



UCL

Uni-variate Prediction Time-Series Foundation Models in Finance

by

Filip Topic

September 2024

Dissertation Supervisor: Ramin Okhrati Dissertation submitted in part in

fulfilment of the Master of Science degree in Banking and Digital Finance

Institute of Finance and Technology University College London

ABSTRACT

This dissertation explores the application of Time-series Foundation models in predicting financial time-series data and compares their performance with traditional time-series forecasting methods such as ARIMA, Prophet, and the simple naive autoregressor. Time-series Foundation Models, built on large-scale machine learning architectures, have shown potential in capturing complex temporal dependencies and patterns, making them promising candidates for financial forecasting. The study evaluates these models across multiple financial datasets across different types of data such as Stock indices, Commodity, Exchange rate and card-spending volumes. Through rigorous empirical analysis, this research provides insights into the strengths and limitations of Time-series Foundation Models, offering guidance on their practical use in financial forecasting. The results highlight the relative performance of these models in varying time-series contexts, ultimately contributing to the broader understanding of Time-series Foundation Model performance in the Finance domain. The results of this research serve as a guide to forecasters in Financial institutions as they provide insight into what type of time-series data do the Time-series Foundation Models perform well.

Keywords: Time-series Foundation Models, Finance, time-series prediction

ACKNOWLEDGEMENTS

I want to thank my University Supervisor Prof. Ramin Okhrati and Andrea Bassani from Natwest for providing me critical guidance throughout the entire research process as well as all the professors and teaching staff at UCL Institute of Finance and Technology for equipping me with knowledge and skills which were greatly needed to tackle such a complex topic.

CONTENTS

Abstract	1
Acknowledgement	2
List of Figures	7
List of Tables	8
1 Introduction	9
1.1 Research Objectives	9
1.2 Dissertation Organization	10
2 Background	11
2.1 The Transformer.....	11
2.1.1 Encoder.....	12
2.1.2 Decoder.....	14
2.1.3 Summary.....	15
2.2 Transformer-based models.....	16
2.2.1 NLP Transformer-based models.....	16
2.2.2 Computer vision Transformer-based models.....	16
2.2.3 Time-series Transformer-based models	16
2.3 Time-series Foundation models (TSFM)	17
3 Literature review	20
3.1 TimeGPT-1 [20]	20
3.1.1 Architecture.....	20
3.1.2 Training Data	21
3.2 Lag-Llama [44]	21
3.2.1 Tokenization	21
3.2.2 Architecture.....	21
3.2.3 Training data	22
3.3 Time Series Foundation Models in Finance	23
4 Methodology and Data	24

4.1	Time-series prediction evaluation	24
4.2	Data	25
4.2.1	Data Pre-processing.....	26
4.2.2	Time-series data characteristics	27
4.2.3	Data exploration	28
4.3	Experiment Design	29
4.3.1	Time Series Cross Validation (TSCV)	29
4.3.2	Time Series Foundation Models.....	31
4.3.3	Benchmark models	32
4.3.4	Experiment.....	32
4.3.5	Experiment configurations	33
4.3.6	Data-parameter experimentation phase	34
4.3.7	Aggregation phase	35
4.3.8	Fine-tuning parameter experimentation phase.....	36
4.3.9	Model-parameter aggregation phase.....	36
5	Results and Discussion	37
5.1	Results	37
5.1.1	Data-parameter experimentation phase results	37
5.1.2	Model-parameter experimentation phase results	47
5.2	Limitations	49
5.2.1	Information disparity between fine-tuned TSFMs and benchmark models + zero-shot TSFMs.....	49
5.2.2	Statistical significance	49
5.2.3	Computational resources limitations	50
5.2.4	TimeGPT-1 API calls	50
5.2.5	TimeGPT-1 pre-training data	51
5.2.6	Data availability.....	51
5.2.7	Fine-tuning disparity between Lag-Llama and TimeGPT-1.....	51
5.3	Discussion.....	52
5.3.1	Overall results discussion.....	52
5.3.2	General Lag-Llama great MDA performance.....	53
5.3.3	Inter-model fine-tuning performance discussion	53

5.3.4	Lag-Llama fine-tuning discussion.....	54
6	Conclusion and Future work	57
6.1	Conclusion	57
6.2	Future work	58
A	History of time-series forecasting	66
A.1	Brief History of time-series prediction in Finance.....	66
A.2	Brief History of modern time-series prediction methods.....	67
A.2.1	Statistical models	67
A.2.2	Machine Learning (ML) models.....	68
A.2.3	Deep Learning models.....	69
B	Background	70
B.1	The transformer.....	70
B.1.1	Token.....	70
B.1.2	LSTM and GRU Hidden state.....	70
C	Lag-Llama	71
C.1	LLaMA	71
D	Methodology	72
D.1	Time-series prediction	72
D.1.1	Prediction error.....	72
D.1.2	MDA	72
D.1.3	MES	72
D.2	Data	73
D.2.1	CHAPS data.....	73
D.2.2	Data sourcing	73
D.2.3	Partial autocorrelation	73
D.3	Experiment design	74
D.3.1	ARIMA.....	74
D.3.2	Naive Simple Autoregressor (NSA)	74
D.3.3	Loss	74

D.3.4	Benchmark	74
-------	-----------------	----

LIST OF FIGURES

2.1	Transformer architecture schema [59]	12
2.2	Scaled Dot product attention [59]	14
3.1	Lag-Llama architecture [44]	22
4.1	Breakdown of all the time/series used by trend, presence of cyclical patterns and stationarity	28
4.2	Breakdown of all time-series by type of data and frequency.....	28
4.3	Breakdown of all time-series by time-series features, type of data and frequency.	29
4.4	Schematic view of the TSCV strategy	31
4.5	Schematic view of a single fold from the TSCV technique.....	33
4.6	Schematic view of the Data-parameter experimentation phase.....	35
5.1	Predictions of Lag-Llama with Batch size 5, fine-tune length 200, Max epochs 8 and Learning rate 0.0005 on weekly WTI returns	56
5.2	WTI weekly returns between 2019-01-01 and 2024-01-01.....	56

LIST OF TABLES

4.1	Types of data used.....	26
4.2	Available frequencies of different types of data.....	26
4.3	Time periods used for different frequencies of data.....	26
4.4	Different experiment parameters and their values	34
5.1	Overall results of the exploration phase	40
5.2	Breakdown of results by different Context lengths	40
5.3	Breakdown of results by different Frequency of data	41
5.4	Breakdown of results by different Type of data	42
5.5	Breakdown of "daily" results by different Types of data	43
5.6	Breakdown of "weekly" results by different Types of data	44
5.7	Breakdown of "monthly" results by different Types of data	44
5.8	Breakdown of results by stationarity	45
5.9	Breakdown of results by trend	45
5.10	Breakdown of results by presence of cyclical patterns	46
5.11	Lag-Llama Batch size optimization	48
5.12	Lag-Llama fine-tune length optimization	48
5.13	Lag-Llama Max epochs optimization	48
5.14	Lag-Llama Learning rate optimization.....	48

1. INTRODUCTION

Time-series forecasting is a critical component of financial analysis, driving decision-making and resource allocation in a company. Accurate predictions of future values based on historical trends enable businesses and policymakers to mitigate risks, optimize strategies, and make informed decisions regarding asset purchases and capital allocations. In particular, time-series forecasting carries profound implications for portfolio management, hedging strategies, monetary policy formulation, and consumer credit risk assessments.

In 2017, the introduction of Transformers revolutionized the field of deep learning, initially making waves in natural language processing (NLP). Unlike recurrent architectures such as RNNs and LSTMs, Transformers rely on self-attention mechanisms that allow them to process entire sequences simultaneously rather than sequentially. This ability to model long-range dependencies and capture intricate patterns has sparked interest in their application beyond NLP, including time-series forecasting. Time-series Foundation Models, derived from these architectures, are now emerging as powerful tools for financial forecasting, potentially outperforming traditional models in certain contexts by effectively capturing the non-linear and dynamic nature of financial markets.

1.1 Research Objectives

This dissertation aims to evaluate the performance of Time-series Foundation models in predicting financial time-series on prominent financial data such as Stock index prices, Oil prices, Exchange rates, and bank card spending volumes. The key research objectives are:

- To assess the predictive accuracy of Time-series Foundation models in comparison with traditional models such as ARIMA, Prophet, and naive autoregressor.
- To analyze the ability of these models to capture complex patterns, seasonality, cyclicity and trend inherent in financial time-series.
- To provide insights into the implications of deploying and using these models in financial forecasting.

1.2 Dissertation Organization

This dissertation is structured into six chapters, this being the first chapter. Chapter 2 provides an in-depth exploration of the Transformer architecture, Transformer-based models, and the Transformer-led revolution of the field of AI. Chapter 3 reviews the relevant literature, with an emphasis on two prominent Time-series Foundation Models that are central to this research. These models are examined in terms of their architecture and capabilities. Chapter 4 outlines the methodology of the research. It describes the financial datasets used, the preprocessing steps, general characteristics of time-series and the experimentation process. Chapter 5 presents the results of the experiments and offers a detailed discussion on the performance of the Time-series Foundation models in comparison with traditional methods. It examines the findings in light of the research objectives and provides insights into their practical implications. Finally, Chapter 6 concludes the dissertation by summarizing the key contributions and discussing potential avenues for future work, including improvements and expansions of the research.

2. BACKGROUND

For a brief and interesting history on time-series forecasting, see Appendix A.1, A.2.1, A.2.2.

2.1 The Transformer

In 2017, Vaswani et al. [59] introduced a ground-breaking model called the "Transformer". Transformer is a model designed to solve the sequence-to-sequence mapping problems. In a sequence-to-sequence problem we have a sequence which consists of tokens (see Appendix B.1.1 for details about tokens) coming from a finite vocabulary¹ which are in a specific order, and we have an output sequence which is again a sequence of tokens in a specific order. The task of a sequence-to-sequence model is to learn the correct relationship between the input and output sequences so that when the model is presented with an unseen input sequence, it can predict what the correct output sequence is. In essence, sequence-to-sequence model's goal is to learn the "meaning" of the relationship between the input and output sequence. An example of sequence-to-sequence tasks are:

1. Language translation - where an input sequence could be a sentence in our native language and the output sequence would be the correct translation of that sentence in a foreign language.
2. Chatbots - where an input sequence could be a question and the output sequence the correct answer to that question.
3. Time-series forecasting - where an input sequence is a certain time-series and the output sequence is the future values of that time-series.

Before the Transformer, there have been many ML models which attempted to solve these tasks. Sutskever et al. (2014) [53] used multilayer LSTM² (type of RNN³ model) for this task. LSTM [25] and GRU⁴ [10] used to be state-of-the art sequence-to-sequence

¹The vocabulary doesn't always need to be finite, as we will see later on.

²LSTM stands for Long Short Term Memory cell

³RNN stands for Recurrent Neural Network. It is a type of Artificial Neural Network architecture.

⁴GRU stands for Gated Recurrent Unit. GRU is also a type of RNN.

models, however they have some key issues. They require large memory⁵, their nature is sequential⁶, and they suffer from an information bottleneck due to the fixed size of the hidden state (see Appendix B.1.2 for hidden state explanation). These problems were addressed by the Transformer.

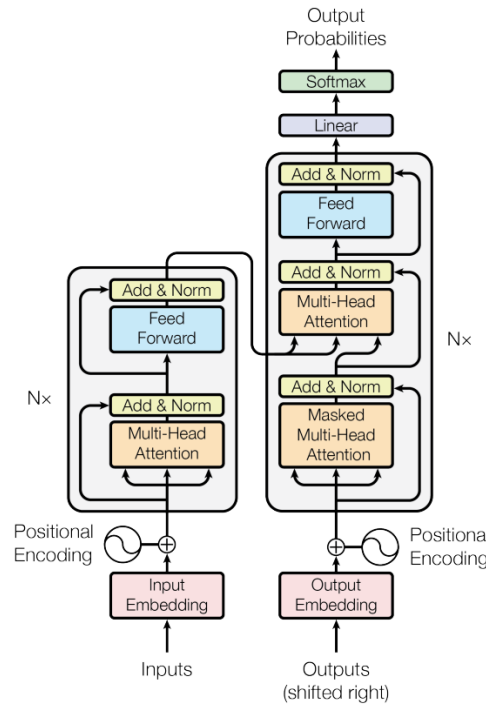


Figure 2.1: Transformer architecture schema [59]

Transformer consists of an Encoder and a Decoder. Figure 2.1 is a schematic representation of the Transformer model. Left part is the Encoder and the right part is the Decoder.

Encoder

The encoder is a stack of multiple identical layers⁷, each layer consisting of two sub-layers:

1. Multi-head self-attention layer
2. Feed-forward neural network

⁵This means that it is difficult to parallelize the training process across training examples and smaller batch sizes had to be used.

⁶This means that it is impossible to parallelize training within the training examples.

⁷In the original paper there are 6 of these layers in the encoder stack

Multi-head Self-attention layer consists of N so-called "attention heads"⁸. An attention head is essentially a series of arithmetic operations performed on the input sequence: Let's say we have a sequence of tokens $S = s_1, s_2, \dots, s_n$ the dimension of input tokens is 1 by d_{model} ⁹. An attention head consists of weight matrices¹⁰ W^Q , W^K and W^V ¹¹ of the dimension (d_{model} by d_k) where $d_k = d_{\text{model}}/N$. Multiplying these weight matrices by token s_i from the input sequence creates Q_i , K_i and V_i vectors respectively of dimension 1 by d_k . For each token i , dot product of Q_i and K_j (for all integer j , $0 < j < n+1$) is calculated: $\text{dot}_{i,1}$, $\text{dot}_{i,2}$, ..., $\text{dot}_{i,n}$. These dot-products are then scaled¹², and a softmax layer is applied¹³ to all these dot products. Each dot product $\text{dot}_{i,j}$ then multiplies the corresponding V_j to get the series of V vectors: $V_{i,1}$, $V_{i,2}$, ..., $V_{i,n}$. All these V vectors are then element-wise summed up to Z_i (of dimension 1 by d_k) which represents the self-attention output for the token i . This is repeated for all n tokens and creates a series of vectors Z_1, Z_2, \dots, Z_n (one for each token) which are then vertically stacked into matrix Z of dimension n by d_k . This matrix Z is the output of that single self-attention head. If we vectorize this process¹⁴, we can express attention as

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V$$

See Figure 2.2:

This process goes on in parallel in every head of the multi-head self-attention layer. Finally, the Z matrices of each self-attention head Z_1, Z_2, \dots, Z_N are horizontally concatenated to produce the output of the whole multi-head self-attention layer L : Z_L (of dimension n by Nd_k). Multi-head self-attention layer has a weight matrix W_O of dimension Nd_k by d_{model} which gets multiplied by Z_L to produce the final output of the multi-head self-attention layer L .

Before the output of the multi-head self-attention layer goes into a feed-forward neural

⁸In the Original paper, $N = 8$.

⁹In the original paper, $d(\text{model}) = 512$.

¹⁰These weight matrices are all learnable weights.

¹¹ Q , K , and V stand for Query, Key and Value.

¹²Scaling factor is $1 / \text{square root of } d(k)$. According to the paper, this is done so that the weights are more stable during training.

¹³Applying Softmax on a series of numbers means scaling all of them so that they sum up to 1.

¹⁴When we need to do the same arithmetic operations multiple times with different vectors, we can simply concatenate these vectors into matrices and do the operations just once with these matrices.

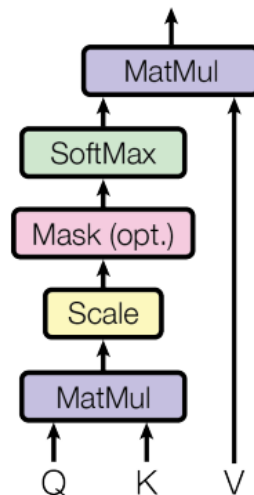


Figure 2.2: Scaled Dot product attention [59]

network¹⁵, a residual connection¹⁶ [23] and layer normalization¹⁷ [3] are applied. Residual connection and layer normalization are applied after every sub-layer of the Encoder and Decoder.

Decoder

The Decoder has a similar structure as the Encoder with a few notable differences.

Same as the Encoder, Decoder has N layers. However, each layer has three sub-layers:

1. masked self-attention layer
2. encoder-decoder attention layer
3. feed forward layer

In principle, masked self-attention layer works similarly as the multi-head attention layer in the Encoder, however, Q , K , and V come from the already generated output,

¹⁵This network has two hidden layers with ReLU activation of the first layer and linear activation of the 2nd layer. In the original paper, dimensionality of hidden layers is 2048.

¹⁶Residual connection in this context means adding the input of the layer to the output of that layer before passing to the next layer. This is used to mitigate the vanishing gradient problem and to stabilize the training process.

¹⁷Layer normalization is a scaling technique that normalizes the the data across features within each data-sample individually, as opposed to batch normalization which normalizes each feature across multiple data-samples.

and the attention is calculated only for the current and previously generated outputs i.e. for an already generated token i , masked self-attention head only calculates Z_1, Z_2, \dots, Z_i . Otherwise, the process is the same as in the Encoder multi-head self-attention sub-layer¹⁸. This makes the Decoder particularly suitable for autoregressive tasks as it predicts tokens sequentially based on which tokens have been predicted before, as opposed to the Encoder which would predict a sequence of tokens which it thinks would have the highest probability of being true overall without regard for the order of the sequence. (See Figure 2.1.)

Encoder-decoder attention layer also works similarly as the multi-head self-attention layer in the Encoder, however in this layer Q comes from the previous masked self-attention layer but K and V come from the Encoder. This allows every position in the Encoder to attend all positions in the input. (See Figure 2.1.)

Summary

Even though the Attention mechanism has been around for a while [4], the self-attention model, as implemented in the Transformer, was a revolutionary step - as it solved the issue of long range dependancies¹⁹. This concept enables the model to learn the intrinsic structure of the sequence it is presented with and therefore builds a certain level of actual "understanding" of the inputs rather than simple "fitting".

Another revolutionary aspect of the Transformer is the positional encoding. Positional encoding injects information about the absolute and relative position of each token in the sequence thus allowing the model to learn about the importance of relative positions of tokens. Positional encoding is done by summing each input token with a positional vector of equal dimension. Value of each number in the positional vector is a sin (or could be cos) function of the position of the token it is being added to as well as of the position inside the positional vector.

Finally, and probably most importantly, the way Transformer processes inputs is

¹⁸i.e. Encoder also stacks the outputs of each head, applies the W_o matrix, the residual connection and layer normalization.

¹⁹Earlier attention models had issues relating tokens which were too far apart due to attention being applied only on hidden states (hidden states are a feature of RNN models) rather than on inputs directly.

inherently parallelizable - allowing significant scaling benefits to be exploited and very large models to be built. This is beneficial as larger models can learn more intricate patterns in the data and they allow for richer vector representations of the input data²⁰. Richer representations of the input data allow for better understanding of more nuanced differences between similar tokens.

2.2 Transformer-based models

Many have built on the success of the Transformer and have come up with their own models which retain many concepts from the original Transformer.

NLP Transformer-based models

Among the most famous NLP transformer-based models are the BERT (2018) [15], T5 (2019) [45] which themselves have been built upon many times (2019 alBERT [31], 2019 roBERTa [34], 2019 distilBERT [47], 2020 mT5 [65], 2024 flan-T5 [11]), however, OpenAI's GPT²¹ series and LLaMA series [58] take the crown as the most famous Transformer-based NLP models by far.

Computer vision Transformer-based models

Another field which saw a lot of Transformer-based work is Computer Vision (CV) where among the most famous Transformer-based models are 2020 ViT [18], 2020 DETR [7], and 2021 swin-Transformer [36].

Time-series Transformer-based models

It was a matter of time Transformer-based architectures appeared in time-series prediction tasks. Li et al. (2019) [32] wrote a pioneering work on transformers (LogSparse Transformer) in time-series forecasting where they addressed the issue of the

²⁰Vector representations in the original Transformer had the dimension 512, whereas some modern transformer-based models support dimensions in the thousands.

²¹GPT stands for Generative Pre-trained Transformer

quadratic space complexity²² of the Transformer²³ and proposed a more task-appropriate version of self-attention²⁴. Zhou et al. (2021) [70] proposed a model (Informer) which also reduces computational and space complexity to $O(L \cdot \log(L))$ using "ProbSparse" self-attention mechanism²⁵ and self-attention distilling²⁶. Zhou et al. (2021) [71] incorporate a seasonal-trend decomposition approach [12] [60], along with Fourier analysis into the transformer-based model (FEDformer)²⁷. This approach saw great improvements in time-series prediction. Improvements were in terms of lower prediction errors as well as in terms of the distribution of predictions being closer to the distribution of ground-truth values according to the Kolmogorov-Smirnov distribution test²⁸ [39].

2.3 Time-series Foundation models (TSFM)

Foundation Models (FM) are a class of deep models which are pre-trained on a large amount of data - thus being able to generalize²⁹ well as they have been taught many diverse patterns. FMs saw success on the fields of CV and NLP so it is only natural that development of FMs start developing on the time-series front. And so was the case: Since 2022, over 50 models, which can be classified as Time-series Foundation Models, have been released. If we analyze the TSFM landscape, according to the taxonomy proposed by Liang et al. (2024) [33], we can see that the TSFM is a very diverse class of models. TSFMs can be classified according to:

- Type of time-series it is working with, which can be:

1. **Standard time-series** - simple sequence of any number of data-points, each associated with a time-stamp, ordered chronologically.

²²Quadratic space complexity of a model means that the amount of working memory a model uses grows quadratically with the size of the input of the model. This is especially problematic for longer inputs for which the required memory could explode.

²³They used heuristic approach to reduce the storage complexity to $O(L \cdot \log(L))$.

²⁴They proposed so-called "convolutional self-attention" which brings local context awareness to Q-K matching

²⁵ProbSparse self-attention uses Query Sparsity Measurement - a metric which helps determine which Qs are more "dominant" so that the attention can be calculated for only those Qs. This is an upgrade from LogSparse Transformer, which used a heuristic to determine Which QK pairs will be calculated.

²⁶Distilling means that the outputs of each the self-attention layers are smaller than their inputs - thus reducing the number of calculations in every subsequent layer relative to the previous one.

²⁷FED stands for Frequency Enhanced Decomposition.

²⁸Authors of the FEDformer claim even though Informer predictions were good, they don't have the same distribution as the ground-truth time-series.

²⁹Generalization is the ability of a model to perform well on data that it hasn't been trained on.

2. **Spatial time-series** - standard time-series but with a spatial dimension as well.
 3. **Trajectory** - a sequence of time-stamped locations which describe the movement of an object in space.
 4. **Event sequence** - a chronologically ordered sequence of events within a specific context.
- Model architecture, which could be:
 1. **Transformer-based**.
 2. **Non-Transformer-based** (usually MLP, RNN or CNN -bases models).
 3. **Diffusion-based** [50] [24] - self-supervised models usually aused for image generation. They are trained by adding gaussian blur to the training-sample in a markov process manner and then applying (usually) CNN-based model to learn how to un-blur the same sample..
 - The nature of the models' pre-training. which could be:
 1. **Pre-trained LM** (Large Model). These models use an LM (usually Large Language Model) which has already been trained on some other type of sequential data, for another task and they just adopt it for time-series. Adoption strategies usually involve changing the way the inputs are tokenized in order to work better with time-series data.
 2. **Self-supervised models** - which can be further split into Generative, Contrastive and Hybrid models. Generative models are those which have been trained to reconstruct the input data. They are particularly useful for tasks that involve generating a "missing" part of the input data (such as time-series forecasting). Contrastive models are those that have been trained to distinguish between "similar" and "dissimilar" pairs of data. As such, they are more appropriate for classification tasks (In time-series analysis, they are usually used for anomaly detection). As the name suggests, Hybrid models use mix of these two strategies..
 3. **Fully supervised**.

- Capabilities to adapt to a new time-series, which could be:

1. **Fine-tuning.**
2. **Zero-shot learning.**
3. **Prompt engineering.**
4. **Tokenization.**

In this dissertation, I will only consider a small subset of TSFMs; ones that belong to a class of standard time-series, Transformer-based, Self-supervised (Generative) models with fine-tuning capability. However, this class of TSFMs is itself very large [33] (in no particular order: PatchTST [42], MOIRAI [63], Lag-Llama [44], TimeSiam [17], Timer [35], TimesFM [14], UniTS [19], TimeGPT-1 [20], Chronos [1], MTSMAE [55]) so the efforts of this dissertation will be focused on two pioneering examples, trained on vast collection of time-series data spanning multiple domains [33]: Lag-Llama and TimeGPT-1.

3. LITERATURE REVIEW

3.1 TimeGPT-1 [20]

TimeGPT-1 is closed-source model - meaning that there is a lot of opacity to model's architecture, parameters and training data.

Architecture

TimeGPT-1 has encoder-decoder structure with multiple layers, each having residual connections and layer normalization. It utilizes self-attention mechanism based on the original Transformer [59].

Special capability of TimeGPT-1 is multivariate time-series forecasting - meaning that it is able to take into account "special events" and exogenous features when making predictions on a target time-series. Special events and exogenous features are time-series which are assumed to have some influence on the target time-series. Example of the former being bank holidays if the target time-series is number of flights per day and example of the latter being price of petrol if we are forecasting number of cars on the road. However, in order to do this kind of forecasting, we would have to know the actual future values of the exogenous variables which is almost impossible. TimeGPT-1 approaches this by separately forecasting the exogenous variables into the future and then basing the forecast of the target time-series on its own forecast of the exogenous features. In my opinion this is a flawed approach, inferior to simple uni-variate prediction because it brings in additional dimensions of uncertainty - whether the exogenous variables forecast are correct and whether the forecasted relationship between exogenous variables and target time-series is correct (essentially making a forecast based on a forecast). In my opinion this approach should be only used if we have high degree of certainty in forecast of future exogenous variables and the relationship between them and the target variable.

Training Data

Authors claim TimeGPT-1 has been trained on a data-set containing over 100B data-points - the largest collection of publically available time-series according to their knowledge. Data comes from the domains of finance, economics, demographics, healthcare, weather, IoT sensor data, energy, web traffic, sales, transport, and banking. Due to diversity of data domains, the training set contains time-series with many different characteristics: seasonalities, cycles, trends, noise, outliers. Having been trained on such a diverse dataset, TimeGPT-1 has very strong robustness and generalization capabilities.

3.2 Lag-Llama [44]

Tokenization

Lag-Llama constructs "lagged features" from the past values of the time-series according to a set of frequency-dependant set of lag indices $L_{\text{indices}}=1, \dots, N$. The set of lagged features of a point y_t at timestamp t is then $\{y_{t-L_{\text{indices}}[i]}, \text{ for all integers } i \text{ such that } 0 < i < N+1\}$. Lag-Llama also constructs "date-time features" F which for a point y_t in time t contain information about second-of-the-minute, minute-of-the-hour, ..., month-of-the-year. Final tokenization is done by simply concatenating the date-time features and lagged features into a single vector. One of Lag-Llama's hyperparameters, "Context Length" cl determines how many timepoints it will consider when making a prediction. When making a prediction for y_{t+1} , Lag-Llama will use $\{y_{t-cl+1}, y_{t-cl+2}, \dots, y_t\}$ in order to make a prediction. It is important to keep in mind that even though Lag-Llama only uses cl previous points for making a prediction, due to how the tokenization process works - specifically the lagged features, tokens contain information from timestamps older than cl time-points in the past.

Architecture

Lag-Llama is an open-source, decoder-only TSFM based on LLaMA [58] architecture with reduced number of parameters³⁰ (2.5m) . Similarly as LLaMA, Lag-Llama utilizes concepts

³⁰Original LLaMA series comes in sizes of 6.7B, 13.0B, 32.5B and 65.2B parameters.

such as RMSnorm [69], RoPE [52] and SwiGLU [48] activation function[41] (See Appendix C.1). On top of N decoder layers, Lag-Llama has a "parametric distribution head". Parametric distribution head is a module which predicts the distribution parameters of the next prediction, given the set distribution. Authors of the paper have chosen student-t distribution for the distribution head, meaning that at inference time, the model predicts the degrees of freedom, mean and scale [51] of the t-distribution from which it randomly draws n samples - idea being that this sample of n predictions gives a probabilistic insight into the model's prediction. Since the t distribution is symmetric, the mean of the distribution is equal to the mode which implies the mean is most likely to be the correct forecast - therefore it is used as the prediction. See Figure 3.1:

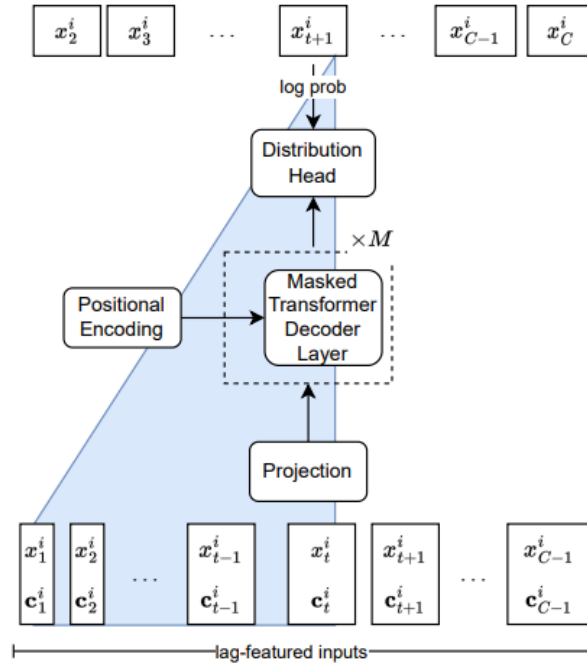


Figure 3.1: Lag-Llama architecture [44]

Training data

Lag-Llama is pre-trained on diverse set of time-series domains such as energy, transportation, economics, nature, air quality and cloud operations. Only Finance-related data that was used is exchange rate data up until 2016.

3.3 Time Series Foundation Models in Finance

There has been work done on time-series prediction using Foundation models in Finance. However, none of it was of uni-variate nature as all used some form of exogenous data and none of it was with the use of dedicated Time Series Foundation Models. Yu et al. (2023) [66] successfully used GPT-4 and LLaMA 13B LLMs with instruction-based fine-tuning, one-shot and few-shot inference on company profile, finance/economy news and index price data to predict NASDAQ-100 returns. However, they were predicting bins and movement direction of future values instead of actual values. Chen et al. (2023) [9] used ChatGPT-informed Graph Neural Network (GNN) with prompt engineering and financial news headlines to predict direction of stock movements. Xie et al. (2023) [64] used chatGPT with chain-of-thought - enhanced zero-shot prompting with twitter data to predict stock price movement. Wimmer et al. (2023) [61] used CLIP (pre-trained CV model) [43] in combination with LSTM to predict German Stock Market movements using Open, High, Low and Close prices. To my knowledge, this is the first paper strictly examining the performance of univariate TSFMs on time-series prediction in finance.

4. METHODOLOGY AND DATA

This Research aims to answer 3 main questions:

1. How good are TSFMs, namely Lag-Llama and TimeGPT-1, on financial time-series data?
2. In which cases do they perform better/worse?
3. How to improve TSFMs performance?

Methodology of this research is designed in a way to try best answer these three questions in a scientific manner.

4.1 Time-series prediction evaluation

Time-series predictions are evaluated using the traditional regression metrics due to the continuous nature of the target variables. All regression metrics are based on some variation of the prediction error (D.1.1). From the prediction error, following metrics are derived, and commonly used in time-series prediction tasks in Finance[26]: RMSE (Root Mean Square Error), MAE (Mean Absolute Error), MAPE (Mean Absolute Percentage Error) and R^2 . An additional metric that will be used in this research is MDA (Mean Directional Accuracy) (D.1.2). See [13] for Directional accuracy calculation. Finally, I suggest a new metric called MES (Mean Equal Sign) to be used in scenario where the time-series being predicted is financial returns. For a set of n time-series predictions $\{y_{\text{predicted}}^1, y_{\text{predicted}}^2, \dots, y_{\text{predicted}}^n\}$ and ground-truth values $\{y_{\text{actual}}^1, y_{\text{actual}}^2, \dots, y_{\text{actual}}^n\}$, I define MES as $\text{mean}\{\text{sign}(y_{\text{predicted}}^1 \times y_{\text{actual}}^1), \text{sign}(y_{\text{predicted}}^2 \times y_{\text{actual}}^2), \dots, \text{sign}(y_{\text{predicted}}^n \times y_{\text{actual}}^n)\}$. This metric is introduced because MDA alone doesn't fully capture the model's ability to predict direction of underlying asset value movement (see D.1.3). MDA and MES should be especially important in the field of Finance as we are not necessarily only interested in predicting the actual values of an asset, but rather we are interested in predicting the direction in which the price will go next time step.

Error-based evaluation metrics are scale-dependant as the errors themselves are inherently scale-dependant. Since the experiments were ran on different types of data of

different frequencies, each experiment is on a different scale, therefore, the raw evaluation metrics for different experiments cannot be used for model comparison across different time-series. It could be argued that MAPE accounts for this as it calculates absolute errors as a percentage of ground-truth actual values, however, due to different time-series having different time-series features such as noise, trend and seasonality, some time-series are inherently more difficult to predict than the others hence the MAPE would be on a different scale on different time-series that reason. Similar case could be made for R^2 as it scales the sum of square errors (SSE) by the total sum of squares (TSS), however, it also doesn't account for the fact that some time-series are inherently more predictable than the others. A perfect example of this is CHAPS data and Stock index data: most models score quite high R^2 on CHAPS data (some even > 0.7), whereas pretty much all score negative R^2 on Stock index data. Due to these reasons, in every experiment, models are given a rank from 1 to 7 (because there are 7 models in each experiment: 2 zero-shot TSFMs, 2 fine-tuned TSFMs and 3 benchmark models - as we will see later on) for each evaluation metric, which denotes how well each model did (relative to other models) on that specific experiment. This is a common practice in time-series results evaluation and is widely employed - most famously in the M-series competitions (prestigious time-series forecasting competition) [38] [37]. Other scale-invariant metrics such as RMSSE³¹ and MASE³² [28] were considered, however, due to the way they are calculated, could facilitate unfair comparison if time-series are non-stationary or have heteroskedastic errors. Example of the ranking system is if a model has MSE Rank 1 on experiment E, it means that the model was the best performing model in experiment E, according to the MSE metric. Since the scale of the ranking system is constant 1-7, it can be used to compare relative model performance of the models on different time-series. However, a flaw of the ranking method is that it doesn't account for exactly how much models are better/worse than each other.

4.2 Data

This project uses 4 types of financial time-series (See table 4.1). Index price data was sourced from Yahoo Finance using yfinance Python package. Commodity and Exchange

³¹Root Mean Square Scaled Error

³²Mean Absolute Scaled Error

rate data was sourced using Alpha Vantage Python API. CHAPS³³ data was sourced from UK Office for National Statistics (see Appendix D.2.1). Depending on the frequency and the type of financial data, there is different availability of it. See table 4.2 for information on which frequencies were available depending on the type of data. See table 4.3 for information on which time-periods were used depending on the frequency of choice. See D.2.2 for exceptions to these tables. All in all, 56 distinct time-series were used for this research.

Stock Index	Commodity	Exchange Rate	CHAPS
S&P 500	WTI	USD/GBP	Aggregate
FTSE 100			Delayable
DOWJ			Social
NASDAQ			Staple
			Work-related

Table 4.1: Types of data used

Frequency	Type of Data			
	Stock Index	Commodity	Exchange Rate	CHAPS
Daily	YES	YES	YES	YES
Weekly	YES	YES	YES	YES
Monthly	YES	YES	NO	NO

Table 4.2: Available frequencies of different types of data

Daily	Weekly	Monthly
2018-01-01 to 2020-01-01	2015-01-01 to 2020-01-01	1987-01-01 to 2024-01-01
2020-01-01 to 2022-01-01	2017-01-01 to 2022-01-01	
2022-01-01 to 2024-01-01	2019-01-01 to 2024-01-01	

Table 4.3: Time periods used for different frequencies of data

Data Pre-processing

In the field of Finance, there is a concept of "return". At a certain point in time T , the return R_T of an asset A_T is calculated as:

$$R_T = (A_T - A_{T-1})/A_{T-1}$$

³³CHAPS is UK debit and credit card spending index

i.e. the percentage change in the value of that asset between the points in time $T-1$ and T . If we think about each one of the fore-mentioned time-series as prices of an underlying asset, we can represent them as time-series of returns rather than their actual values which is a common practice in financial time-series prediction [26]. This method is not applied to CHAPS data as it wouldn't make sense given the nature of that data. Besides practical applications or working with asset returns rather than with raw values, this procedure is done to make the time-series stationary - which can be beneficial to models' performance.

Time-series data characteristics

Time-series data exhibits several key characteristics that are essential to understand for effective analysis and forecasting. **Trend** is a long-term increase or decrease in the data. It could take many shapes: none, linear, exponential, polynomial. This research uses Kendall's Tau coefficient of correlation [30]. Kendall's Tau measures the strength of the relationship between two ordinal variables. If a time-series of length L has statistically significant Kendall Tau correlation ($p\text{-value} < 0.05$) with a monotonically increasing sequence $\{1, 2, \dots, L\}$, this is a sign that there is a linear trend present in the data. If the square roots of every value in a time-series have statistically significant Kendall Tau correlation with a monotonically increasing sequence $\{1, 2, \dots, L\}$, this is a sign that there is a quadratic trend present in the data. Same logic applies for $\log()$ values of the time-series and exponential trend. However, since majority of our time-series contain negative values, quadratic and exponential trends weren't tested for. **Seasonality** is the presence of recurring patterns at regular intervals in the data. Chosen method for determining the presence of seasonality is by looking at partial autocorrelation (PAC) values (see D.2.3 for explanation) in the data and declaring a seasonal pattern present if there is a significant positive or negative partial autocorrelation present. Determining significance in this context is rather objective so I chose presence of PAC values >0.3 or <-0.3 to be significant. **Stationarity**: A time series is said to be stationary if the statistical properties such as mean and variance remain constant over time. Stationarity is determined using the Dickey-Fuller test [16]. Other characteristics of time-series are **Noise** and **Volatility** however, they are relatively difficult to classify and are scale dependent.

Data exploration

56 distinct time-series were used for this research. 38% of all time-series contain linear trend, 62% contain no trend. 46% of time-series contain cyclical patterns and 54% contain no discernable cyclical pattern. 88% of time-series are stationary whereas only 12% are non-stationary (see Figure 4.1).

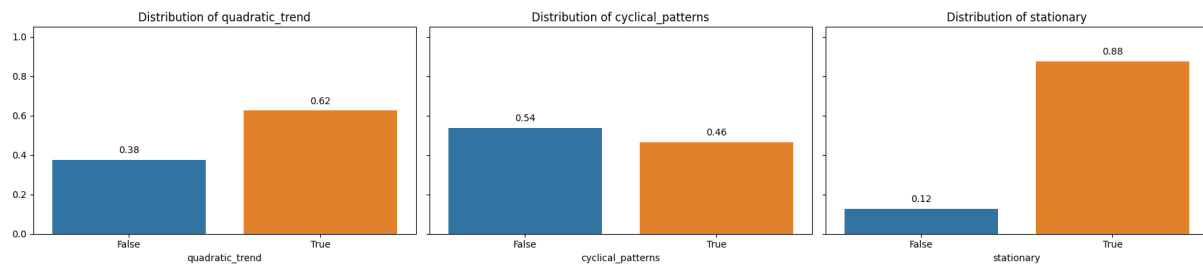


Figure 4.1: Breakdown of all the time/series used by trend, presence of cyclical patterns and stationarity

Half of all time-series used are Stock indices, 27% are CHAPS, 12% Commodity and 11% Exchange rate. 50% of time-series used had daily frequency, 41% had weekly frequency and 0.09% were monthly See Figure 4.2).

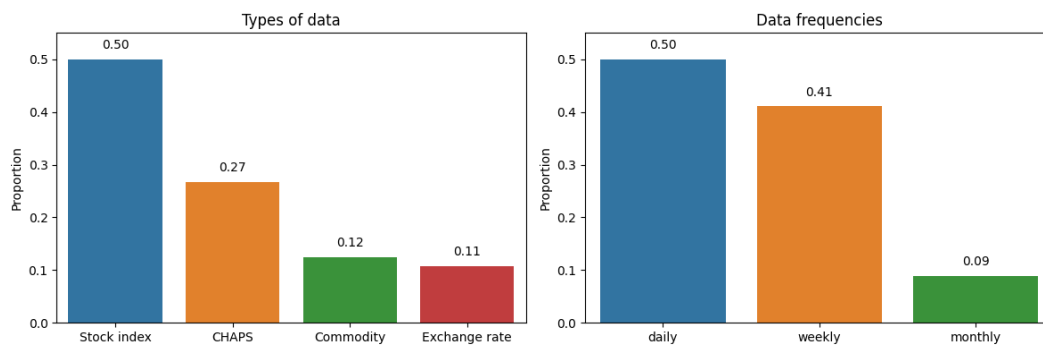


Figure 4.2: Breakdown of all time-series by type of data and frequency.

See Figure 4.3 for breakdown of time-series features by Type of data and Frequency of data. Expected observations are that Commodity, Exchange rate and Stock index (apart from one) data are all stationary and contain no linear trend, and CHAPS data is the only type which contains some non-stationary time-series and some linear trends. This is due to the pre-processing method applied on the non-CHAPS data.

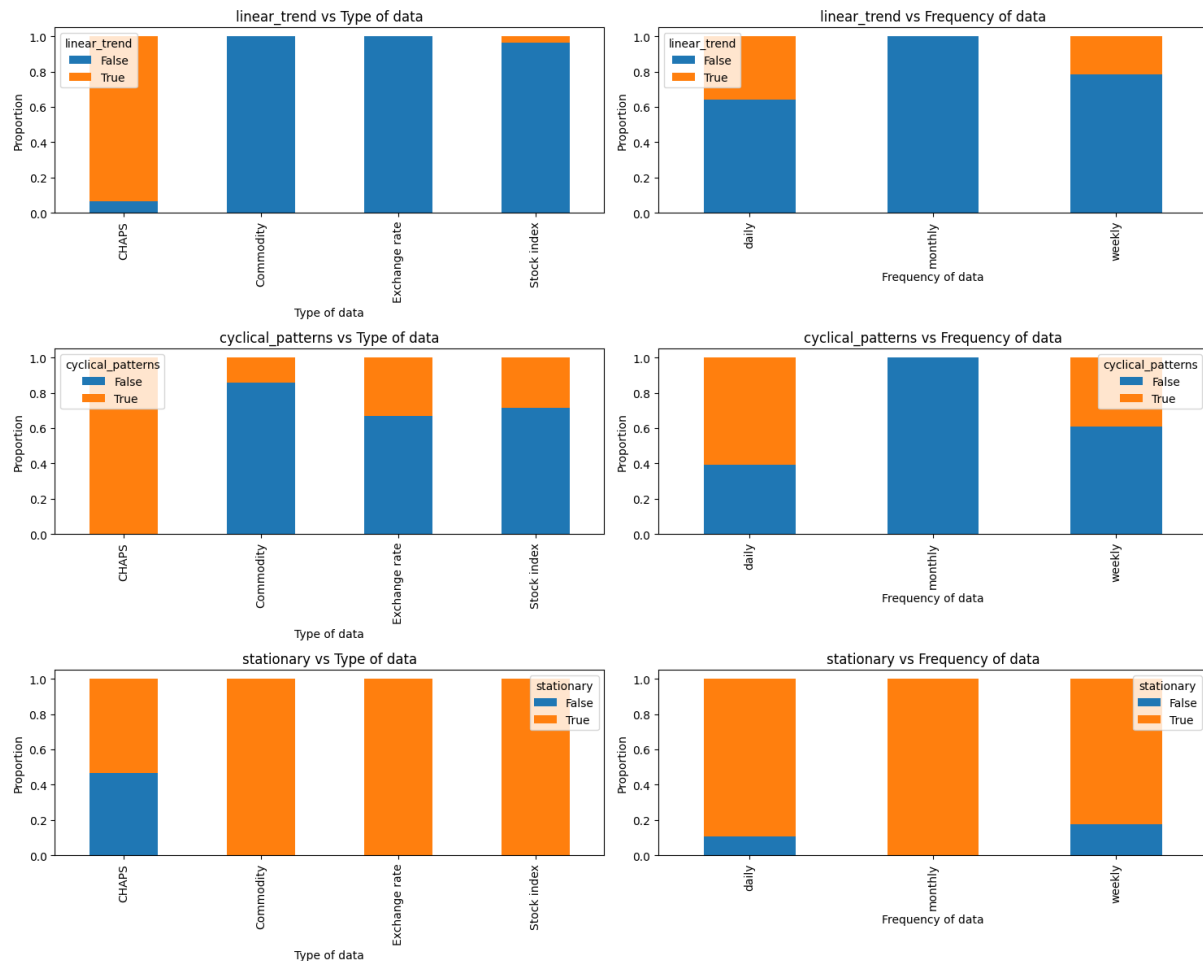


Figure 4.3: Breakdown of all time-series by time-series features, type of data and frequency.

4.3 Experiment Design

Time Series Cross Validation (TSCV)

In a traditional Machine Learning sense, K-fold cross-validation (CV) is a statistical technique used to evaluate the performance of a model by dividing the dataset into K equally sized subsets, or "folds." The model is trained on K-1 "train" folds and tested on the remaining "test" fold, and this process repeats K times, with each fold serving as the

test set once. The results are then averaged to provide a more reliable estimate of model performance. It is important because it helps prevent overfitting by providing multiple evaluations of the model's ability to generalize to unseen data. However, in time-series prediction, where data points are temporally dependent, standard K-fold cross-validation cannot be applied directly since randomly partitioning the data breaks the time order and dependencies. To address this, K-fold cross-validation for time-series is modified by partitioning the data in a way that respects temporal structure. This allows for a valid evaluation of time-series models while still leveraging the benefits of cross-validation.

In this research, TSCV works the following way: say we have a time-series S of length l : $S = \{y_1, y_2, \dots, y_l\}$, we specified the length of the train sequence to be n , and the length of prediction horizon to be h . K-fold TSCV method will sample K time-series $\{F_1, F_2, \dots, F_K\}$ from S where each fold F_i ($1 \leq i \leq K$, for all integer i) consists of the "train" and "test" part. Fold F_i^j which starts from timepoint j is defined as $F_i^j = \{y_j, y_{j+1}, \dots, y_{j+n-1}, y_{j+n}, \dots, y_{j+n+h-1}\}$. The model is trained on $\{y_j, y_{j+1}, \dots, y_{j+n-1}\}$ and tested on $\{y_{j+n}, \dots, y_{j+n+h-1}\}$ (Size constrains being $K \leq (l-n)/h$). This research only considers next-point prediction i.e. $h=1$ due to the nature of financial markets where we are primarily interested in the next value in the time-series instead of multiple future values. This research uses the version of TSCV where l is fixed which is called "rolling window" TSCV. We employ the following fold sampling strategy: say we have two additional cross-validation parameters: f and r . Folds we sample from S are split into r groups: $\{G_1, G_2, \dots, G_r\}$ such that all folds within a group are sequential i.e. the start-point of the following fold is the time-point right after the start-point of the preceding fold, each group containing f folds. The first group of folds G_1 is always $F_1^1, F_2^2, \dots, F_f^f$ (i.e. first f folds are folds with start-points 1, 2, ..., f respectively). This sampling procedure is repeated r times in a way that the time-distance between the first fold of G_k and the first fold of G_{k+1} is equal to the time-distance between G_{k+1} and G_{k+2} for all $1 \leq k \leq r-2$, and so that the distance between two consecutive groups is maximised subject to constraint $rf \leq (l-n)/h$ (this is essentially the same constraint as $K \leq (l-n)/h$). r and f parameter values used in the experiments are $r=6$ (An exception to this is with the experiments ran on the monthly data where $r=8$) and $f=5$. In simple words: we are sampling five consecutive folds (a group), 6 times in a way that the groups are uniformly spread out over the whole time-series. See Figure 4.4:

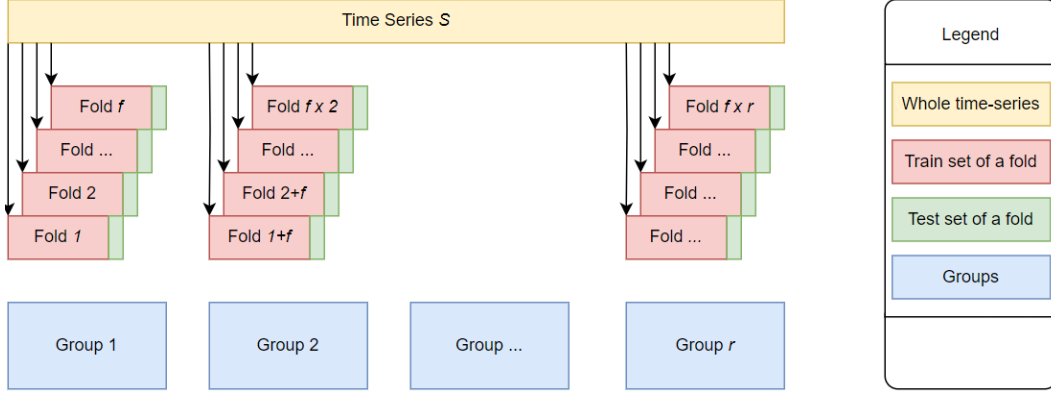


Figure 4.4: Schematic view of the TSCV strategy

Time Series Foundation Models

TSFMs used in this research are Lag-Llama [44] and TimeGPT-1 [20]. Both models have fine-tuning capabilities with accompanying fine-tuning hyper-parameters. Lag-Llama enables the user to choose Batch size BS , Max epochs ME , and Learning rate LR . Batch size refers to the number of training examples processed in one iteration before the model updates its' parameters. A smaller batch size allows for more frequent updates, while a larger batch size processes more data per update. An Epoch is a single pass through the entire fine-tuning dataset during model fine-tuning. Before a fine-tuning epoch i : E_i ($1 < i \leq ME$), model makes a record of all its' parameters. After E_i , model's parameters have changed and the loss (see Appendix D.3.3) for that epoch L_i is recorded. Then after the next epoch E_{i+1} , if $L_{i+1} < L_i$, then the model will keep the updated parameters from E_{i+1} , but if $L_{i+1} > L_i$, then the model will revert to the weights it had before E_{i+1} . On E_1 , model parameters are updated no matter what as there is no previous loss as a reference. After E_1 this procedure repeats for all $i \leq ME$. A higher ME means the model will have more opportunities to learn patterns from the data, but it can also risk overfitting. Learning rate is a hyperparameter that controls how much the model adjusts its weights with respect to the loss gradient; a small learning rate can lead to slow learning, while a large learning rate can cause unstable training (overshooting). These parameters together affect how efficiently and accurately a model learns from data.

TimeGPT-1 enables users only to choose number of "Fine-tune Steps" which corresponds to the number of epochs. TimeGPT uses the Adam optimizer with unknown learning hyperparameters, and batch size. Because we don't have sufficient information

on TimeGPT-1 fine-tuning process and fine-tuning hyperparameters, we cannot assume the same Max epochs would be optimal for Lag-Llama and TimeGPT-1. According to its' authors [20], TimeGPT-1 performance strictly improves with the number of Fine-tune Steps and reaches the plateau at around 100 Fine-tune Steps, whereas Lag-Llama reaches its' full potential with only 4 Max Epochs - as we will see later in this research. This research uses both zero-shot and fine-tuned versions of both Lag-Llama and TimeGPT-1 in order to gain insight into whether the models' performance changes with fine-tuning.

Benchmark models

Calculating the evaluation metrics for time-series prediction is not sufficient to answer whether a time-series prediction model is good or not as some time-series are inherently more and some less predictable. Therefore, the results of the TSFM evaluation need to be put into context by comparing them against widely-used benchmark models in order to make a judgment on how good they are (relative to the benchmark). Chosen benchmark models are: ARIMA (with automatic hyper-parameter selection) [27] (Appendix D.3.1), Naive Simple Autoregressor (Appendix D.3.2) and Meta's Prophet model [56].

Experiment

"Running an experiment" on a certain time-series S involves the following: We employ the TSCV technique, in a manner described before, on S . Each fold F_i ($1 \leq i \leq rf$) represents a slice of S with the train-set length n and prediction horizon $h=1$: $F_i = \{y_1, y_2, \dots, y_n, y_{n+1}\}$. Train set T_i of the Fold F_i is then $TR_i = \{y_1, y_2, \dots, y_n\}$, and the test set is $TE_i = \{y_{n+1}\}$. Firstly, the whole train set TR_i is used for TSFM fine-tuning. Then, only the last cl (Context length, see Section 3.2.1) points of T_i : $\{y_{n-cl+1}, y_{n-cl+2}, \dots, y_n\}$ are used as the context window for fine-tuned TSFMs, zero-shot TSFMs and as the actual train-set for benchmark models (see Section 5.2.1) (see Figure 4.5 for schematic view of a single TSCV fold). After all the models have made their predictions for the fold F_i , their predictions are recorded along with the ground-truth actual value TE_i . This is done for every fold of the TSCV process. Nuance worth mentioning is that fine-tuned Lag-Llama is getting fine-tuned only on the first fold of every group of folds (i.e. every 5 folds) whereas

fine-tuned TimeGPT-1 is fine-tuned on every single fold (see Section 5.2.7). After the models' predictions have been made and recorded for all folds, evaluation metrics for each model's predictions are calculated (see Section 4.1) and also recorded.

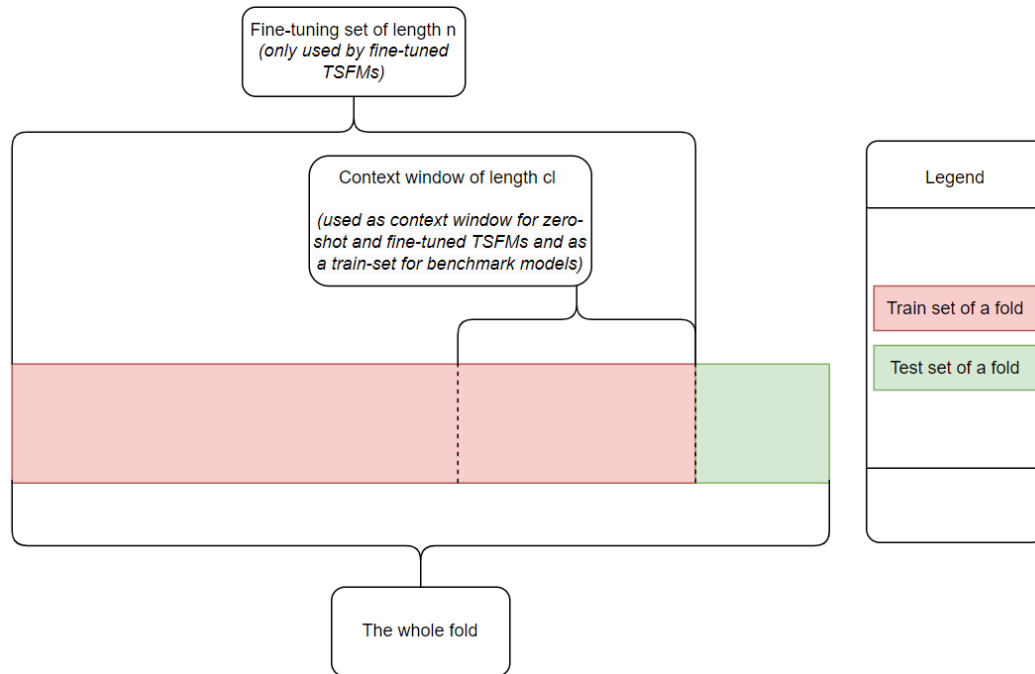


Figure 4.5: Schematic view of a single fold from the TSCV technique

Experiment configurations

This entire research consists of running many experiments, as described above, in different experiment configurations. Experiment-configuration refers to a unique combination of experiment-parameters. Experiment-parameter is a feature of an experiment which has a certain impact on how one or more models perform. Experiment-parameters can be divided into two groups: data-parameters, and fine-tuning parameters. Data-parameters are *Type of Data*, *Frequency of Data* and *Context Length*. Fine-tuning parameters are the following: *Fine-tuning length* denotes the length of the fine-tuning set used for TSFM fine-tuning (length of the red box in Figure 4.5) and is applicable to both Lag-Llama and TimeGPT-1. *Max Epochs*, *Batch Size* and *Learning Rate* are only applicable to Lag-Llama and *Fine-tune Steps* is only applicable to TimeGPT-1. See Table 4.4 for details on experiment-parameter values used. Running the experiments across many different experiment configurations allows us to analyze the results of all the experiments according to different experiment-parameters. Breaking down the results by data-parameters would

answer one of the main question of this research: "In which cases do TSFMs perform better/worse?". Breaking down the results by fine-tuning parameters would answer another main question of this research: "How to improve TSFMs' performance?". Breaking down the results by two (or more) data-parameters would give more granular insight into how models perform in specific cases which is very relevant because it is unreasonable to assume model would perform the same on time-series of same frequency but different type or time-series of same type but different frequencies. In order to get full granular insight into all joint effects of experiment-parameters, experiments need to be run across every single experiment configuration. Since we have one experiment-parameter with 4 different values and seven experiment-parameters with 3 different values, this gives the total of $4 \times 3^6 = 2916$ different experiment-configurations³⁴ to run - which is computationally unfeasible given that we run all the models, apart from TimeGPT-1, locally.

Type of experiment-parameter	Experiment-parameters	Parameter values			
Data-parameters	Type of data	Stock index	Commodity	Exchange rate	CHAPS
	Frequency of data	Daily	Weekly	Monthly	
	Context length	32	64	128	
Fine-tuning hyper-parameters	Fine-tuning length	200	128	64	
	Batch size	5	10	20	
	Max epochs	4	8	16	
	Learning rate	0.0005	0.005	0.00005	
	Fine-tune steps	100			

Table 4.4: Different experiment parameters and their values

Data-parameter experimentation phase

Since we are computationally restricted, we divide the whole experimentation process into phases, the first phase being the Data-parameter experimentation phase. In this phase, we employ a heuristic approach to experimentation where we keep all the fine-tuning hyper-parameters (last 5 rows of Table 4.4) fixed (respective fine-tuning hyper-parameter values being the ones in the first column of the last four rows of Table 4.4) and run experiments with different combinations of data-parameters. See Figure 4.6:

³⁴This is not even accounting for the fact that Stock index is represented by 4 different time-series and that CHAPS is represented by 5 different time-series.

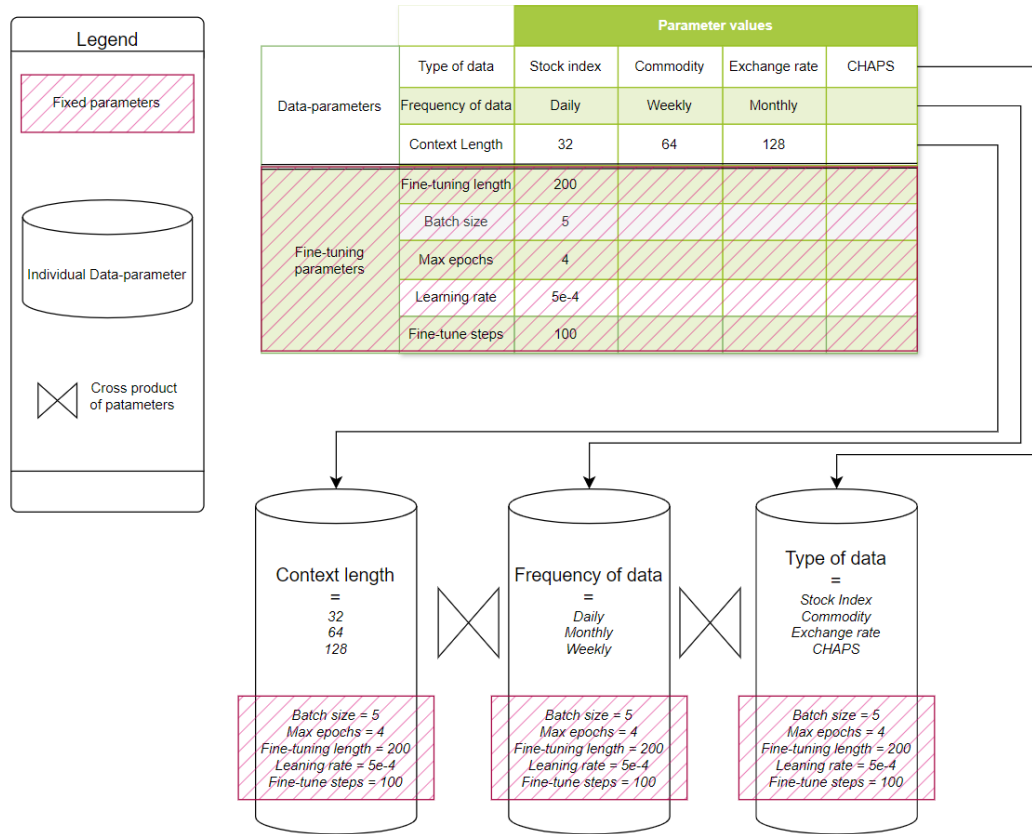


Figure 4.6: Schematic view of the Data-parameter experimentation phase

Aggregation phase

After having recorded and aggregated the results (model evaluations) and experiment-configurations in the Data-parameter experimentation phase, we group the aggregated results by any single, or combination of data-parameters and take the mean of the evaluation metric ranks (according to Section 4.1) in order to understand how models perform relative to each other, with different data-parameters. This phase constitutes the main body of the whole research.

One drawback of this whole research is that it was conducted on a relatively small collection of time-series of few different types (see limitations in Sections 5.2.6 and 5.2.3) hence with the analysis described so far, it is difficult to extrapolate any conclusions for types of data or frequencies which we haven't used in this experiment. A potential solution to this problem is to look at underlying characteristics of the time-series that were used in this experiment (as described in Section 4.2.2) because these characteristics are the factors which fundamentally influence the performance of time-series prediction

models. Time-series used in this research were classified according to presence of a trend, stationarity and presence of cyclical pattern. Breaking down the results of all experiments by these features can provide an insight into which underlying characteristics of time-series cause the models to perform well or poor. By incorporating this type of analysis into this research, a future user of TSFMs can simply perform tests to determine the mentioned characteristics of a time-series they want to predict (whether or not that time-series was used in this research or not), and this by consulting this research, they will know whether they should use TSFMs or not.

Fine-tuning parameter experimentation phase

Last phase in the heuristic experimentation approach is to pick (fix) a single data-configuration and run different experiment configurations with fixed data-configuration and variable model-configuration. The strategy for choosing the data-configuration for each TSFM's fine-tuning hyper-parameter optimization is by taking the data-configuration on which the TSFM performed really well on, but with some room for improvement in order to be able to tell whether the TSFM can improve with different fine-tuning hyperparameters. After we have chosen the data-configuration for Lag-Llama, we keep it fixed and run experiments across different Model-parameter configurations. Conceptually, this process is identical to the process in the Data-parameter experimentation phase, but this time the data-parameters are fixed, and we vary the fine-tuning hyper-parameters in each experiment. Important to note is that in this phase we only experiment with Lag-Llama fine-tuning parameters (see Section 5.2.4)

Model-parameter aggregation phase

Similarly as in the Aggregation phase, we can group the aggregated results of the Fine-tuning parameter experimentation phase by fine-tuning parameters and take the mean of the evaluation metric ranks in order to understand Lag-Llama performance changes with different fine-tuning parameters on the data-configuration that we chose. Doing this type of analysis enables us to answer the question "How to improve TSFMs' performance?" i.e. gives insight into which fine-tuning parameters it works best.

5. RESULTS AND DISCUSSION

5.1 Results

As outlined in Section 4.1, results of the Data-parameter experimentation phase (Section 4.3.6) are evaluated using the metric ranking system and aggregated as outlined in Section 4.3.7. These results can be found in this section in Figures 5.1 to 5.10. Results are broken down by different Data-parameters and time-series characteristics. For each combination of Data-parameters / time-series characteristics, the best model's performance has been highlighted in purple according to each individual evaluation metric average rank. In some results tables, highlight is missing from the "mes" column as few different models have the same performance in that group and highlighting it would bring confusion.

Data-parameter exploration phase results sub-section shows the main body of this project's research. However it is important to stress that the results from Model-parameter experimentation phase are not included in this section. The reason for that is because the specific time-series chosen for Model-parameter experimentation phase is specifically chosen because of Lag-Llama's great performance on that specific time-series. Hence, adding the results from this set of experiments to experiments from Data-parameter experimentation phase would bring bias to the whole corpus of results. Results from Model-parameter experimentation phase are separated from the results of Data-parameter experimentation phase and are addressed separately in Section 5.1.2.

Data-parameter experimentation phase results

Overall, across all experiments in this phase (see Table 5.1), ARIMA model performs the best according to all metrics apart from MDA - according to which, fine-tuned Lag-Llama is the best. This pattern repeats if we break down the results by the Context length (see Table 5.2).

Across all **Context lengths** (128, 32 and 64), ARIMA performs best according to all metrics apart from MDA, where fine-tuned Lag-Llama performs best when Context length is 128 (mean rank 2.79) or 32 (mean rank 2.95) and zero-shot Lag-Llama is the best model according to MDA with Context length of 64 (mean rank 2.84).

Breaking down the results by **Frequency** of data (see Table 5.3), ARIMA performs best on "daily" data according to all metrics (mean rank between 1.11 and 2.44) apart from MDA, but fine-tuned Lag-Llama comes on top with "monthly" (mean rank between 1.33 and 2.0) and "weekly" (mean rank between 2.94 and 3.01) data - leading in R^2 , MSE, MAE, RMSE in both of those scenarios with zero-shot Lag-Llama being the best according to MDA (mean rank 2.52 on weekly and 2.2 on monthly data) and ARIMA being the best in MES.

Breaking down the results by **Type of data** (see Table 5.4), Fine-tuned Lag-Llama tends to be the best model on "Commodity", "Exchange rate" and "Stock index" types of data according to R^2 , MSE, MAE, RMSE (mean ranks between 1.67 and 2.78). However ARIMA comes equal or better on "Stock index" in terms of MAE, MAPE and MES (mean ranks between 1.15 and 2.49), and on "Exchange rate" in terms of MAPE (2.11) and MES (1.33).

Breaking down the results of "**daily**" frequency by **Type of data** (see Table 5.5), simple naive autoregressor tends to be the best model on CHAPS (mean ranks between 1.43 and 1.90), fine-tuned Lag-Llama tends to be the best on "Commodity" (mean ranks between 2.11 and 2.89) and ARIMA tends to be best on Stock index (mean ranks between 1.17 and 2.14). On exchange rates, it isn't clear whether ARIMA or Prophet are the best.

Breaking down the results of "**weekly**" frequency by **Type of data** (see Table 5.6), fine-tuned Lag-Llama tends to dominate other models on "commodity" (mean ranks between 1.0 and 2.0), "Exchange rate" (mean ranks between 1.22 and 1.67) and "Stock index" (mean ranks between 1.83 and 1.94), with ARIMA consistently proving itself the best on MES. On CHAPS data, it is unclear whether fine-tuned TimeGPT-1 or the autoregressor is the best as both score best in different metrics.

Breaking down the results of "**monthly**" frequency by **Type of data** (see Table 5.7), it is difficult to say which is the best model on "commodity" type of data as few different models perform best according to different metrics. However, fine-tuned Lag-Llama, significantly outperforms other models on Stock index type of data (mean ranks between 1.17 and 2.17), with zero-shot Lag-Llama being tied first on MDA and ARIMA beating Lag-Llama only on MES.

Breaking down the results by the **stationarity** of the time-series (see Table 5.8), simple naive autoregressor performs best on non-stationary data (mean ranks between 1.76 and 2.95) and ARIMA performs best on stationary data (mean ranks between 1.27 and 2.65). Fine-tuned Lag-Llama's is the best model by MDA (2.78) on stationary data. An interesting observation is that Lag-Llama's relative performance is significantly better on stationary data, while TimeGPT-1 relative performance is significantly better on non-stationary data.

Breaking down the results by **presence of a linear trend** (see Table 5.9), autoregressor again dominates the time-series on data with linear trend (mean ranks between 1.87 and 2.71). However, fine-tuned Lag-Llama is the best model on time-series with no trend present (mean ranks between 2.28 and 2.35) - exception being ARIMA which beats fine-tuned Lag-Llama in MAPE (2.25) and MES (1.32).

Breaking down the results by **presence of a cyclical pattern** (see Table 5.10), fine-tuned Lag-Llama is the best model when there isn't a cyclical pattern present (mean ranks between 2.26 and 2.42) (with ARIMA beating in only in MAPE (1.97) and MES (1.36)). On time-series with a cyclical pattern, ARIMA is the best model (mean ranks between 1.09 and 2.95) apart from fine-tuned TimeGPT-1 which beats it in MDA (3.26).

Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
arima	2.71	2.71	2.68	2.71	2.33	5.09	1.23
autoregressor	5.40	5.40	5.32	5.40	5.06	3.94	3.79
ft_lag_llama	3.38	3.38	3.45	3.38	3.50	2.92	1.92
ft_timeGPT	3.57	3.57	3.64	3.57	3.65	3.58	3.93
lag_llama	4.52	4.52	4.58	4.52	4.67	3.15	2.93
prophet	4.48	4.48	4.52	4.48	4.83	3.32	4.44
timeGPT	3.94	3.94	3.82	3.94	3.96	4.25	3.52

Table 5.1: Overall results of the exploration phase

Context length	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
128	arima	2.68	2.68	2.64	2.68	2.18	4.88	1.50
	autoregressor	5.46	5.46	5.43	5.46	5.23	4.11	3.73
	ft_lag_llama	3.30	3.30	3.23	3.30	3.52	2.79	1.80
	ft_timeGPT	3.71	3.71	3.82	3.71	3.77	3.71	3.91
	lag_llama	4.34	4.34	4.39	4.34	4.50	2.91	2.91
	prophet	4.43	4.43	4.57	4.43	4.82	3.54	3.82
	timeGPT	4.07	4.07	3.91	4.07	3.98	4.39	3.43
32	arima	2.91	2.91	2.91	2.91	2.55	5.46	1.04
	autoregressor	5.38	5.38	5.20	5.38	4.95	3.59	3.79
	ft_lag_llama	3.21	3.21	3.34	3.21	3.34	2.95	2.00
	ft_timeGPT	3.23	3.23	3.32	3.23	3.48	3.38	3.95
	lag_llama	5.05	5.05	5.09	5.05	5.02	3.70	3.09
	prophet	4.52	4.52	4.54	4.52	4.73	3.05	4.80
	timeGPT	3.70	3.70	3.61	3.70	3.93	4.00	3.57
64	arima	2.55	2.55	2.48	2.55	2.25	4.93	1.16
	autoregressor	5.38	5.38	5.32	5.38	5.00	4.13	3.86
	ft_lag_llama	3.61	3.61	3.79	3.61	3.64	3.04	1.96
	ft_timeGPT	3.75	3.75	3.77	3.75	3.71	3.64	3.95
	lag_llama	4.18	4.18	4.25	4.18	4.48	2.84	2.80
	prophet	4.48	4.48	4.45	4.48	4.93	3.38	4.70
	timeGPT	4.05	4.05	3.95	4.05	3.98	4.36	3.55

Table 5.2: Breakdown of results by different Context lengths

Frequency of data	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
daily	arima	2.44	2.44	2.39	2.44	2.32	4.68	1.11
	autoregressor	5.15	5.15	4.98	5.15	4.65	3.83	3.63
	ft_lag_llama	4.04	4.04	4.13	4.04	4.42	3.25	1.71
	ft_timeGPT	3.62	3.62	3.42	3.62	3.62	3.65	3.32
	lag_llama	4.02	4.02	4.25	4.02	3.90	3.83	3.31
	prophet	5.15	5.15	5.27	5.15	5.51	2.93	4.67
	timeGPT	3.57	3.57	3.56	3.57	3.57	4.15	3.29
monthly	arima	2.20	2.20	2.40	2.20	2.07	5.53	1.33
	autoregressor	6.33	6.33	6.33	6.33	6.27	4.80	4.33
	ft_lag_llama	1.33	1.33	2.00	1.33	1.13	2.73	2.07
	ft_timeGPT	3.73	3.73	3.53	3.73	3.27	4.40	5.13
	lag_llama	3.13	3.13	2.47	3.13	4.73	2.20	2.20
	prophet	5.93	5.93	5.93	5.93	5.73	2.60	5.93
	timeGPT	5.33	5.33	5.33	5.33	4.80	4.27	4.60
weekly	arima	3.16	3.16	3.09	3.16	2.39	5.49	1.36
	autoregressor	5.51	5.51	5.51	5.51	5.29	3.88	3.87
	ft_lag_llama	3.01	3.01	2.94	3.01	2.90	2.57	2.14
	ft_timeGPT	3.46	3.46	3.93	3.46	3.78	3.30	4.42
	lag_llama	5.43	5.43	5.43	5.43	5.58	2.52	2.64
	prophet	3.33	3.33	3.29	3.33	3.80	3.96	3.84
	timeGPT	4.09	4.09	3.81	4.09	4.26	4.36	3.57

Table 5.3: Breakdown of results by different Frequency of data

Type of data	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
CHAPS	arima	2.96	2.96	2.78	2.96	2.76	2.78	1.00
	autoregressor	2.11	2.11	1.67	2.11	1.62	2.78	1.00
	ft_lag_llama	6.04	6.04	6.18	6.04	6.27	4.56	1.00
	ft_timeGPT	2.87	2.87	2.93	2.87	2.93	3.60	1.00
	lag_llama	5.36	5.36	5.71	5.36	5.64	4.80	1.00
	prophet	5.87	5.87	5.78	5.87	5.82	4.13	3.67
	timeGPT	2.80	2.80	2.96	2.80	2.96	2.98	1.00
Commodity	arima	3.19	3.19	3.14	3.19	3.00	5.48	1.95
	autoregressor	6.14	6.14	6.29	6.14	5.81	3.90	3.90
	ft_lag_llama	2.38	2.38	2.24	2.38	2.00	3.14	2.29
	ft_timeGPT	4.05	4.05	4.10	4.05	4.33	4.29	4.48
	lag_llama	4.00	4.00	3.57	4.00	4.19	2.71	3.71
	prophet	4.76	4.76	4.81	4.76	4.81	2.95	5.62
	timeGPT	3.48	3.48	3.86	3.48	3.86	3.81	4.00
Exchange rate	arima	3.17	3.17	2.78	3.17	2.11	6.00	1.33
	autoregressor	6.50	6.50	6.39	6.50	6.44	4.50	4.28
	ft_lag_llama	2.78	2.78	2.56	2.78	2.28	2.50	3.33
	ft_timeGPT	2.83	2.83	3.17	2.83	4.11	2.78	3.83
	lag_llama	5.11	5.11	5.28	5.11	4.94	2.94	4.39
	prophet	3.28	3.28	3.94	3.28	4.22	2.94	4.44
	timeGPT	4.33	4.33	3.89	4.33	3.89	4.89	4.22
Stock index	arima	2.37	2.37	2.49	2.37	1.98	6.04	1.15
	autoregressor	6.75	6.75	6.80	6.75	6.42	4.45	5.15
	ft_lag_llama	2.32	2.32	2.49	2.32	2.65	2.08	2.02
	ft_timeGPT	3.98	3.98	4.00	3.98	3.77	3.56	5.39
	lag_llama	4.08	4.08	4.07	4.08	4.20	2.42	3.46
	prophet	3.92	3.92	3.89	3.92	4.43	3.06	4.56
	timeGPT	4.58	4.58	4.26	4.58	4.55	4.90	4.60

Table 5.4: Breakdown of results by different Type of data

Frequency of data	Type of data	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
daily	CHAPS	arima	2.67	2.67	2.43	2.67	2.43	2.23	1.00
		autoregressor	1.90	1.90	1.50	1.90	1.43	2.33	1.00
		ft_lag_llama	5.57	5.57	5.77	5.57	5.90	5.03	1.00
		ft_timeGPT	3.33	3.33	3.27	3.33	3.27	3.77	1.00
		lag_llama	5.03	5.03	5.57	5.03	5.47	5.30	1.00
		prophet	6.57	6.57	6.47	6.57	6.53	3.87	5.00
		timeGPT	2.93	2.93	3.00	2.93	2.97	3.10	1.00
	Commodity	arima	3.00	3.00	3.11	3.00	3.44	6.67	1.22
		autoregressor	6.78	6.78	6.67	6.78	5.22	3.67	2.89
		ft_lag_llama	2.89	2.89	2.78	2.89	3.33	2.11	2.22
		ft_timeGPT	3.67	3.67	3.22	3.67	4.11	3.44	4.67
		lag_llama	3.67	3.67	3.33	3.67	3.22	3.56	5.56
		prophet	4.67	4.67	4.56	4.67	5.44	2.56	5.56
		timeGPT	3.33	3.33	4.33	3.33	3.22	5.00	4.56
	Exchange rate	arima	3.22	3.22	2.56	3.22	1.89	6.00	1.11
		autoregressor	7.00	7.00	6.78	7.00	6.89	3.89	4.22
		ft_lag_llama	3.89	3.89	3.67	3.89	3.33	2.22	2.22
		ft_timeGPT	3.56	3.56	3.44	3.56	4.56	4.11	5.33
		lag_llama	3.22	3.22	3.56	3.22	2.89	3.56	4.33
		prophet	3.11	3.11	4.67	3.11	5.22	2.00	4.56
		timeGPT	4.00	4.00	3.33	4.00	3.22	4.56	4.78
	Stock index	arima	1.92	1.92	2.14	1.92	2.06	5.89	1.17
		autoregressor	7.00	7.00	7.00	7.00	6.64	5.11	5.86
		ft_lag_llama	3.08	3.08	3.22	3.08	3.72	2.31	2.06
		ft_timeGPT	3.86	3.86	3.58	3.86	3.56	3.50	4.42
		lag_llama	3.47	3.47	3.56	3.47	3.03	2.75	4.42
		prophet	4.61	4.61	4.61	4.61	4.75	2.47	4.19
		timeGPT	4.06	4.06	3.89	4.06	4.25	4.72	4.50

Table 5.5: Breakdown of "daily" results by different Types of data

Frequency of data	Type of data	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
weekly	CHAPS	arima	3.53	3.53	3.47	3.53	3.40	3.87	1.00
		autoregressor	2.53	2.53	2.00	2.53	2.00	3.67	1.00
		ft_lag_llama	7.00	7.00	7.00	7.00	7.00	3.60	1.00
		ft_timeGPT	1.93	1.93	2.27	1.93	2.27	3.27	1.00
		lag_llama	6.00	6.00	6.00	6.00	6.00	3.80	1.00
		prophet	4.47	4.47	4.40	4.47	4.40	4.67	1.00
		timeGPT	2.53	2.53	2.87	2.53	2.93	2.73	1.00
	Commodity	arima	4.11	4.11	3.44	4.11	2.56	5.33	2.44
		autoregressor	5.44	5.44	5.89	5.44	6.11	3.78	4.22
		ft_lag_llama	2.00	2.00	1.67	2.00	1.00	3.56	2.67
		ft_timeGPT	4.33	4.33	4.89	4.33	5.22	4.56	4.56
		lag_llama	4.56	4.56	4.56	4.56	5.22	2.00	2.67
		prophet	4.56	4.56	4.67	4.56	3.89	3.67	5.33
		timeGPT	3.00	3.00	2.89	3.00	4.00	3.11	3.22
	Exchange rate	arima	3.11	3.11	3.00	3.11	2.33	6.00	1.56
		autoregressor	6.00	6.00	6.00	6.00	6.00	5.11	4.33
		ft_lag_llama	1.67	1.67	1.44	1.67	1.22	2.78	4.44
		ft_timeGPT	2.11	2.11	2.89	2.11	3.67	1.44	2.33
		lag_llama	7.00	7.00	7.00	7.00	7.00	2.33	4.44
		prophet	3.44	3.44	3.22	3.44	3.22	3.89	4.33
		timeGPT	4.67	4.67	4.44	4.67	4.56	5.22	3.67
	Stock index	arima	2.78	2.78	2.86	2.78	1.94	6.08	1.19
		autoregressor	6.64	6.64	6.75	6.64	6.28	3.69	4.86
		ft_lag_llama	1.94	1.94	1.94	1.94	2.08	1.83	1.92
		ft_timeGPT	4.22	4.22	4.64	4.22	4.08	3.47	6.33
		lag_llama	5.03	5.03	5.03	5.03	5.14	2.17	2.86
		prophet	2.53	2.53	2.50	2.53	3.67	3.75	4.53
		timeGPT	4.86	4.86	4.28	4.86	4.81	5.14	4.69

Table 5.6: Breakdown of "weekly" results by different Types of data

Frequency of data	Type of data	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
monthly	Commodity	arima	1.00	1.00	2.33	1.00	3.00	2.33	2.67
		autoregressor	6.33	6.33	6.33	6.33	6.67	5.00	6.00
		ft_lag_llama	2.00	2.00	2.33	2.00	1.00	5.00	1.33
		ft_timeGPT	4.33	4.33	4.33	4.33	2.33	6.00	3.67
		lag_llama	3.33	3.33	1.33	3.33	4.00	2.33	1.33
		prophet	5.67	5.67	6.00	5.67	5.67	2.00	6.67
		timeGPT	5.33	5.33	5.33	5.33	5.33	2.33	4.67
	Stock index	arima	2.50	2.50	2.42	2.50	1.83	6.33	1.00
		autoregressor	6.33	6.33	6.33	6.33	6.17	4.75	3.92
		ft_lag_llama	1.17	1.17	1.92	1.17	1.17	2.17	2.25
		ft_timeGPT	3.58	3.58	3.33	3.58	3.50	4.00	5.50
		lag_llama	3.08	3.08	2.75	3.08	4.92	2.17	2.42
		prophet	6.00	6.00	5.92	6.00	5.75	2.75	5.75
		timeGPT	5.33	5.33	5.33	5.33	4.67	4.75	4.58

Table 5.7: Breakdown of "monthly" results by different Types of data

Stationary data	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
FALSE	arima	3.19	3.19	3.00	3.19	3.00	3.38	1.00
	autoregressor	2.24	2.24	1.76	2.24	1.76	2.95	1.00
	ft_lag_llama	6.57	6.57	6.57	6.57	6.57	3.90	1.00
	ft_timeGPT	2.48	2.48	2.71	2.48	2.71	3.71	1.00
	lag_llama	5.81	5.81	5.86	5.81	5.86	4.19	1.00
	prophet	5.29	5.29	5.24	5.29	5.24	4.29	2.71
	timeGPT	2.43	2.43	2.86	2.43	2.86	3.05	1.00
TRUE	arima	2.65	2.65	2.63	2.65	2.23	5.33	1.27
	autoregressor	5.86	5.86	5.82	5.86	5.53	4.08	4.19
	ft_lag_llama	2.92	2.92	3.01	2.92	3.06	2.78	2.05
	ft_timeGPT	3.72	3.72	3.77	3.72	3.79	3.56	4.35
	lag_llama	4.34	4.34	4.39	4.34	4.50	3.00	3.21
	prophet	4.36	4.36	4.41	4.36	4.77	3.18	4.69
	timeGPT	4.16	4.16	3.96	4.16	4.12	4.42	3.88

Table 5.8: Breakdown of results by stationarity

Trend	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
linear	arima	2.64	2.64	2.51	2.64	2.53	2.98	1.00
	autoregressor	2.22	2.22	1.89	2.22	1.87	2.71	1.29
	ft_lag_llama	6.18	6.18	6.20	6.18	6.09	4.69	1.07
	ft_timeGPT	2.91	2.91	3.04	2.91	3.11	3.64	1.13
	lag_llama	5.38	5.38	5.56	5.38	5.53	4.73	1.29
	prophet	5.84	5.84	5.64	5.84	5.67	3.98	3.58
	timeGPT	2.82	2.82	3.16	2.82	3.20	2.93	1.29
no trend	arima	2.74	2.74	2.74	2.74	2.25	5.86	1.32
	autoregressor	6.57	6.57	6.57	6.57	6.23	4.39	4.71
	ft_lag_llama	2.35	2.35	2.45	2.35	2.55	2.28	2.24
	ft_timeGPT	3.80	3.80	3.85	3.80	3.85	3.55	4.96
	lag_llama	4.21	4.21	4.22	4.21	4.35	2.57	3.54
	prophet	3.98	3.98	4.11	3.98	4.52	3.08	4.76
	timeGPT	4.35	4.35	4.07	4.35	4.24	4.73	4.33

Table 5.9: Breakdown of results by trend

Cyclical patterns	Model	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
FALSE	arima	2.51	2.51	2.56	2.51	1.97	5.87	1.36
	autoregressor	6.53	6.53	6.59	6.53	6.40	4.14	4.82
	ft_lag_llama	2.42	2.42	2.41	2.42	2.26	2.26	2.39
	ft_timeGPT	3.77	3.77	3.96	3.77	3.84	3.86	5.04
	lag_llama	4.40	4.40	4.36	4.40	4.62	2.56	3.58
	prophet	4.07	4.07	4.03	4.07	4.38	3.13	4.68
	timeGPT	4.30	4.30	4.10	4.30	4.53	4.53	4.21
TRUE	arima	2.95	2.95	2.82	2.95	2.74	4.19	1.09
	autoregressor	4.10	4.10	3.85	4.10	3.51	3.71	2.60
	ft_lag_llama	4.47	4.47	4.65	4.47	4.94	3.69	1.38
	ft_timeGPT	3.33	3.33	3.27	3.33	3.44	3.26	2.65
	lag_llama	4.67	4.67	4.83	4.67	4.72	3.83	2.19
	prophet	4.95	4.95	5.08	4.95	5.35	3.54	4.17
	timeGPT	3.53	3.53	3.50	3.53	3.31	3.92	2.72

Table 5.10: Breakdown of results by presence of cyclical patterns

Model-parameter experimentation phase results

Time-series chosen for Lag-Llama fine-tuning hyper-parameter optimization was Commodity of weekly frequency, period of 01-01-2019 until 01-01-2024, and Context length of 32. Since there are 4 Lag-Llama fine-tuning parameters we are experimenting with (see Table 4.4), each having 3 possible values, $3^4 = 81$ experiments were conducted.

Since working on a single time-series we can look at the actual performance metrics of fine-tuned Lag-Llama instead of the ranks. Breaking down the fine-tuning experimentation results by Batch size (see Table 5.11), we can tell that Lag-Llama benefits from smaller batch size - arguably³⁵ performing best with Batch size of 5 (15.47% better than other Batch sizes in terms of RMSE), which isn't surprising as lower batch sizes mean more frequent model-weight updates hence more "granular" learning. Breaking down the performance by length of the fine-tune set (see Table 5.12), we can tell that Lag-Llama benefits from higher fine-tune lengths as its' performance was arguably best with the highest fine-tune length: 200 (In terms of RMSE, 16.60% better than Fine-tune length of 128 and 94.41% better than Fine-tune length of 64). Table 5.13 suggests that Lag-Llama performs better with lower Max epochs as it performed overall better with 4 Max epochs than with 8 or 16 (in terms of RMSE, 15.47% better than both other cases). Finally, regarding the Learning rate (see Table 5.14), we cannot tell certainly which learning rate works best for Lag-Llama as it seems to have performed better with the lower (0.00005) and higher one (0.005) than with the default Learning rate of 0.0005.

³⁵Higher batch sizes only outperformed the low batch size in terms of MDA and MES.

model	Batch size	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
ft_lag_llama	20	-0.0113	0.1652	0.0883	0.2625	1.0000	0.5441	0.5333
	10	-0.0113	0.1652	0.0883	0.2625	1.0000	0.5517	0.5333
	5	-0.0105	0.1355	0.0781	0.2219	1.0000	0.5078	0.5212

Table 5.11: Lag-Llama Batch size optimization

model	Fine-tune length	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
ft_lag_llama	200	-0.0069	0.0015	0.0321	0.0392	1.0000	0.5141	0.4667
	128	-0.0018	0.0022	0.0333	0.0470	1.0000	0.5632	0.5000
	64	-0.0251	0.4919	0.1995	0.7014	1.0000	0.5249	0.6333

Table 5.12: Lag-Llama fine-tune length optimization

model	Max epochs	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
ft_lag_llama	16	-0.0113	0.1652	0.0883	0.2625	1.0000	0.5517	0.5333
	8	-0.0113	0.1652	0.0883	0.2625	1.0000	0.5785	0.5333
	4	-0.0105	0.1355	0.0781	0.2219	1.0000	0.4796	0.5212

Table 5.13: Lag-Llama Max epochs optimization

model	Learning rate	R ²	MSE	MAE	RMSE	MAPE	MDA	MES
ft_lag_llama	0.00005	-0.0069	0.0015	0.0321	0.0392	1.0000	0.5575	0.4667
	0.0005	-0.0110	0.1539	0.0844	0.2471	1.0000	0.5327	0.5287
	0.005	-0.0069	0.0015	0.0321	0.0392	1.0000	0.5632	0.4667

Table 5.14: Lag-Llama Learning rate optimization

5.2 Limitations

Information disparity between fine-tuned TSFMs and benchmark models + zero-shot TSFMs

Theoretically, fine-tuned TSFMs are systematically advantaged over benchmark models and zero-shot TSFMs as they have access to more information as the set that they are fine-tuned on is larger than the size of the context window which zero-shot TSFMs and benchmark models use (they are using the entire train set TR_i of the fold F_i for fine-tuning, whereas the benchmark models and zero-shot only have access to the last cl data-points from TR_i). However, this is unavoidable as the technical implementation of fine-tuning in TSFMs is such that the length of the fine-tuning set has to be larger than its' context length - so there is bound to be information disparity between zero-shot and fine-tuned TSFMs. If we simply let benchmark models train on the whole TR_i , then zero-shot TSFMs would be systematically disadvantaged in comparison to benchmark models as they would have access to less information than them. The final option is to let zero-shot TSFMs have the context length of same length as the length of the fine-tuning set, however, then the comparison between zero-shot and fine-tuned TSFMs wouldn't be appropriate as fine-tuned TSFMs would have shorter context lengths and the amount of direct information going into zero-shot and fine-tuned TSFMs would be different.

Statistical significance

Unfortunately, due to the nature of this experiment, it is difficult to quantify the statistical significance of the results within time-series. The way that statistical significance of time-series prediction evaluation metrics is calculated is by looking at their performance across different time-series and then calculating the mean and variance of the performance metrics and then calculating confidence intervals. In this case, this doesn't make much sense because, as we saw, time-series used in this research are inherently very diverse in terms of their characteristics, and there simply isn't enough time-series (see Section 5.2.3) with similar characteristics so that hypothesis testing could be done. Instead, reaching statistical significance has been attempted through the use of the TSCV (see Figure 4.4)

where the 6 groups of 5 individual predictions are spread uniformly across each time-series in order to cover the largest period possible so that the likelihood of getting an accidental localized result is minimized.

Computational resources limitations

The main limitation of this research is the lacking computational resources. The code which runs the experiments takes a relatively long time to execute (around 10-45 minutes for each experiment depending on Max epochs and Batch size) because Lag-Llama is run locally. This limitation led to multiple "shortcuts" being taken in this research process - such as:

- Limiting the number of time-series used in this research.
- Limiting the parameter-space across which was experimented.
- Using the heuristic approach with fixed model-parameters and varying data-parameters, and vice-versa in the Model-parameter experimentation phase.
- Not fine-tuning Lag-Llama before every prediction but rather every five predictions. (see Section 5.2.7)
- Not conducting multiple repeats of the same experiment in order to gain statistical confidence in results.
- Doing "only" 30 prediction per each time-series.

TimeGPT-1 API calls

Second limitation of this research is the TimeGPT-1 API call limit. With paid subscription, it is limited to 10,000 calls per month - which was raised to 30,000 per my request. However, 30,000 was not enough to fully conduct the research. Each TimeGPT-1 prediction requires 2 API calls - which means that per experiment we are using 4 API calls (2 for zero-shot TimeGPT-1 and 2 for fine-tuned version). 4 calls per prediction times 30 predictions equals 120 API calls. This means that with the 30,000 limit, I was able to make 250

experiments. After the code testing and a few failed attempts, there was barely enough API calls to finish the research up to this level. due to this, unfortunately few avenues of research had to be left out such as timeGPT-1 fine-tuning parameter optimization and timeGPT-1 exogenous variables research.

TimeGPT-1 pre-training data

Third limitation is that we don't have public information on timeGPT-1 pre-training data and therefore we cannot be sure that timeGPT-1 wasn't already pre-trained on the data we used to evaluate its' performance (data leakage). I sent a query to timeGPT-1 team (Nixtla labs) but they refused to disclose the information. However, judging by timeGPT-1 overall performance, it doesn't seem it has already been pre-trained on the data which was used. One solution to this would've been running the experiments only on most recent data, however, this was subject to the API limitation and data availability limitation.

Data availability

Fourth limitation of this research is data availability. AlphaVantage API doesn't provide monthly Exchange rate data before 2015 which isn't enough to conduct experiments on. CHAPS monthly data isn't available before 2020 which also wasn't enough to conduct proper experiment. In addition Alphavantage API doesn't provide commodity price information for commodities of higher than "monthly" frequency unless for "WTI" and "BRENT".

Fine-tuning disparity between Lag-Llama and TimeGPT-1

Theoretically, fine-tuned versions of TSFMs should be fine-tuned on every single fold as there being distance between the fine-tuning data and the timepoint of prediction increases the likelihood that the model is learning outdated patterns. However, as Lag-Llama is run locally, fine-tuning process is taking a very long time and it was unfeasible to fine-tune it before every single prediction. So a compromise had to be made and fine-tune it only every 5 predictions.

5.3 Discussion

Overall results discussion

Between TimeGPT-1 and Lag-Llama, Lag-Llama seems to be the superior model if we consider the fine-tuned versions. Most interesting observations are that fine-tuned Lag-Llama performs best on data which are of "Commodity", "Exchange rate" and "Stock index" type (see Table 5.4). This can be explained through the features of these types of data and frequencies. Looking at Figure 4.3, we can see those types of data are the ones that mostly do not contain cyclical patterns - "Commodity" containing the least cyclical patterns of them all and also being the type of data on which fine-tuned Lag-Llama performs best on in terms of the rank difference between it and the next best model³⁶. We are seeing a similar kind of pattern with results breakdown by "frequency". Lag-Llama is the best model on data of "monthly" and "weekly" frequencies (see Table 5.3) - which if we look at Figure 4.3 we can see are the frequencies with least cyclical patterns. Similarly as before - the frequency which has less cyclical patterns: "monthly" is the one on which fine-tuned Lag-Llama performs comparatively better on than of "weekly" frequency³⁷. These facts point to a conclusion that fine-tuned Lag-Llama performs really well on data with no cyclical pattern and quite poorly on data with a cyclical pattern. This is further evidenced in Table 5.10 where we see fine-tuned Lag-Llama is arguably the best model on non-cyclical time series and one of the worst on cyclical time series. If we look at "trend", from Type of data - point of view, evidence is very strong that presence of linear trend negatively affects the performance of Fine-tuned Lag-Llama as Commodity, Exchange rate and Stock index almost all contain no trend and fine-tuned Lag-Llama performs excellent. From Frequency perspective, we can conclude the same as monthly frequency contains no trend and fine-tuned Lag-Llama performs by far the best on that data. Where daily frequency in almost 40% cases contains linear trend and it also seems to be the frequency on which fine-tuned Lag-Llama performs worst. Same case can be said for stationarity of time-series.

³⁶Difference (RMSE) between fine-tuned Lag-Llama and next best model was 0.81 on commodity, 0.05 on exchange rate and 0.05 on stock index.

³⁷Difference (RMSE) between fine-tuned Lag-Llama and next best model was 0.87 on monthly data but only 0.15 on weekly data.

TimeGPT-1 models' performance is a little bit disappointing. Not in many cases they outperform other models and in the cases when they do, it is hard to find pattern which could explain why they do - leading to believe that their good performance on those specific occasions could be accidental. Only time TimeGPT model appears the best model (fine-tuned version) in more than one metric is in the weekly CHAPS data - scoring 1.93 in R^2 , MSE and RMSE, whereas the 2nd best model came close behind and scored 2.53 in those respective metrics. An interesting observation is that when breaking down the results by time-series characteristics (Tables 5.8 to 5.10), effect of time-series characteristics on fine-tuned TimeGPT-1 performance is diametrically opposite than on fine-tuned Lag-Llama. Whereas fine-tuned Lag-Llama works relatively better with stationary data, fine-tuned TimeGPT-1 works relatively better with non-stationary data. Same goes for cyclical and trend.

General Lag-Llama great MDA performance

Another interesting observation is that both fine-tuned and zero-shot versions of Lag-Llama perform quite well in terms MDA across the board with fine-tuned Lag-Llama being the best overall (see Table 5.1), and either one of them being the best in all cases of Context length (Table 5.2), with monthly and weekly frequencies (Table 5.3), Commodity, Exchange rate and Stock index data (Table 5.4). From the previous paragraph, we know that fine-tuned Lag-Llama performs best with stationary data with no trend and no cyclical patterns. If we look at MDA metric only, we can see that in all those cases (Tables 5.8, 5.9, 5.10) fine-tuned Lag-Llama is the best model by each of those time-series characteristics, and zero-shot Lag-Llama is always the 2nd best model. This is unlikely to be a coincidence and likely points to the fact that general Lag-Llama model has superior capability to predict to direction of the movement of a time-series given stationarity, no trend and no cyclical pattern.

Inter-model fine-tuning performance discussion

Across all results, Lag-Llama seems to universally benefit from fine-tuning, only exception to this being "daily" and "weekly" CHAPS and "daily" Exchange rate data where zero-shot

and fine-tuned versions performs actually better (Tables 5.5, 5.6). Looking at time-series characteristics, zero-shot Lag-Llama performs better than fine-tuned version on non-stationary data (Table 5.8), on data with linear trend (Table 5.9), and almost equally on time-series with a cyclical pattern (Table 5.10). This is a very interesting because one would expect that generally a model's performance would generally improve with fine-tuning, however, according to these results, on time-series with these characteristics; Lag-Llama's performance actually worsens with fine-tuning. This phenomenon brings sheds some light on Lag-Llama's fine-tuning capabilities and shows which kind of characteristics it is able to learn (stationarity, no trend, no cyclical) and which it isn't (non-stationarity, linear trend and cyclical). Overall, TimeGPT seems to benefit from fine-tuning only slightly. Only cases where TimeGPT-1 benefits significantly from fine-tuning is on "weekly" Exchange rate data, and all "monthly" data (Tables 5.6, 5.7). Zero-shot TimeGPT-1 even beat the fine-tuned version on "daily" CHAPS, "daily" and "weekly" Commodity data and "weekly" Stock index data (Tables 5.5, 5.6). A very intriguing observation is that on time-series characteristics where fine-tuned TimeGPT performs well compared to other models (non-stationary and linear trend), the zero-shot TimeGPT-1 actually outperforms fine-tuned version (Tables 5.8, 5.9) and comes very close on cyclical data (Table 5.10). This is complete opposite fine-tuning behaviour than Lag-Llama's whos' fine-tuned version was dominating the zero-shot version on characteristics it was good at (stationary, no trend, no seasonality) compared to characteristics it was bad at.

Lag-Llama fine-tuning discussion

Fine-tuning - wise, Lag-Llama seems to perform best with higher fine-tune lengths, and lower Batch sizes and Max epochs. However there are a few caveats. Fine-tune length of 128 produced a better R^2 (73.91% better) than the fine-tune length of 200 meaning that increasing the fine-tune length ad infinitum might not result in strictly better performance. This is likely because the nature of financial time-series data. Patterns don't stay constant over time and having longer fine-tuning lengths means the model is learning older and less relevant patterns which might have explained why fine-tune length of 128 was at least in one metric better than fine-tune length of 200. This extends to the fine-tune length of 64 which had significantly higher MES than both 128 and 200 fine-tune lengths. Max

epochs - wise, in theory higher Max epochs shouldn't lead to detriment in performance as the model discards the new weights in every epoch if that epoch's loss was higher than previous epoch's. Even though it is not certain, having higher Max epochs could give the model more chances to overfit the fine-tuning data - hence make worse predictions. Last interesting observation of the fine-tuning results is that There are a lot of results that look very similar. Looking at the Tables 5.11 and 5.13, they look nearly identical with some minor differences in the MDA column. This could be an indication that there were mistakes in results recording process, however, after detailed investigation, I concluded that there was no mistake, rather an interesting phenomena was occurring. In majority of the experiments in this set, Lag-Llama was making predictions which were in the order of magnitude $\times 10^{-14}$ (near zero) (see Figure 5.1). In every single experiment, it did make different predictions because it had a different set of fine-tuning parameters, however, because the scale of the predictions was so small, the differences in the predictions were often only in the 14th, 15th, 16th decimal place - meaning that they would have the same evaluation metrics if we're considering only 4 decimal places (see Tables 5.11 to 5.14). Likely reason for that is that weekly WTI returns time-series looks inherently predictable, therefore, model tries to avoid making large mistakes therefore makes safe predictions. Furthermore around at around 1/3 of the time-series, there was a period of very large outliers. Where WTI returns sharply fell and then rose to over 350% (see Figure 5.2). Having used this volatile time-period in fine-tuning in most cases (depending on fine-tune-length and which fold of the TSCV), it is likely that the model-weights got stuck in a local minimum - hence similar predictions - hence near identical evaluation metrics in many cases. However, even though the the model's predictions have degraded to predicting near-0 all the time, we can still see that in every single scenario, the model has $>50\%$ MDA, meaning that it has more than 50% chance of predicting the correct direction of the time-series movement in every configuration of fine-tuning parameters.

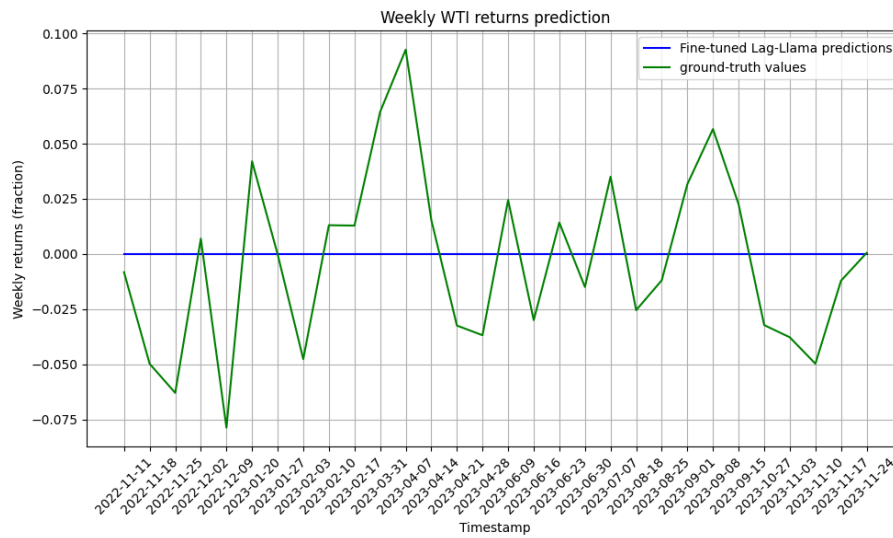


Figure 5.1: Predictions of Lag-Llama with Batch size 5, fine-tune length 200, Max epochs 8 and Learning rate 0.0005 on weekly WTI returns

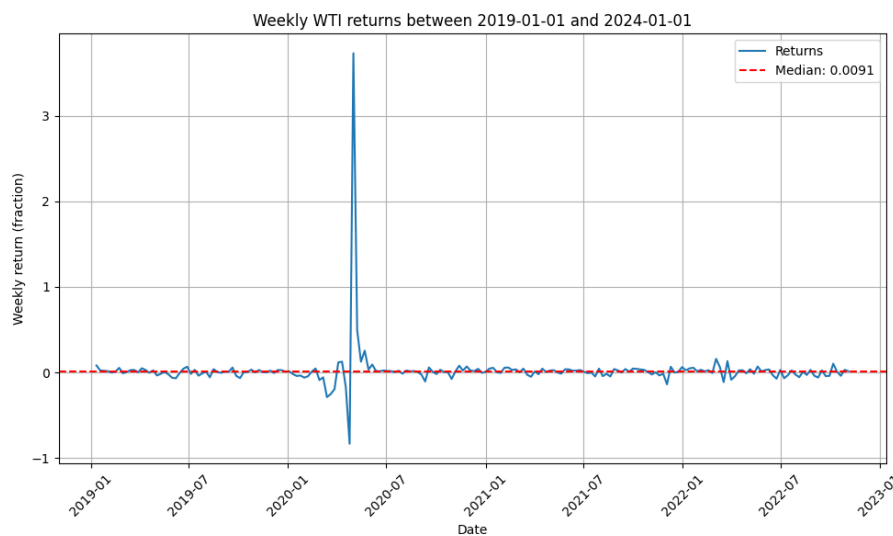


Figure 5.2: WTI weekly returns between 2019-01-01 and 2024-01-01

6. CONCLUSION AND FUTURE WORK

6.1 Conclusion

Experiments were conducted across different types of data, frequencies and Context lengths - providing an interesting breakdown of results and how the TSFMs behave on different time-series, with different frequencies, and what are some of the fundamental time-series characteristics which cause them to perform better or worse in relation to other models. Many different combinations of Lag-Llama fine-tuning parameters were experimented with on a single time-series and we've shed light how the fine-tuned model behaves in the face of a very unpredictable time-series with extreme outliers during this process (resorts to near-median prediction). Furthermore, We've went into a deep analysis of TSFM's fine-tuning behaviour and uncovered the strange phenomenon of how the same time-series characteristics which make fine-tuned lag-Llama perform worse, make the zero-shot version perform better and vice versa. We have also discovered how time-series with characteristics beneficial to TimeGPT-1 are actually seeing the zero-shot TimeGPT-1 perform even better than its' fine-tuned version and vice-versa. Finally we point out the unique ability of both fine-tuned and zero-shot Lag-Llama to perform in terms of directional accuracy.

Overall, have seen good performance from Lag-Llama on certain financial time-series. Overall, the fine-tuned version outperforms TimeGPT-1 and the benchmark models on time-series with no cyclical patterns and no trend. Crucial for Lag-Llama's performance is that the data it is predicting is stationary. Furthermore, both fine-tuned and zero-shot versions of Lag-Llama have proved better than all other models in cases where time-series meet mentioned conditions - which points to the model's inherent ability to learn from complex patterns in the data. TimeGPT-1 hasn't showed great performance overall, managing to come out on top on a very small set of experiments (weekly CHAPS data). Another interesting results of this research is that TimeGPT-1 and Lag-Llama are complementary - meaning that the time-series on which Lag-Llama performs worse on, TimeGPT-1 performs better on, and vice-versa.

6.2 Future work

In a way, this research has opened more questions than it closed. First avenue of future research is to expand the volume of experiments conducted to more different asset classes such as option prices, stock trading volumes, etc., and more different frequencies (minutely, hourly, quarterly), and over more different time-periods. This research could also include more broad classification of time-series characteristics which could include testing for quadratic, exponential, damped trend as well as make distinction between additive and multiplicative seasonality of the time-series. This type of analysis could also include additional time-series features mentioned in this research but not tested for such as the volatility of time-series and the structure of errors.

Another avenue of research could be towards fine-tuning optimization. Analysis which is missing from this research due to certain limiting factors (see Section 5.2.3), is between optimal fine-tuning hyper-parameters and the time-series characteristics, i.e. how to choose the optimal fine-tuning hyper-parameters based on the characteristics of the time-series we are trying to predict. Furthermore, this avenue could also look into TimeGPT-1 fine-tuning optimization as well as use of exogenous variables which are aspects missing from this research due to low API call limit (Section 5.2.4).

Lastly, an interesting avenue could be towards hybrid models. This kind of research could involve time-series decomposition by trend, cycles and errors where Lag-Llama could be used to predict the components of a time-series which are stationary and have no trend and cyclical patterns, and another model (TimeGPT-1 perhaps) could be used to predict the components which are opposite.

REFERENCES

- [1] A. F. Ansari, L. Stella, C. Turkmen, X. Zhang, P. Mercado, H. Shen, O. Shchur, S. S. Rangapuram, S. P. Arango, S. Kapoor, et al. Chronos: Learning the language of time series. *arXiv preprint arXiv:2403.07815*, 2024.
- [2] V. Assimakopoulos and K. Nikolopoulos. The theta model: a decomposition approach to forecasting. *International journal of forecasting*, 16(4):521–530, 2000.
- [3] J. Ba. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [4] D. Bahdanau. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [5] G. Box and G. Jenkins. *Time Series Analysis: Forecasting and Control*. Holden-Day series in time series analysis and digital processing. Holden-Day, 1970.
- [6] R. G. Brown. *Exponential smoothing for predicting demand*. Little, 1956.
- [7] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.
- [8] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794, 2016.
- [9] Z. Chen, L. N. Zheng, C. Lu, J. Yuan, and D. Zhu. Chatgpt informed graph neural network for stock movement prediction. *arXiv preprint arXiv:2306.03763*, 2023.
- [10] K. Cho. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [11] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, Y. Li, X. Wang, M. Dehghani, S. Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- [12] R. B. Cleveland, W. S. Cleveland, J. E. McRae, I. Terpenning, et al. Stl: A seasonal-trend decomposition. *J. off. Stat*, 6(1):3–73, 1990.

- [13] M. Costantini, J. C. Cuaresma, and J. Hlouskova. Forecasting errors, directional accuracy and profitability of currency trading: The case of eur/usd exchange rate. *Journal of Forecasting*, 35(7):652–668, 2016.
- [14] A. Das, W. Kong, R. Sen, and Y. Zhou. A decoder-only foundation model for time-series forecasting. *arXiv preprint arXiv:2310.10688*, 2023.
- [15] J. Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [16] D. A. Dickey and W. A. Fuller. Distribution of the estimators for autoregressive time series with a unit root. *Journal of the American statistical association*, 74(366a):427–431, 1979.
- [17] J. Dong, H. Wu, Y. Wang, Y. Qiu, L. Zhang, J. Wang, and M. Long. Timesiam: A pre-training framework for siamese time-series modeling. *arXiv preprint arXiv:2402.02475*, 2024.
- [18] A. Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [19] S. Gao, T. Koker, O. Queen, T. Hartvigsen, T. Tsiligkaridis, and M. Zitnik. Units: Building a unified time series model. *arXiv preprint arXiv:2403.00131*, 2024.
- [20] A. Garza and M. Mergenthaler-Canseco. Timegpt-1. *arXiv preprint arXiv:2310.03589*, 2023.
- [21] R. Gençay and M. Qi. Pricing and hedging derivative securities with neural networks: Bayesian regularization, early stopping, and bagging. *IEEE transactions on neural networks*, 12(4):726–734, 2001.
- [22] J. Hasanhodzic and A. W. Lo. Can hedge-fund returns be replicated?: The linear case. *The Linear Case (August 16, 2006)*, 2006.
- [23] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

- [24] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [25] S. Hochreiter. Long short-term memory. *Neural Computation MIT-Press*, 1997.
- [26] Z. Hu, Y. Zhao, and M. Khushi. A survey of forex and stock price prediction using deep learning. *Applied System Innovation*, 4(1):9, 2021.
- [27] R. J. Hyndman and Y. Khandakar. Automatic time series forecasting: the forecast package for r. *Journal of statistical software*, 27:1–22, 2008.
- [28] R. J. Hyndman and A. B. Koehler. Another look at measures of forecast accuracy. *International journal of forecasting*, 22(4):679–688, 2006.
- [29] R. J. Hyndman, A. B. Koehler, R. D. Snyder, and S. Grose. A state space framework for automatic forecasting using exponential smoothing methods. *International Journal of forecasting*, 18(3):439–454, 2002.
- [30] M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1-2):81–93, 1938.
- [31] Z. Lan. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [32] S. Li, X. Jin, Y. Xuan, X. Zhou, W. Chen, Y.-X. Wang, and X. Yan. Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *Advances in neural information processing systems*, 32, 2019.
- [33] Y. Liang, H. Wen, Y. Nie, Y. Jiang, M. Jin, D. Song, S. Pan, and Q. 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*, pages 6555–6565, 2024.
- [34] Y. Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [35] Y. Liu, H. Zhang, C. Li, X. Huang, J. Wang, and M. Long. Timer: Transformers for time series analysis at scale. *arXiv preprint arXiv:2402.02368*, 2024.

- [36] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- [37] S. Makridakis and M. Hibon. The m3-competition: results, conclusions and implications. *International journal of forecasting*, 16(4):451–476, 2000.
- [38] S. Makridakis, E. Spiliotis, and V. Assimakopoulos. M5 accuracy competition: Results, findings, and conclusions. *International Journal of Forecasting*, 38(4):1346–1364, 2022.
- [39] F. J. Massey Jr. The kolmogorov-smirnov test for goodness of fit. *Journal of the American statistical Association*, 46(253):68–78, 1951.
- [40] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5:115–133, 1943.
- [41] J. A. Miller, M. Aldosari, F. Saeed, N. H. Barna, S. Rana, I. B. Arpinar, and N. Liu. A survey of deep learning and foundation models for time series forecasting. *arXiv preprint arXiv:2401.13912*, 2024.
- [42] Y. Nie, N. H. Nguyen, P. Sinthong, and J. Kalagnanam. A time series is worth 64 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*, 2022.
- [43] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [44] K. Rasul, A. Ashok, A. R. Williams, A. Khorasani, G. Adamopoulos, R. Bhagwatkar, M. Biloš, H. Ghonia, N. V. Hassen, A. Schneider, et al. Lag-llama: Towards foundation models for time series forecasting. *arXiv preprint arXiv:2310.08278*, 2023.
- [45] A. Roberts, C. Raffel, K. Lee, M. Matena, N. Shazeer, P. J. Liu, S. Narang, W. Li, and Y. Zhou. Exploring the limits of transfer learning with a unified text-to-text transformer. *Google, Tech. Rep.*, 2019.

- [46] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- [47] V. Sanh. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [48] N. Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- [49] E. SLUTSKY. on a case where the law of large numbers applies to mutually dependent quantities. the extension of a theorem by magoichirô watanabe. *Tohoku Mathematical Journal, First Series*, 28:26–32, 1927.
- [50] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- [51] Student. The probable error of a mean. *Biometrika*, pages 1–25, 1908.
- [52] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [53] I. Sutskever. Sequence to sequence learning with neural networks. *arXiv preprint arXiv:1409.3215*, 2014.
- [54] I. Svetunkov, N. Kourentzes, and J. K. Ord. Complex exponential smoothing. *Naval Research Logistics (NRL)*, 69(8):1108–1123, 2022.
- [55] P. Tang and X. Zhang. Mtsmae: Masked autoencoders for multivariate time-series forecasting. In *2022 IEEE 34th International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 982–989. IEEE, 2022.
- [56] S. J. Taylor and B. Letham. Forecasting at scale. *The American Statistician*, 72(1):37–45, 2018.
- [57] P. Temin. Price behavior in ancient babylon. *Explorations in Economic History*, 39(1):46–60, 2002.

- [58] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [59] A. Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- [60] Q. Wen, J. Gao, X. Song, L. Sun, H. Xu, and S. Zhu. Robuststl: A robust seasonal-trend decomposition algorithm for long time series. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 5409–5416, 2019.
- [61] C. Wimmer and N. Rekabsaz. Leveraging vision-language models for granular market change prediction. *arXiv preprint arXiv:2301.10166*, 2023.
- [62] H. Wold. *A study in the analysis of stationary time series*. PhD thesis, Almqvist & Wiksell, 1938.
- [63] G. Woo, C. Liu, A. Kumar, C. Xiong, S. Savarese, and D. Sahoo. Unified training of universal time series forecasting transformers. *arXiv preprint arXiv:2402.02592*, 2024.
- [64] Q. Xie, W. Han, Y. Lai, M. Peng, and J. Huang. The wall street neophyte: A zero-shot analysis of chatgpt over multimodal stock movement prediction challenges. *arXiv preprint arXiv:2304.05351*, 2023.
- [65] L. Xue. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*, 2020.
- [66] X. Yu, Z. Chen, Y. Ling, S. Dong, Z. Liu, and Y. Lu. Temporal data meets llm—explainable financial time series forecasting. *arXiv preprint arXiv:2306.11025*, 2023.
- [67] G. U. Yule. On the time-correlation problem, with especial reference to the variate-difference correlation method. *Journal of the Royal Statistical Society*, 84(4):497–537, 1921.
- [68] G. U. Yule. Vii. on a method of investigating periodicities disturbed series, with special reference to wolfer’s sunspot numbers. *Philosophical Transactions of the*

Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character, 226(636-646):267–298, 1927.

- [69] B. Zhang and R. Sennrich. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [70] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, and W. Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115, 2021.
- [71] T. Zhou, Z. Ma, Q. Wen, X. Wang, L. Sun, and R. Jin. Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. In *International conference on machine learning*, pages 27268–27286. PMLR, 2022.

A. HISTORY OF TIME-SERIES FORECASTING

A.1 Brief History of time-series prediction in Finance

When the agricultural revolution took place (circa 10,000 BC), for the first time in history, humans have started producing more food and items than they required just to survive. This gave birth to the concept of "storing" things and eventually trading them for other things which were perceived to be of greater value. When sufficient amount of people started transacting goods in this manner, the concept of "market price" took hold. It did not take long before humans realized that they didn't have to work in the field the whole day to amass wealth, but that they could do it through trade alone. They could buy items from one person at a low price and sell to another at a high price (arbitrage) or buy it now at a low price and hopefully sell later at a high price (asset appreciation). Opportunities to exploit the first option became more rare and less profitable as more and more people started engaging in trade and eliminating these arbitrage opportunities. However, the 2nd option has captured the imagination of traders ever since 10,000BC as no one seemed to be able to tell with certainty whether the price of a good would indeed go up in the future or not. As humans learned to write and keep records, they started noting down the prices of goods that were being traded along with the time when those prices were prevailing. We have records of ancient Babylonians keeping historic records of end-of-month barley, dates, cuscutea, sesame, cardamom and wool prices on clay tablets spanning some 400 years [22]. Investigating these prices, we can see that they follow a random walk pattern with exogenous shocks (such as the death of Alexander the Great) [57] and they seem almost impossible to predict with a naked eye, however, it is not impossible that someone might have tried and unknowingly became the first person in history to work on time-series prediction in the field of Finance.

Over the centuries, methods to solve the problem of time-series prediction in finance seem to have developed universally across many different cultures. Christopher Kurz, a 16th century merchant from Antwerp, thought about the idea of "trend" and "seasonality" in agricultural prices and claimed he could predict prices 20 days in advance. In 18th century Japan, Munehisa Honma - a man dubbed "God of the markets" by his contemporaries, developed a principle stating that when prices become extremely high, they must come

down again which we know today to be the principle of "mean reversion". In the late 19th century, Charles Dow, co-founder of Wall Street Journal and Dow Jones Company³⁸ popularized and coined the term "technical analysis" (a collection of methods that employ math and statistics to predict future prices) which later evolved into quantitative analysis - which is used today in institutions engaged in future price prediction.

A.2 Brief History of modern time-series prediction methods

Statistical models

Although the field of future price prediction has moved far beyond simply considering just the historical prices³⁹, there has been notable work on the general field of univariate⁴⁰ time-series prediction. It is hard to pinpoint the exact start of modern time-series prediction, but I think a good candidate would be Udny Yule's⁴¹ 1927 [68] paper which is considered first to have used simple autoregressor model (AR)⁴² to predict Wolfer's Sunspot numbers⁴³. Even though Slutsky (1927) [49] and Udny (1921) [67] demonstrated the use of Moving Average as a way to show trend and cyclicity in time-series data, Herman Wold's work (1938)⁴⁴ [62] indirectly invented the Auto Regressive Moving Average (ARMA) model⁴⁵. ARMA model is theoretically sound if the time-series being predicted is stationary⁴⁶, which is not the case with all time-series. In 1970, George Box and Gwilym Jenkins introduced probably the most famous time-series prediction model of all times: ARIMA [5]. ARIMA(p, i, q) model is an ARMA(p, q) model that employs the method of

³⁸Dow Jones Company is the publisher of the prominent DOWJ stock market index

³⁹Quantitative analysis now-adays considers wide palette of data and information

⁴⁰Univariate time-series prediction is a practice of using only the past values of a certain time-series in order to predict its future value(s).

⁴¹Udny studied at UCL under a well-known (and controversial) professor Karl Pearson - father of mathematical statistics and generally one of the most well-known personas in statistics.

⁴²An AR(p) time-series prediction model predicts future values based on a linear combination of the previous p observations in the series, assuming that the current value depends on its own past values and a stochastic error term.

⁴³Certain time-series data from the field of Astronomy.

⁴⁴Wold's Decomposition Theorem states that any stationary time series can be decomposed into a deterministic part (AR component) and a stochastic part (MA component)

⁴⁵The MA(q) part of the ARMA(p, q) model forecasts future values by modeling the current value as a linear combination of the past q error terms (shocks) and a stochastic error term.

⁴⁶A time series is said to be stationary if its statistical properties, such as mean, variance, and autocorrelation, remain constant over time

differencing⁴⁷ on the date before training and prediction.

Another direction of innovation in the field of univariate time-series prediction was towards "Exponential Smoothing" (ES). Concept of ES originates from the work of Brown 1956 [6] and has been thoroughly built upon since. Most notable invention was by Hyndman et al. in 2002 [29] - implementing exponential smoothing within a state-space framework, giving us the ETS⁴⁸ model. Most recent culmination in the ES family of models was by Ivan Svetunkov in 2022 [54] with the CES model⁴⁹. Final honorable mention among the time-series prediction models of "statistical" nature is the "Theta" model. Vassilis Assimakopoulos, and Konstantinos Nikolopoulos (2000) [2] Introduced the Theta model which decomposes a time-series into its' cyclical and trend components which are then forecasted separately and finally combined⁵⁰.

Machine Learning (ML) models

As the hardware technology advanced in the 21st century, so did the time-series models as the scientists have finally gotten machines with enough computational power to run more complex algorithms. The most famous of these algorithms being the Artificial Neural Network (ANN). Although the idea (in its primitive form⁵¹) dates back all the way to McCulloch and Pitts (1943) [40], ANN first real ancestor was the "Perceptron"⁵², invented by a psychologist - Frank Rosenblatt in 1958 [46]. The perceptron had ability to take continuous inputs and proper learning algorithm - known even to this day as the "Rosenblatt algorithm". Through the rest of the 20th century, there have been general

⁴⁷Differencing is a technique in time-series analysis used to transform a non-stationary series into a stationary one by deducting the preceding observation from each observation. The order "i" of the differencing refers to the ammount of times this procedure is applied to a time-series.

⁴⁸ETS stands for Error Trend Seasonality. ETS is a family of 18 models. ETS model can account for errors being additive or multiplicative, and trend or seasonality being none, additive or multiplicative. ETS models with no multiplicative errors or seasonality have their equivalent withing the ARIMA family of models

⁴⁹CES stands for Complex Exponential Smoothing model. CES method models the time-series as a series of complex numbers where the real part is the prediction and imaginary parts are errors. CES models can also be expressed in state-space form to adress seasonality. CES advantages are that it is flexible and has been empirically shown to perform well.

⁵⁰These components are called "theta lines". Theta model can use exponential smoothing or an AR model to predict these two lines. Predictions are finally combined using weighted average where weights could be equal or adjusted to favor either the trend line or the cycle line.

⁵¹The original vision of an ANN was based on simple binary inputs with fixed thresholds activation - a far cry from modern ANNs

⁵²Perceptron can be thought of as a single-neuron ANN. ANNs are comprised of one or more neuron layers, each layer containing one or more neurons.

periods of great interest and disinterest in ANNs for many reasons, but ANNs finally came on their own in the 21st century when computers become powerful enough to be able to train larger versions of these models⁵³ in reasonable time⁵⁴. Naturally it was not long before people realized ANN-type models can also be used for time-series prediction⁵⁵ which led to many variations of the ANN architecture to be invented. Bayesian regularization ANN has been used in Finance [21]. Another successful variation of ANN is the Radial Basis Function ANN (RBFANN).

Chen and Guestrin (2016) [8] introduced the Xgboost model which has

Deep Learning models

⁵³Size of an ANN matters when dealing with very complex tasks (tasks with many input features).

⁵⁴Another reason for explosion in popularity was the beginning of the internet era and sudden availability of large ammounts of data which are needed for ANN training

⁵⁵It is important to note that univariate time-series prediction ANNs are simply regression ANNs where lags of the data we're trying to predict are explanatory features (inputs) to the model

B. BACKGROUND

B.1 The transformer

Token

A token is a multi-dimensional numeric vector representation of an element in the sequence. Reason why sequence of tokens is used as input to a sequence-to-sequence model is because ML models can only "understand" numbers - therefore they can't work with some type of sequences (eg. human language sentences). Reason why tokens are vectors rather than simple numbers is because they are designed to represent the semantic meaning of a real-world element they represent and the semantic meaning of a real-world element could change depending on the context. Therefore different dimensions of a token account for different meaning of that element in different contexts.

LSTM and GRU Hidden state

The hidden state is an intermittent sequence within a sequence-to-sequence model. It is a product of applying the encoder on the input sequence. The output sequence is generated by applying the decoder to the hidden state. Having a hidden state of fixed size means that some information will inherently get lost when a very long input sequence is fed into the model.

C. LAG-LLAMA

C.1 LLaMA

RMSnorm is a layer normalization technique. Zhang and Sennrich (2019) demonstrate that centering component of the LayerNorm technique is dispensable (and computationally expensive) and propose their own layer normalization strategy called RMSnorm which delivers similar benefits (more stable training and better model convergence) but with lower computational cost.

RoPE stands for Rotary Positional Encoding. This is a method for positional encoding. RoPE enables valuable properties, including the flexibility of sequence length, decaying inter-token dependency with increasing relative distances, and the capability of equipping linear self-attention with relative position encoding.

SwiGLU activation function is a combination of swish and GLU activation functions. With beta hyperparameter equal to 1, it is of similar shape as ReLU but completely smooth (continuous) - therefore "behaving better" during model training.

D. METHODOLOGY

D.1 Time-series prediction

Prediction error

Prediction error is the difference between the value a model predicts and the actual (ground-truth) value.

MDA

MDA is time-series specific. Since time-series data is inherently ordered (as opposed to regular regression data), we can measure the direction in which a time-series is moving at each time-step. Hence we can compare the direction of movement of our predicted time-series against the actual time-series and we can calculate in what percentage of cases did the two time-series move in the same direction (up or down).

MES

Reason why MDA doesn't fully capture the model's ability to predict direction of underlying asset value movement is because the models are working with returns time-series, therefore MDA measures whether the model accurately predicts the directional change of returns and not the actual direction of movement of the underlying values. Imagine scenario where $y_{\text{actual}}^T = 1\%$, $y_{\text{actual}}^{T+1} = 0.5\%$ and $y_{\text{predicted}}^T = 0.7\%$, $y_{\text{predicted}}^{T+1} = -0.5\%$. In this case MDA would account this as a correct directional guess as the model predicted y would go down in period $T+1$, and y actually did go down at $T+1$. However, if we inspect this case in greater detail, at the moment $T+1$, the value of the underlying asset which y represents, went up at $T+1$, as the return in that period is still positive, however the model actually implicitly predicted the price to go down. This is a nuance which MDA doesn't capture, but MES does. This metric makes sense if the values of the underlying time-series are centered around 0 (mean = 0), which returns in financial markets are in the short-term. However, this metric can be generalized to for time-series of any mean μ if we define it as $\text{MES} = \text{mean}\{\text{sign}(y_{\text{predicted}}^1 \times (y_{\text{actual}}^1 - \mu)), \text{sign}(y_{\text{predicted}}^2 \times$

$$(y_{\text{actual}}^2 - \mu), \dots, \text{sign}(y_{\text{predicted}}^n \times (y_{\text{actual}}^n - \mu))\}$$

D.2 Data

CHAPS data

CHAPS data source:

<https://www.ons.gov.uk/economy/economicoutputandproductivity/output/datasets/ukspendingoncredit>

Data sourcing

Exceptions:

- Daily CHAPS data not available 2018-01-01 to 2020-01-01.
- Weekly CHAPS data not available 2017-01-01 to 2022-01-01 and 2015-01-01 to 2020-01-01.
- Weekly CHAPS data used was available only 2020-01-01 to 2024-01-01.
- Monthly CHAPS and Exchange rate data was not available in sufficient quantity to conduct experiment.

Partial autocorrelation

Partial autocorrelation measures the correlation between a time series and its lagged values, while controlling for the influence of intermediate lags. In contrast, autocorrelation simply measures the correlation between the time series and its lagged versions without accounting for the effects of other lags. Partial autocorrelation provides a clearer picture of the direct relationship between observations and their lagged counterparts by isolating the specific influence of each lag. This makes it particularly useful for detecting seasonality in time-series data because it helps identify the direct impact of observations at seasonal intervals. By highlighting the most relevant lags for seasonality, partial autocorrelation allows for a more precise analysis of recurring patterns, which is harder to achieve using autocorrelation alone.

D.3 Experiment design

ARIMA

Originally, this package is for R. I used the equivalent Python implementation:
<https://github.com/alkaline-ml/pmdarima>.

Naive Simple Autoregressor (NSA)

Given a training time-series sequence of length n : $T = \{y_1, y_2, \dots, y_n\}$, and a prediction horizon of length h , the prediction of NSA will be: $P = \{y_n \text{ repeated } h \text{ times}\}$, i.e. NSA predicts the future value(s) the same as the last value.

Loss

Loss in the context of fine-tuning is a measure of model's performance on the current epoch it has trained on. Lower the loss, the better the model's performance in that epoch.

Benchmark