baselines.

Table 1. MAE for Darts datasets. We also include the naive baseline that predicts the last values in the context repeatedly.

	gp	arima	TCN	N-BEATS	N-HiTS	llmtime(ZS)	TimesFM(ZS)	NAIVE
AirPassengersDataset	34.67	24.03	54.96	97.89	59.16	34.37	14.75	81.45
AusBeerDataset	102.05	17.13	30.90	10.39	34.23	16.13	10.25	96.35
GasRateCO2Dataset	2.27	2.37	2.64	2.63	3.85	3.50	2.69	2.29
MonthlyMilkDataset	30.33	37.19	70.86	33.64	32.73	9.68	22.46	85.71
SunspotsDataset	53.74	43.56	51.82	73.15	49.93	47.34	50.88	48.24
WineDataset	4552.06	2306.70	3287.14	4562.02	3909.51	1569.32	2462.11	4075.28
WoollyDataset	649.98	588.78	1158.79	903.01	382.09	808.73	917.10	1210.33
HeartRateDataset	5.65	5.56	5.49	6.57	6.10	6.21	5.44	5.92

A.2.2. MONASH

In Table 2 we present the actual MAE numbers that are behind the main Figure 2a.

It should be noted that though the scaled MAE metric accounts for scale variation across datasets, it does not account for within-dataset scale variations. Indeed, in datasets like cif 2016 and covid deaths, there are time-series with different order of magnitudes. The mean MAE for these datasets does not reflect average performance but only the performance of the higher scale time-series. Therefore in Table 3 we also present the normalized msMAPE metrics for TimesFM and all baselines.

In Figure 10, we present some examples of our zero-shot forecasts.

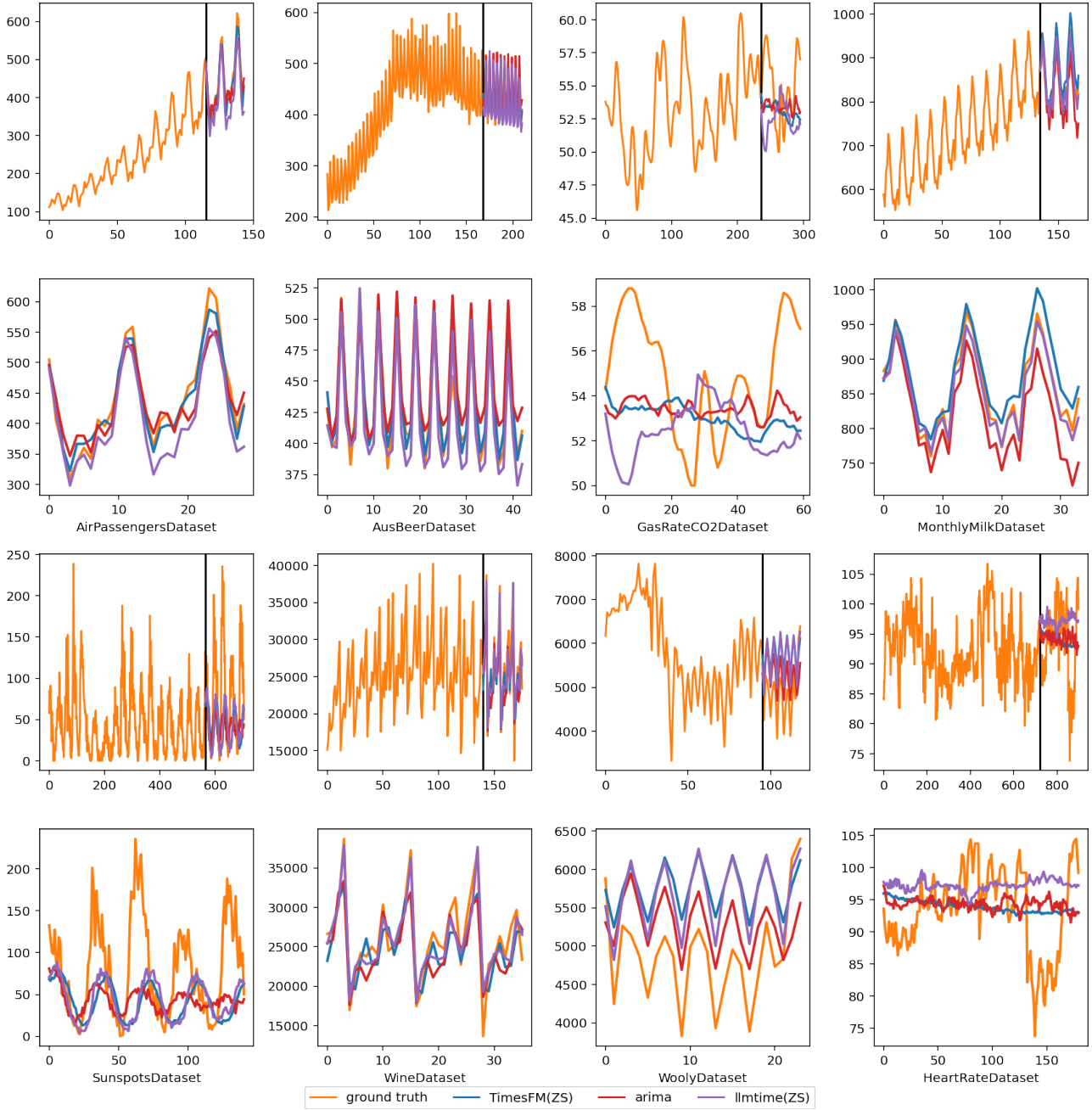


Figure 9. Forecasts visualized on all Darts datasets. The second row plots zoom in on the prediction horizon for the sake of clarity.

A.2.3. INFORMER

We present the MAE on the last split of the test set for all dataset, horizon pairs considered in Table 4. Owing to expensive evaluations for llmtime, the results are reported on the last test window of the original test split, as done in (Gruver et al., 2023).

A.2.4. SYNTHETIC DATA VISUALIZATIONS

In this section we aim to investigate how TimesFM generalizes to some common temporal patterns that are potentially out of distribution from the synthetic parts of its pretraining dataset. We present some examples in Figure 11. We also compare how ARIMA and llmtime behaves on these examples.

Table 2. We present the mean MAE results for our methods along size monash baselines. We also include the naive baseline that predicts the last values in the context repeatedly.

Dataset	llmtime(ZS)	SES	Theta	TBATS	ETS	(DHR-)ARIMA	PR	CatBoost	FFNN	DeepAR	N-BEATS	WaveNet	Transformer	TimesFM(ZS)	NAIVE
australian electricity demand	459.96	659.60	665.04	370.74	1282.99	1045.92	247.18	241.77	258.76	302.41	213.83	227.50	231.45	523.12	659.60
bitcoin	1.75e18	5.33e18	5.33e18	9.9e18	1.10e18	3.62e18	6.66e18	1.93e18	1.45e18	1.95e18	1.06e18	2.46e18	2.61e18	1.97e18	5.32e18
pedestrian counts	70.20	170.87	170.94	222.38	216.50	635.16	44.18	43.41	46.41	44.78	66.84	46.46	47.29	58.35	170.88
weather	2.32	2.24	2.51	2.30	2.35	2.45	8.17	2.51	2.09	2.02	2.34	2.29	2.03	2.17	2.36
nn5 daily	9.39	6.63	3.80	3.70	3.72	4.41	5.47	4.22	4.06	3.94	4.92	3.97	4.16	4.09	8.26
nn5 weekly	15.91	15.66	15.30	14.98	15.70	15.38	14.94	15.29	15.02	14.69	14.19	19.34	20.34	14.33	16.71
tourism yearly	140081.78	95579.23	90653.60	94121.08	94818.89	95033.24	82682.97	79567.22	79593.22	71471.29	70951.80	69905.47	74316.52	93572.42	99456.05
tourism quarterly	14121.09	15014.19	7656.49	9972.42	8925.52	10475.47	9092.58	10267.97	8981.04	9511.37	8640.56	9137.12	9521.67	10037.84	15845.10
tourism monthly	4724.94	5302.10	2069.96	2940.08	2004.51	2536.77	2187.28	2537.04	2022.21	1871.69	2003.02	2095.13	2146.98	2916.02	5636.83
cif 2016	715086.33	581875.97	714818.58	855578.40	642421.42	469059.49	563205.57	603551.30	1495923.44	3200418.00	679034.80	5998224.62	4057973.04	896414.41	386526.37
covid deaths	304.68	353.71	321.32	96.29	85.59	85.77	347.98	475.15	144.14	201.98	158.81	1049.48	408.66	555.24	2653.98
fred md	2013.49	2798.22	3492.84	1989.97	2041.42	2957.11	8921.94	2475.68	2339.57	4264.36	2557.80	2508.40	4666.04	1823.16	2825.67
traffic hourly	0.03	0.03	0.03	0.04	0.03	0.04	0.02	0.02	0.01	0.01	0.02	0.02	0.01	0.01	0.03
traffic weekly	1.17	1.12	1.13	1.17	1.14	1.22	1.13	1.17	1.15	1.18	1.11	1.20	1.42	1.13	1.19
saugeenday	28.63	21.50	21.49	22.26	30.69	22.38	25.24	21.28	22.98	23.51	27.92	22.17	28.06	24.53	21.50
us births	459.43	1192.20	586.93	399.00	419.73	526.33	574.93	441.70	557.87	424.93	422.00	504.40	452.87	408.49	1152.67
hospital	24.62	21.76	18.54	17.43	17.97	19.60	19.24	19.17	22.86	18.25	20.18	19.35	36.19	19.34	24.07
solar weekly	2049.09	1202.39	1210.83	908.65	1131.01	839.88	1044.98	1513.49	1050.84	721.59	1172.64	1996.89	576.35	1184.01	1729.41

Table 3. We present the mean msMAPE results for our methods along size monash baselines. We also include the naive baseline that predicts the last values in the context repeatedly.

Dataset	llmtime(ZS)	SES	Theta	TBATS	ETS	(DHR-)ARIMA	PR	CatBoost	FFNN	DeepAR	N-BEATS	WaveNet	Transformer	TimesFM(ZS)	NAIVE
australian electricity demand	15.30	22.07	22.18	13.70	44.22	28.75	8.76	8.35	9.22	12.77	7.74	8.07	8.68	17.27	22.07
bitcoin	22.70	30.31	40.36	20.12	20.65	31.11	21.48	29.78	21.00	21.16	31.95	22.16	23.32	42.00	192.46
pedestrian counts	53.02	121.39	122.08	119.76	148.48	138.58	40.29	45.54	39.30	36.10	54.15	34.22	36.16	50.73	121.35
weather	34.05	50.85	56.19	58.06	51.47	57.98	106.01	59.12	38.17	36.46	50.87	40.13	35.35	50.85	36.81
nn5 daily	47.43	35.38	21.93	21.11	21.49	25.91	30.20	24.04	23.30	23.75	28.47	22.65	23.18	23.46	48.09
nn5 weekly	12.28	12.24	11.96	11.62	12.29	11.83	11.45	11.67	11.49	11.52	10.93	14.95	14.83	11.01	13.26
tourism yearly	48.28	34.10	31.93	33.94	36.52	33.39	46.92	31.54	33.73	34.06	30.24	28.80	34.66	30.14	42.14
tourism quarterly	25.25	27.41	15.37	17.16	15.07	16.58	15.86	16.53	16.20	15.29	14.45	15.56	16.97	19.92	31.68
tourism monthly	33.53	36.39	19.89	21.20	19.02	19.73	21.11	21.10	20.11	18.35	20.42	18.92	19.74	22.77	40.40
cif 2016	18.63	14.94	13.04	12.19	12.18	11.69	12.32	14.86	12.31	13.54	11.71	18.82	12.55	20.00	15.41
covid deaths	20.65	15.35	15.57	8.71	8.64	9.26	18.34	15.40	18.58	34.18	32.36	14.50	40.71	15.89	134.79
fred md	9.75	8.72	9.72	7.97	8.40	7.98	30.77	9.16	8.99	8.33	8.32	9.08	11.26	8.47	8.76
traffic hourly	11.70	8.73	8.73	12.58	9.84	11.72	5.97	7.94	4.30	4.16	5.16	5.22	4.10	4.76	8.73
traffic weekly	12.85	12.40	12.48	12.80	12.63	13.45	12.46	12.90	12.64	13.13	12.31	13.21	15.18	12.48	12.98
saugeenday	58.53	35.99	35.97	37.34	67.50	37.55	45.32	35.55	39.32	40.22	56.03	37.02	56.62	43.31	35.99
us births	4.48	11.77	5.82	3.81	4.05	5.17	5.75	4.23	5.55	4.13	4.17	4.88	4.36	3.94	11.35
hospital	21.68	17.94	17.27	17.55	17.46	17.79	17.56	18.04	18.29	17.41	17.72	17.51	20.04	17.86	21.55
solar weekly	37.33	24.59	24.76	19.05	22.93	17.87	21.65	29.35	21.52	15.00	24.05	32.50	12.26	24.20	32.80

Table 4. MAE for ETT datasets for prediction horizons 96 and 192. Owing to expensive evaluations for lllmtime, the results are reported on the last test window of the original test split.

Dataset	llmtime(ZS)	PatchTST	FEDFormer	AutoFormer	Informer	TimesFM(ZS)
ETTh1 (Horizon=96)	0.42	0.41	0.58	0.55	0.76	0.37
ETTh1 (Horizon=192)	0.50	0.49	0.64	0.64	0.78	0.49
ETTh2 (Horizon=96)	0.33	0.28	0.67	0.65	1.94	0.28
ETTh2 (Horizon=192)	0.70	0.68	0.82	0.82	2.02	0.58
ETTm1 (Horizon=96)	0.37	0.33	0.41	0.54	0.71	0.25
ETTm1 (Horizon=192)	0.71	0.31	0.49	0.46	0.68	0.24
ETTm2 (Horizon=96)	0.29	0.23	0.36	0.29	0.48	0.28
ETTm2 (Horizon=192)	0.31	0.25	0.25	0.30	0.51	0.24

In particular, we generate these curves by (i) and (ii) linear trend + ARMA(2, 1), (iii) sum of 5 sines of different periods, (iv) and (v) a sine curve scaled linearly, (vi) a sine curve with a linear trend, (vii) a sine curve capped linearly, and (viii) minimum of two sines with a linear trend.

We notice that, compared to (Auto)ARIMA and lllmtime, TimesFM is more capable of following trends and nuanced seasonal patterns. For example, Plot 3 shows the sum of 5 sine curves which is out of the distribution of our synthetic dataset. However TimesFM still correctly identifies the seasonal pattern, likely because similar pattern occurred in the context. Plot 4 and 5 demonstrate that TimesFM correctly identifies the multiplicative trend which seems hard for either lllmtime or ARIMA.

It is worth pointing out that TimesFM is practically the fastest and the easiest forecasting method to run here. In order for

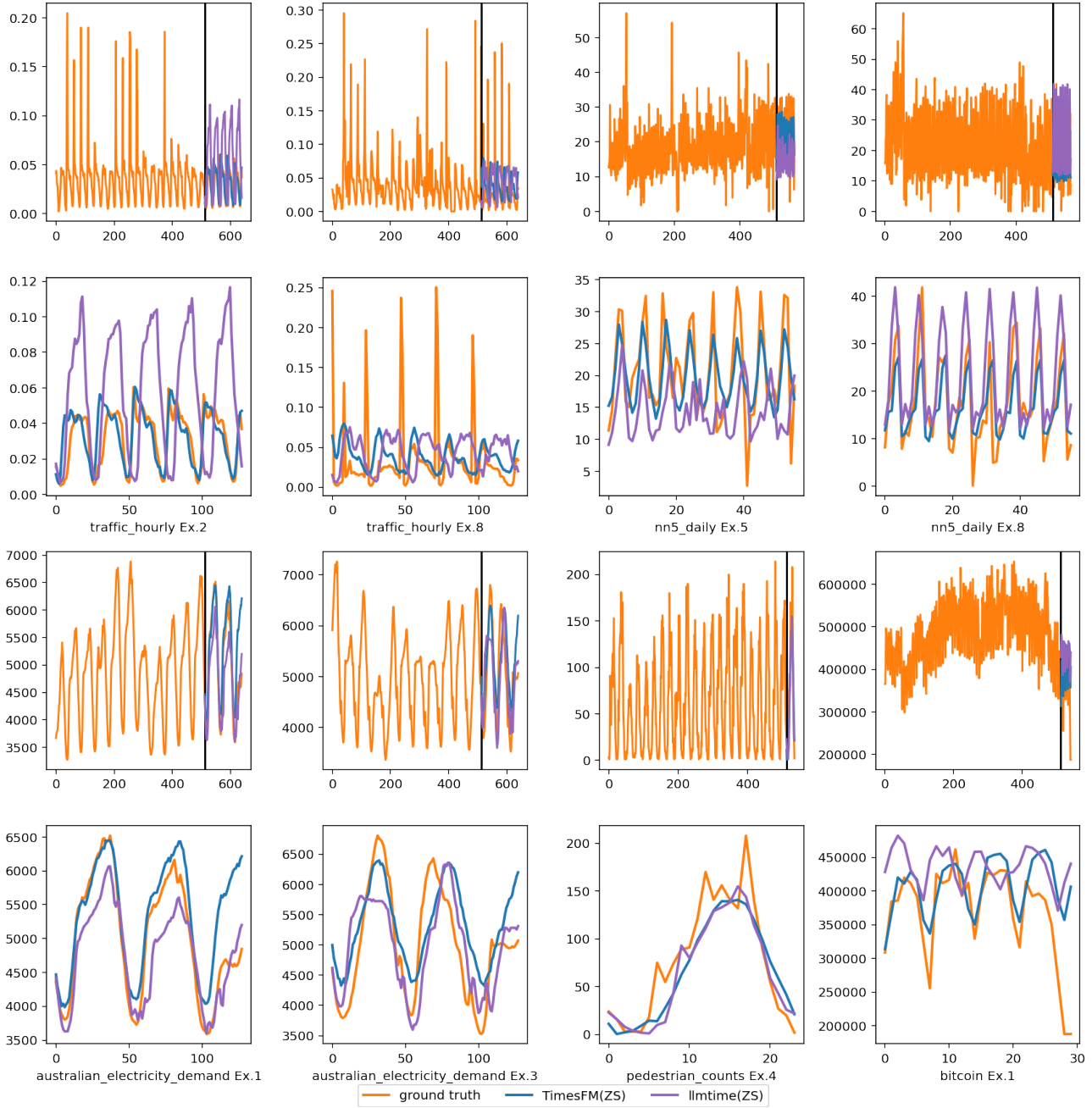


Figure 10. Forecasts visualized on a few Monash datasets. The second row plots zoom in on the prediction horizon for the sake of clarity.

(Auto)ARIMA to work the best one needs to properly identify the seasonal length and the trend. Even so the per time-series optimization takes way longer than TimesFM simply decoding.

A.3. More Details on Models

We now present implementation details about TimesFM and other baselines.

Monash Baselines. The raw metrics for the Monash baselines are directly taken from Tables 9 and 11 of the supplementary material of the original paper (Godaheva et al., 2021). For llmtime, we use the precomputed outputs provided by the authors of (Gruver et al., 2023).

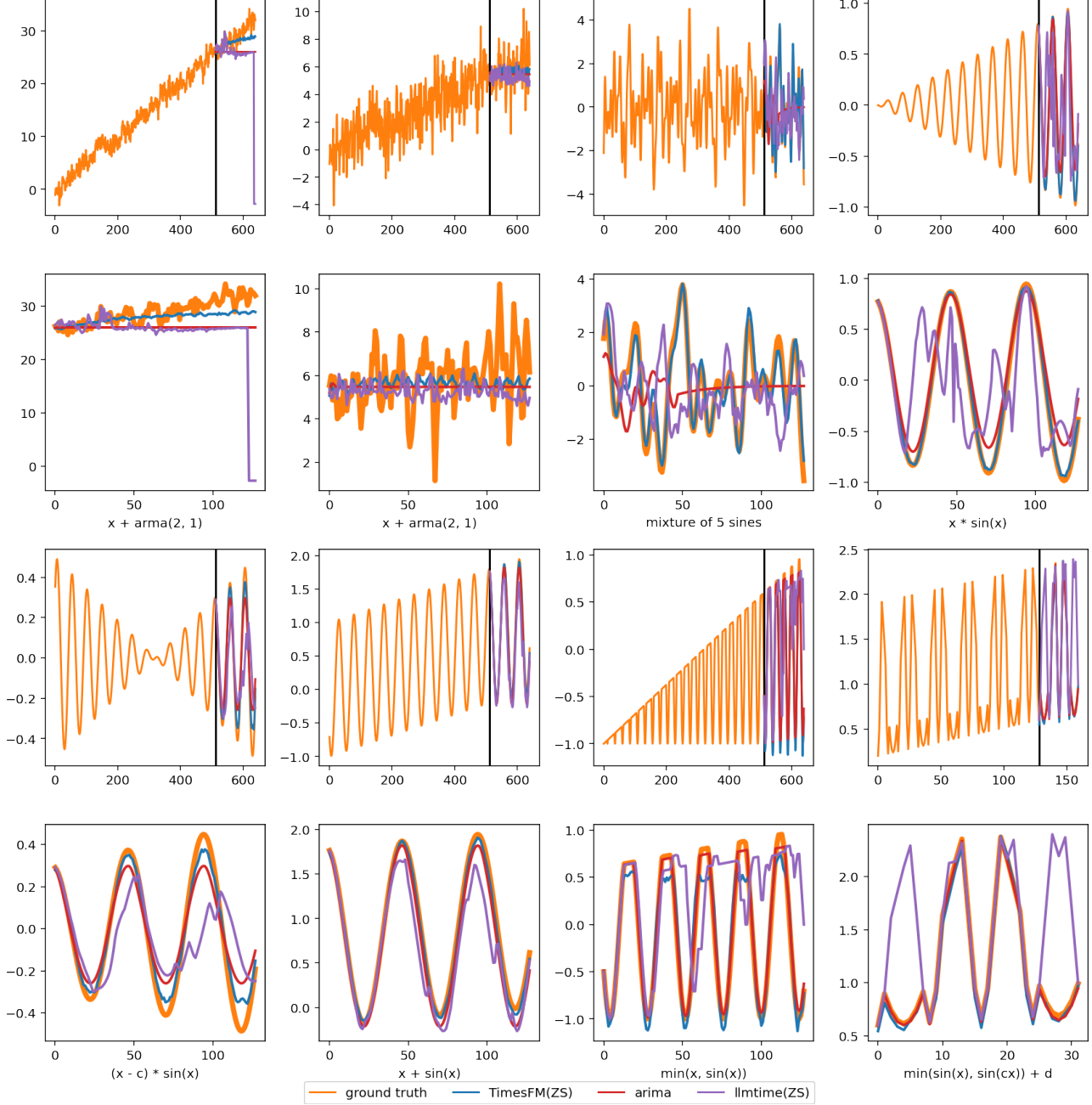


Figure 11. Forecasts visualized on synthetic curves. The second row plots zoom in on the prediction horizon for the sake of clarity.

Darts Baselines. For all the Darts baselines we use the precomputed outputs provided by the authors of (Gruver et al., 2023). For more details please see Section C.1 in that paper.

TimesFM. For our main 200M model we use 16 attention heads, 20 layers, a input patch length of 32 and output patch length of 128. The model dimension is set to 1280. We train with layer norm and a cosine decay learning rate schedule with peak learning rate of $5e-4$. The hyper-parameters of TimesFM for various sizes are provided in Table 5. Note that the settings are for the base models and not ablation models. The hidden dims of both the residual block and the FFN in the transformer layers are set as the same as model dimensions. We keep layer norm in transformer layers but not in the residual blocks.

Table 5. Hyper-parameters for TimesFM

	num_layers	model_dims	output_patch_len	input_patch_len	num_heads	dropout
Size						
200M	20	1280	128	32	16	0.2
70M	10	1024	128	32	16	0.2
17M	10	512	128	32	16	0.2

Informer Baselines. For FEDFormer (Zhou et al., 2022), Autoformer (Wu et al., 2021), Informer (Zhou et al., 2021) and PatchTST (Nie et al., 2022) we use the original hyperparameters and implementation. The results presented in the main paper are obtained on the last test window of length horizon length as stated in the llmtime (Gruver et al., 2023) paper.

We generate the llmtime predictions using the code provided by the authors⁶ but adapted to the ETT datasets. Note that as of January 2024, OpenAI has discontinued access to GPT-3, therefore we had to use the GPT-3.5-Turbo model.

A.4. Date Features

As we mentioned earlier, since we are building a single pre-trained model, we cannot have dataset specific dynamic or static covariates during training time. However, the datetime column is ubiquitous in all time-series data, so we can technically have date derived features like day of the week, month of the year etc processed into a vector at each time-point t , denoted by $\mathbf{x}_t \in \mathbb{R}^r$.

If so, the learning task can be rewritten as

$$f : (\mathbf{y}_{1:L}, \mathbf{x}_{1:L+H}) \longrightarrow \hat{\mathbf{y}}_{L+1:L+H}.$$

There are many options to incorporate these features into the model, one being to directly concatenate them after the time-points in each patch. For this paper we decide to focus on the univariate time-series input, and will investigate this enhancement in the future.

A.5. Synthetic Data

We create the synthetic data to reflect common time-series patterns using traditional statistical models. We start with four simple times series patterns:

- Piece-wise linear trends (I), where the number of the piece-wise linear components is randomly chosen between 2 and 8.
- ARMA(p, q) (II), where $1 \leq p, q \leq 8$ and the corresponding coefficients are generated from either a multivariate Gaussian or a uniform, then normalized.
- Seasonal patterns. In particular we create the sine (III) and the cosine (IV) waves of different random periods between 4 and 96 time-points and time delays.

We then randomly enable / disable these four components (I) - (IV), generate their time-series of length 1024 respectively, and sum them up using uniformly sampled random weights to create each times series in the synthetic datasets. We also choose to apply the trend multiplicatively 50% of the times the trend component is chosen.

⁶https://github.com/ngruver/llmtime/blob/main/experiments/run_monash.py