

Harnessing LLMs for Temporal Data - A Study on Explainable Financial Time Series Forecasting

Xinli Yu, Zheng Chen, Yanbin Lu

{xinliy, zgchen, luyanbin}@amazon.com

Abstract

Applying machine learning to financial time series has been an active area of industrial research enabling innovation in market insights, risk management, strategic decision-making, and policy formation. This paper explores the novel use of Large Language Models (LLMs) for explainable financial time series forecasting, addressing challenges in cross-sequence reasoning, multi-modal data integration, and result interpretation that are inherent in traditional approaches. Focusing on NASDAQ-100 stocks, we utilize public historical stock data, company metadata, and economic/financial news. Our experiments employ GPT-4 for zero-shot/few-shot inference and Open LLaMA for instruction-based fine-tuning. The study demonstrates LLMs' ability to generate well-reasoned decisions by leveraging cross-sequence information and extracting insights from text and price time series. We show that our LLM-based approach outperforms classic ARMA-GARCH and gradient-boosting tree models. Furthermore, fine-tuned public LLMs, such as Open-LLaMA, can generate reasonable and explainable forecasts, although they underperform compared to GPT-4.

1 Introduction

The rapid advancements in Machine Learning (ML) and Artificial Intelligence (AI) technologies over the past few years have opened up numerous opportunities and challenges across various domains, including the realm of financial markets (Kovalerchuk and Vityaev, 2005; Bahrammirzaee, 2010; Qi and Xiao, 2018). In particular, the task of financial time series forecasting, a key element in industrial risk management, market insights, strategic decision-making and policy formation, has witnessed significant technological innovations, from statistical/econometric time series techniques (Härdle et al., 1997; Andersen et al., 2009; Chen et al., 2011; Patton, 2012), to machine learning techniques (Kim, 2003; Yoo et al., 2005; Krollner et al.,

2010), to deep learning (Dingli and Fournier, 2017; Júnior and Nievola, 2018; Sezer et al., 2020; Leung and Zhao, 2021; Lara-Benítez et al., 2021). Despite these advancements, there are several inherent challenges associated with the deployment of ML/AI models in finance.

One challenge lies in the realm of **cross-sequence reasoning and inference**, a vital aspect for understanding temporal patterns and making accurate predictions. The current approaches include time-series correlation analysis (Plerou et al., 1999; Gopikrishnan et al., 2000; Conlon et al., 2009; Chen et al., 2018) and clustering (Rani and Sikka, 2012; Babu et al., 2012; Aghabozorgi et al., 2015). Deep learning has recently been leveraged to learn from the complex latent dependencies among time series (Hua et al., 2019; Maulik et al., 2020; Song and Fujimura, 2021; Nguyen and Quanz, 2021). Despite these advancements, existing methods have yet to effectively capture the intricate dependencies characteristic of time series data. The varying design, implementation, and data requirements of these methods further creates a barrier for their widespread application in the field.

Another notable hurdle involves handling **complex multi-modal financial temporal data** that extends beyond numeric sequences. The data may encapsulate diverse sources such as historical news, financial knowledge graphs, social media activities, and various other market indicators. There has been recent effort leveraging *statistical inference* (Kanungasukkasem and Leelanupab, 2019), RNN/CNN with text embedding (Vargas et al., 2017), *graph neural networks* (Cheng et al., 2022), etc. to integrate the complex information.

Last but of utmost importance, the issue of **interpretability and explainability** poses significant challenges to the trustworthiness of machine learning and deep learning models. The majority of existing deep learning models operate as black boxes, offering little insight into their decision-making

processes. This lack of transparency sometimes raises concerns about the result reliability and impedes user trust. This is particularly relevant in sensitive fields like finance, where substantial investments and assets are at stake. There is recent study trying to understand deep-learning based predictions through attention scores (Hsieh et al., 2021), but such insight is still not readily human readable and still requires considerable interpretation effort.

The recent advancement of *Large Language Models* (LLMs) (Brown et al., 2020a; Touvron et al., 2023b; Brown et al., 2020b; OpenAI, 2023a) potentially lend us a powerful tool to address all above challenges in a unified, flexible way.

First, **LLMs can learn complex relations among sequences**. LLMs are the most powerful Transformer-based models, and there has been abundant researches showing Transformer-based models capable of learning the underlying complex relations among textual sequences (Yun et al., 2019; Rong et al., 2020; Zhang et al., 2020; Dwivedi and Bresson, 2020; Ying et al., 2021) and solving quantitative problems (Wei et al., 2022; Lewkowycz et al., 2022; Imani et al., 2023). It is reasonable to expect the potential of LLMs understanding complex dependencies among numeric time series augmented by temporal textual sequences.

Secondly, **LLMs have demonstrated outstanding reasoning and inference capability over multi-modal data**. By design, LLMs are proficient at learning from a broad spectrum of data sources and types. They are trained on a vast amount of texts from the internet, encompassing a wide range of topics, styles, and formats. This equips them to handle diverse input data, such as numerical, textual, structured data (Wu et al., 2023; Shen et al., 2023). This multi-modal data handling capability could be particularly useful for financial forecasting, where crucial information often comes from disparate sources, such as numerical market data, textual news articles, and social media posts.

Lastly, **LLMs are natural explainers that generate human readable explanations providing insight into a decision**. One of the key advantages of LLMs is their ability to generate natural language text that is coherent, contextual, and comprehensive. This allows them to provide human-readable explanations for their decisions (Zhao et al., 2023). Furthermore, through Chain-of-Thoughts (COT) or step-by-step thinking (Wei et al., 2022; Zhang et al., 2023; Lightman et al., 2023), beyond a few

sentences of explanation, LLMs can even generate detailed step-by-step reasoning to reveal the decision-making process.

The following summarizes the main contributions of this paper,

- This paper takes a novel exploration to study LLMs' potential to the valuable task of explainable financial time series forecasting. For this paper, we focus on the NASDAQ-100 stock price time series. To the best of our knowledge, there is not yet public studies on this topic to date.
- We experiment with a combination of zero-shot/few-shot inference techniques with the state-of-the-art AI model GPT-4 (OpenAI, 2023a), and instruction-based fine-tuning using Open LLaMA (Geng and Liu, 2023). Our experiment results also show that the technique of chain-of-thoughts helps boost the performance in most of the experiments.
- We compare our proposed LLM approaches with existing methods, including an ARMA-GARCH model and a gradient-boosting tree model. We show even zero-shot inference using GPT-4 can outperform a boosting-tree model with about ~ 300 features.

2 Related Works

The field of financial time series forecasting has been a subject of extensive research, with various methodologies being proposed over the years.

2.1 Traditional Statistical/Econometric Methods

Traditional statistical/econometric methods have long been the cornerstone of financial time series forecasting. Techniques such as ARMA-GARCH models have been widely used due to their ability to capture dependencies and volatility clustering in financial time series (Drost and Nijman, 1993; Francq and Zakoian, 2004; Andersen et al., 2009; Henneke et al., 2011). These models have been extended and modified in various ways to better capture the complexities of financial markets (Tang et al., 2003; Ghahramani and Thavaneswaran, 2006; Hossain and Nasser, 2011; Ma and Yu, 2013). Other popular statistical/econometric methods for financial time series include Vector Autoregressive Models (VAM) (Zivot and Wang, 2006), State-Space Models and the Kalman Filter (De Jong and

Zehnwirth, 1983), Diffusion Models (Fan, 2005), Vector Error Correction Model (VECM) (Johansen, 1995), Dynamic Stochastic General Equilibrium (DSGE) (Smets and Wouters, 2003), etc.

2.2 Machine Learning Techniques

With the advent of machine learning, a variety of models have been applied to financial forecasting. Decision trees, support vector machines, etc., have been actively studied for financial time series prediction (Trafalis and Ince, 2000; Yang et al., 2002; Pai and Lin, 2005; Wang and Chan, 2006; Tsai and Wang, 2009; Li and Liao, 2017). More recently, deep learning techniques, such as Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs), and Transformer models, have been applied to this task, demonstrating their ability to capture complex, non-linear relationships in the data (Dingli and Fournier, 2017; Júnior and Nievola, 2018; Sezer et al., 2020; Leung and Zhao, 2021; Lara-Benítez et al., 2021).

2.3 Large Language Models

The recent development of Large Language Models (LLMs) has opened up new possibilities for financial time series forecasting. LLMs, such as GPT-3 (Brown et al., 2020b) and GPT-4 (OpenAI, 2023a), LLaMA (Touvron et al., 2023a) (including Alpaca (Taori et al., 2023), Vincuna (Chiang et al., 2023)), have demonstrated remarkable capabilities in reasoning and understanding complex dependencies in the heterogeneous data, and the ability to generate human-readable explanations for their decisions (Zhao et al., 2023; Wei et al., 2022; Zhang et al., 2023; Lightman et al., 2023). However, the application of LLMs in financial time series forecasting with explanation is still a relatively unexplored area, and this paper aims to contribute to this emerging field.

3 Methodology

For this paper, we study the NASDAQ-100 stock price time series, supplemented by metadata about the stock company and relevant financial news data concerning both the specific stock and the broader financial/economic landscape. Our primary focus is on forecasting weekly/monthly *stock returns* (defined as the percentage change in stock price from the beginning to the end of the week/month) with accompanying explanations¹.

¹Weekly/monthly stock return forecast aligns better with LLMs' expertise at strategic reasoning and decision-making.

We demonstrate our structured design of prompts for LLMs and apply the state-of-the-art GPT-4 model (OpenAI, 2023b) for zero-shot and few-shot inference tasks. For fine-tuning, we utilize the publicly available Open LLaMA (Geng and Liu, 2023). We also incorporate the technique of Chain of Thoughts (COT) (Wei et al., 2022; Lightman et al., 2023), which has been found to enhance the effectiveness of LLMs in other research studies.

3.1 Data

3.1.1 Stock Price Data

We download daily NASDAQ-100 stock price data from Yahoo Finance². We first normalize the numeric price time series as a percentage-change time series, and then categorize them into bins. For example, for weekly forecasting, we categorize the price change between this week and last week into 12 bins "D5+", "D5", "D4", "D3", "D2", "D1", "U1", "U2", "U3", "U4", "U5", "U5+", where "D5+" means price dropping more than 5%, "D i " ($i=5,4,3,2,1$) means price dropping between $(i-1)\%$ and $i\%$, "U5+" means price rising more than 5%, "U i " ($i=1,2,3,4,5$) means price rising between $(i-1)\%$ and $i\%$. The number of bins might vary at different granularity. For example, for monthly inference, we allow i be up to 10, and there is corresponding "D10+", "U10+" bins.

3.1.2 Company Profile Data

We use GTP-4 to generate company description, general positive/negative factors that might impact the company's stock price. See Appendix Figure 1 for an example of the prompt to ask GPT-4 to generate the company profile, and the GPT-4 response.

3.1.3 Finance/Economy News Data

We use *Google Custom Search API* to obtain stock top-5 news stories on a weekly basis for each NASDAQ-100 stock. After that, we use GPT-4 to generate a summary and extract keywords from each obtained news article. An example of prompt and GPT-4 response is shown in Appendix Figure 2. A similar method is applied to obtain weekly top-5 news stories about macro economy and finance.

To reduce input size, We further generate *meta summary & keywords* for each week using GPT-4, given all the top story summaries and keywords of the week. An example of the meta summary &

It is also much less demanding on LLMs' latency.

²Using the stock list of Jun 01, 2023, and download using yfinance package pypi.org/project/yfinance/

keywords is shown in Appendix Figure 3. They look similar to the example in Appendix Figure 2, but much condensed. We use the meta summary & keywords for further experiments and evaluation.

3.2 Instruction-Based Zero-shot/Few-shot Inference with LLMs

In zero-shot and few-shot inference, LLMs demonstrate their ability to generate responses either without any additional examples (*zero-shot*) or based on a minimal number of examples beyond the original training set (*few-shot*). In our zero-shot/few-shot inference experiment, we utilize an instruction-based prompt. The structure of our prompt includes instructions, the company profile, a historical temporal news summary/keywords sequence intermixed with the categorized stock price time series, and *cross-sequence few-shot learning examples* (see Appendix Figure 4 for example).

- We provide few-shot learning examples from stocks similar to the subject of interest. This design multi-purposes the few-shot examples to enable the LLM **consider cross-sequence information from other stocks**. To identify similar stocks, we query GPT-4 with a query "List top 3 NASDAQ stocks most similar to AAPL". A typical response is like "MSFT, GOOGL, AMZN"³. Here we in fact implicitly **leverage LLM inherent knowledge of financial entities and concepts**.
- There are other tweaks to the prompt structure. For instance, we divided the instruction into two parts, positioning them at the beginning and end of the prompt. This aids the model in better recognizing the task: to predict next week's summary & keywords, rather than summarizing historical data. The predicted summary & keywords serve as the explanation for the stock return prediction.

We also experimented the *Chain-of-Thoughts* approach (Wei et al., 2022; Zhang et al., 2023; Lightman et al., 2023), i.e., the idea of "**step-by-step thinking**", by appending the instruction "Can you reason step by step before finalizing the output?" to the end of the prompt. To our surprise, this notably improved the performance by a few points

³We manually checked the "similar-stocks" responses by GPT-4 for NASDAQ-100 stocks and can confirm the results are logical; no hallucination was observed. The 09/2021 knowledge cut-off does not pose a concern, as all NASDAQ-100 stocks have established for some time prior to this date.

(see Section 4.2). The result of the step-by-step thinking process in response to Appendix Figure 4 is illustrated in Appendix Figure 5, where it is evident that GPT-4 identifies a previously overlooked crucial point about "earnings reports" when explicit reasoning steps are generated.

3.3 Instruction-based Fine-tuning with Open LLaMA

We perform instruction-based fine-tuning using Open LLaMA 13B model to see how well a publicly available model could perform in comparison to GPT-4, especially after fine-tuning. The Open LLaMA 13B model, in its zero-shot inference, typically tends to replicate portions of the prompt rather than executing the prompt instructions effectively. Therefore, it is incapable of properly handling instruction-based prompts as shown in Appendix Figure 4 without undergoing a process of fine-tuning. Therefore we focus on fine-tuning with the Open LLaMA model in this paper.

Instruction-based fine-tuning has been recently shown to be effective in guiding the model's training process with specific directives (Taori et al., 2023; Peng et al., 2023). We created a dataset of 30K weekly forecasting plus 7K monthly forecasting, derived from 5-year historical data spanning from Jun 2017 to June 2022. Unlike GPT-4 that supports up to 8K token size, we need to compress the prompt into 1K tokens for fine-tuning Open LLaMA, due to model and hardware constraints. For each fine-tuning example, we employ GPT-4 to condense the full historical meta news summary/keywords (e.g. from week 8 to the last week as shown in Appendix Figure 4) into a single, even more concise summary/keywords pair. Simultaneously, the "Company Profile" and "Forecasting Examples" sections of the prompt are also respectively condensed into more succinct summary paragraphs.

While it would be ideal for Open LLaMA to manage its own end-to-end experiment, including the task of prompt compression for fine-tuning, we still resort to using GPT-4 right now. This is due to Open LLaMA 13B model's zero-shot summarization capability is considerably inferior to those of GPT-4 in practice. The summaries and keywords extracted by Open LLaMA 13B model often fall short of usability.

Once fine-tuned, the Open LLaMA 13B model demonstrates a much more satisfactory comprehen-

sion of the instruction, resulting in the generation of a forecast and an accompanying explanation that appears coherent. This is illustrated in Appendix Figure 6. As per the result in section 4.2, when it comes to binary classification, the Open LLaMA model’s performance is competitive compared to GPT-4. However, we’ve noticed that the Open LLaMA model has a tendency to produce more extreme predictions, such as U5+ or D5+, which result in a relatively higher squared error.

4 Experiments

4.1 Experiment Setup

4.1.1 Data Time Window

The details of the data used in the experiments is as described in Section 3.1. We focus on NASDAQ-100 stock return forecasting for this paper.

- The *training/fine-tuning time window* contains 5-year data from 06/12/2017 to 06/05/2022. There are in total 26K data examples in this time window (260 weeks and 100 stocks), and they are used for training the baseline models, and fine-tuning the Open LLaMA 13B model.
- The *evaluation time window* has 52 weeks spanning from 06/06/2022 to 06/04/2023. The evaluation of baseline models, the zero/few-shot inference with GPT-4, and the evaluation of fine-tuned Open LLaMA 13B model, are based on the 5.2K evaluation examples (52 weeks, 100 stocks) in this time window.

4.1.2 Baseline Models

To evaluate the performance of our approach, we include a heuristic baseline using the most-frequent historical bin (i.e. the most frequent bin from historical weeks before the target week to forecast) as the prediction, an ARMA-GARCH model ($p = q = 1$) (Tang et al., 2003; Ma and Yu, 2013), and a gradient-boosting tree model (Natekin and Knoll, 2013) implemented by LightGBM package (Ke et al., 2017). These baseline models are trained on the training/fine-tuning data time window, and evaluated on the evaluation time window.

For the gradient-boosting tree model, we include the following features. There are total about 300 features for the tree.

1. Historical price time series available in the daily stock price data, including open, close, min, max prices, and the daily trading volume.

2. The average, medium, min, max, and stddev of a rolling window of size 2, 5, 10, 30, 60, 90 for the above time series.
3. The stock sector information and historical earnings are obtained from Alpha Vantage⁴.

4.1.3 Evaluation Metrics

We perform weekly and monthly stock return forecasting with the baselines and LLM-based methods. We treat 4 weeks as one month for convenience, and therefore there are 13 "month"s in the 52-week evaluation time window.

To evaluate the performance of our forecasting models, we employ three metrics.

- Binary precision assesses the model’s ability to correctly predict the general direction of stock price movement, i.e., "Up" (U) or "Down" (D).
- Bin precision, on the other hand, evaluates the model’s accuracy in predicting the exact bin from a full list of bins such as "D5+", "D5", "D4", ..., "D1", "U1", ..., "U5", "U5+".
- The MSE of consecutive bin ordinals (e.g., -6 for "D5+", -5 for "D5", ..., 0 for "U1", ..., 4 for "U5", 5 for "U5+") is used to measure the average squared differences between the model’s predictions and the actual values. This metric helps to understand the model’s tendency to make extreme forecasts when its predictions are incorrect.

To evaluate the quality of the forecasting explanation (the predicted next-week/month summary/keywords), we employ ROGUE-1 and ROGUE-2 scores to compare with the actual summary/keywords by GPT-4 extracted from the actual top news of the next week/month.

4.2 Performance Evaluation

Our experiment results are summarized in Table 1 and 2. Table 1 provides a comparative analysis of our LLM-based methods and the baseline models in terms of their performance in forecasting stock returns. Table 2, on the other hand, evaluates the quality of the explanations generated by the LLMs.

In summary, our results show the effectiveness of LLMs in financial time series forecasting, with "GPT-4 few-shot with COT" consistently showing

⁴<https://www.alphavantage.co/documentation/>

	Weekly			Monthly (Every 4 Weeks)		
	Binary Precision	Bin Precision	MSE	Binary Precision	Bin Precision	MSE
Most-Frequent Historical Bin	50.7%	16.4%	43.5	51.4%	17.2%	155.1
ARMA-GARCH	52.4%	11.1%	22.1	50.5%	6.2%	90.1
Gradient Boosting Tree Model	60.8%	26.4%	24.3	56.4%	17.7%	85.6
GPT-4 Zero-Shot	64.5%	31.2%	20.5	64.8%	26.0%	60.1
GPT-4 Few-Shot	65.8%	32.7%	20.6	65.3%	26.5%	58.2
GPT-4 Few-Shot w/ COT	66.5%	35.2%	18.7	69.5%	28.6%	50.4
Open LLaMA (13B) Fine-Tuned	62.2%	26.5%	23.3	60.1%	22.6%	63.3
Open LLaMA (13B) Fine-Tuned w/ COT	64.7%	30.7%	21.0	62.2%	24.4%	63.5

Table 1: Performance comparison between the baseline models and LLMs for stock price weekly/monthly forecast.

	Weekly				Monthly (Every 4 Weeks)			
	ROUGE-1 (S)	ROUGE-2 (S)	ROUGE-1 (K)	ROUGE-2 (K)	ROUGE-1 (S)	ROUGE-2 (S)	ROUGE-1 (K)	ROUGE-2 (K)
GPT-4 Zero-Shot	0.2212	0.0675	0.1295	0.0447	0.2528	0.0665	0.1335	0.0657
GPT-4 Few-Shot	0.2242	0.0526	0.1304	0.0454	0.2450	0.0634	0.1348	0.0644
GPT-4 Few-Shot w/ COT	0.2414	0.0543	0.2083	0.0869	0.2645	0.0758	0.2450	0.1025
Open LLaMA (13B) Fine-Tuned	0.2053	0.0395	0.0927	0.0324	0.2242	0.0474	0.1167	0.0520
Open LLaMA (13B) Fine-Tuned w/ COT	0.2371	0.0434	0.1123	0.0425	0.2436	0.0536	0.1356	0.0834

Table 2: Explanation quality evaluation using ROGUE scores, using the GPT-4 summary/keyword extraction of each week’s true top news from google search as the ground truth.

the best performance in both prediction accuracy and explanation quality. The results also highlight the technique of Chain-of-Thoughts (COT) consistently boosts performance, and the potential of instruction-based fine-tuning with publicly available LLMs like Open LLaMA to achieve reasonable performance in comparison to GPT-4 through fine-tuning with COT.

4.2.1 Stock Price Forecasting

From the results of Table 1, we observe that both GPT-4 and Open LLaMA 13B model outperform the ARMA-GARCH model and the gradient-boosting tree model in terms of both binary and bin precision. GPT-4, in particular, shows superior performance in both zero-shot and few-shot settings, with the few-shot setting with COT achieving the best performance. In terms of MSE, "GPT-4 few-shot with COT" also achieves the lowest error, indicating that it not only best predicts the direction of the price change but also provides a more accurate estimate of the magnitude of the change.

Open LLaMA 13B model, after fine-tuning, shows competitive performance compared to GPT-4 in terms of binary precision. However, its bin precision is obviously worse, indicating it lacks competitive fine-grained reasoning capability to pick the right bin. It also tends to produce more extreme predictions, resulting in a higher MSE.

4.2.2 Explanation Quality

Table 2 shows the quality of the explanations generated by the LLMs (GPT-4 and fine-tuned Open LLaMA), evaluated using ROUGE-1 and ROUGE-

2 scores for both the summary (S) and keywords (K) of the news. Again, the results show that "GPT-4 few-shot with COT" achieves the highest ROUGE scores, indicating that it generates the most relevant and accurate explanations for the predictions. Open LLaMA, after fine-tuning with COT, also shows reasonable explanation quality in parallel with GPT-4 results without COT.

5 Conclusion

In this study, we explored using Large Language Models (LLMs) to tackle inherent challenges like cross-sequence reasoning, multi-modal signals integration, and result interpretability in financial time series forecasting. In particular, we experimented GPT-4 and Open LLaMA for the NASDAQ-100 stock return predictions. With structured prompts comprising company profile, historical stock price, and financial news data, LLMs generated human understandable explanations and forecasts. The performance of these LLMs surpassed traditional models like ARMA-GARCH and gradient-boosting trees, especially when integrating a step-by-step reasoning process based on the Chain of Thought (COT) approach. Furthermore, our fine-tuning experiments highlighted the viability of tuning a publicly available LLM to also achieve reasonable performance in comparison to GPT-4.

The preliminary results of applying LLMs in explainable financial forecasting are encouraging. This is the first step to develop a LLM-based explainable financial forecast system to assist business decision-making. We envision a future

where financial forecasting is not only more precise but also more comprehensible and transparent, thus transforming financial and business decision-making across the sector.

Limitations

While we present promising initial results for the LLM-based approach for explainable financial time series based on NASDAQ-100 stock returns, the general applicability of our approach to different types of temporal data remains a question of future investigation.

- In a narrower context, the effectiveness of our approach when applied to other stock indices like the S&P 500 or Russell 2000 is yet to be validated. Each of these indices harbors distinct characteristics and diverse company compositions, which may influence the performance of our method.
- In a wider context, the potential of our method to forecast other types of financial temporal data remains unexplored. This includes internal temporal time series such as return-on-investment (ROI), sales, headcounts, and costs from various departments, augmented by related internal documents. It also extends to other public time series data such as company earnings, housing prices, security prices, and futures prices.

The quality of explanations generated by our method is currently evaluated using automatic text similarity metrics, namely ROGUE and BLEU scores. These scores compare predicted explanations to ground-truth news summaries and keywords. However, we have yet to thoroughly examine other aspects of these explanations, including the possibility of hallucination or the presence of specific patterns in explanations associated with failed time series forecasts.

Other signals can be considered for future investigations, such as macro economy time series (Stock and Watson, 1999) and social media data (Mankar et al., 2018; Javed Awan et al., 2021).

Ethics Statement

In conducting our research, we committed to transparency in our methodology, results reporting, and data usage. Our work adhered to the guidelines and principles of responsible AI research. All data

used in this paper is public data or can be obtained through publicly available APIs. No user data is involved in the research.

References

- Saeed Aghabozorgi, Ali Seyed Shirkhorshidi, and Teh Ying Wah. 2015. Time-series clustering—a decade review. *Information systems* 53 (2015), 16–38.
- Torben Gustav Andersen, Richard A Davis, Jens-Peter Kreiß, and Thomas V Mikosch. 2009. *Handbook of financial time series*. Springer Science & Business Media.
- M Suresh Babu, N Geethanjali, and B Satyanarayana. 2012. Clustering approach to stock market prediction. *International Journal of Advanced Networking and Applications* 3, 4 (2012), 1281.
- Arash Bahrammirzaee. 2010. A comparative survey of artificial intelligence applications in finance: artificial neural networks, expert system and hybrid intelligent systems. *Neural Computing and Applications* 19, 8 (2010), 1165–1195.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020b. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020a. Language Models are Few-Shot Learners. *ArXiv abs/2005.14165* (2020).
- Cathy WS Chen, Feng-Chi Liu, and Mike KP So. 2011. A review of threshold time series models in finance. *Statistics and its Interface* 4, 2 (2011), 167–181.
- Zheng Chen, Xinli Yu, Yuan Ling, Bo Song, Wei Quan, Xiaohua Hu, and Erjia Yan. 2018. Correlated anomaly detection from large streaming data. In *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 982–992.
- Dawei Cheng, Fangzhou Yang, Sheng Xiang, and Jin Liu. 2022. Financial time series forecasting with multi-modality graph neural network. *Pattern Recognition* 121 (2022), 108218.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan

- Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* Chat-GPT Quality. <https://vicuna.lmsys.org>
- Thomas Conlon, Heather J Ruskin, and Martin Crane. 2009. Cross-correlation dynamics in financial time series. *Physica A: Statistical Mechanics and its Applications* 388, 5 (2009), 705–714.
- Piet De Jong and Ben Zehnwirth. 1983. Claims reserving, state-space models and the Kalman filter. *Journal of the Institute of Actuaries* 110, 1 (1983), 157–181.
- Alexiei Dingli and Karl Sant Fournier. 2017. Financial time series forecasting-a deep learning approach. *International Journal of Machine Learning and Computing* 7, 5 (2017), 118–122.
- Feike C Drost and Theo E Nijman. 1993. Temporal aggregation of GARCH processes. *Econometrica: Journal of the Econometric Society* (1993), 909–927.
- Vijay Prakash Dwivedi and Xavier Bresson. 2020. A generalization of transformer networks to graphs. *arXiv preprint arXiv:2012.09699* (2020).
- Jianqing Fan. 2005. A selective overview of nonparametric methods in financial econometrics. *Statist. Sci.* (2005), 317–337.
- Christian Francq and Jean-Michel Zakoian. 2004. Maximum likelihood estimation of pure GARCH and ARMA-GARCH processes. *Bernoulli* 10, 4 (2004), 605–637.
- Xinyang Geng and Hao Liu. 2023. *OpenLLaMA: An Open Reproduction of LLaMA*. https://github.com/openlm-research/open_llama
- M Ghahramani and A Thavaneswaran. 2006. Financial applications of ARMA models with GARCH errors. *The Journal of Risk Finance* 7, 5 (2006), 525–543.
- Parameswaran Gopikrishnan, Vasiliki Plerou, Yan Liu, LA Nunes Amaral, Xavier Gabaix, and H Eugene Stanley. 2000. Scaling and correlation in financial time series. *Physica A: Statistical Mechanics and its Applications* 287, 3-4 (2000), 362–373.
- Wolfgang Härdle, Helmut Lütkepohl, and Rong Chen. 1997. A review of nonparametric time series analysis. *International statistical review* 65, 1 (1997), 49–72.
- Jan S Henneke, Svetlozar T Rachev, Frank J Fabozzi, and Metodi Nikolov. 2011. MCMC-based estimation of Markov Switching ARMA–GARCH models. *Applied Economics* 43, 3 (2011), 259–271.
- Altat Hossain and Mohammed Nasser. 2011. Comparison of the finite mixture of ARMA-GARCH, back propagation neural networks and support-vector machines in forecasting financial returns. *Journal of Applied Statistics* 38, 3 (2011), 533–551.
- Tsung-Yu Hsieh, Suhang Wang, Yiwei Sun, and Vasant Honavar. 2021. Explainable multivariate time series classification: a deep neural network which learns to attend to important variables as well as time intervals. In *Proceedings of the 14th ACM international conference on web search and data mining*. 607–615.
- Yuxiu Hua, Zhifeng Zhao, Rongpeng Li, Xianfu Chen, Zhiming Liu, and Honggang Zhang. 2019. Deep learning with long short-term memory for time series prediction. *IEEE Communications Magazine* 57, 6 (2019), 114–119.
- Shima Imani, Liang Du, and Harsh Shrivastava. 2023. Mathprompter: Mathematical reasoning using large language models. *arXiv preprint arXiv:2303.05398* (2023).
- Mazhar Javed Awan, Mohd Shafry Mohd Rahim, Haitham Nobanee, Ashna Munawar, Awais Yasin, and Azlan Mohd Zain. 2021. Social media and stock market prediction: a big data approach. *MJ Awan, M. Shafry, H. Nobanee, A. Munawar, A. Yasin et al., "Social media and stock market prediction: a big data approach," Computers, Materials & Continua* 67, 2 (2021), 2569–2583.
- Søren Johansen. 1995. *Likelihood-based inference in cointegrated vector autoregressive models*. OUP Oxford.
- Norberto Ritzmann Júnior and Julio Cesar Nievola. 2018. A generalized financial time series forecasting model based on automatic feature engineering using genetic algorithms and support vector machine. In *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 1–8.
- Nont Kanungsukkasem and Teerapong Leelanupab. 2019. Financial latent Dirichlet allocation (FinLDA): Feature extraction in text and data mining for financial time series prediction. *IEEE Access* 7 (2019), 71645–71664.
- Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. 2017. Lightgbm: A highly efficient gradient boosting decision tree. *Advances in neural information processing systems* 30 (2017).
- Kyoung-jae Kim. 2003. Financial time series forecasting using support vector machines. *Neurocomputing* 55, 1-2 (2003), 307–319.
- Boris Kovalerchuk and Evgenii Vityaev. 2005. Data mining for financial applications. *Data Mining and Knowledge Discovery Handbook* (2005), 1203–1224.
- Bjoern Krollner, Bruce J Vanstone, Gavin R Finnie, et al. 2010. Financial time series forecasting with machine learning techniques: a survey.. In *ESANN*.
- Pedro Lara-Benítez, Manuel Carranza-García, and José C Riquelme. 2021. An experimental review

- on deep learning architectures for time series forecasting. *International Journal of Neural Systems* 31, 03 (2021), 2130001.
- Tim Leung and Theodore Zhao. 2021. Financial time series analysis and forecasting with Hilbert–Huang transform feature generation and machine learning. *Applied Stochastic Models in Business and Industry* 37, 6 (2021), 993–1016.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. 2022. Solving quantitative reasoning problems with language models. *arXiv preprint arXiv:2206.14858* (2022).
- Wei Li and Jian Liao. 2017. A comparative study on trend forecasting approach for stock price time series. In *2017 11th IEEE International Conference on Anti-counterfeiting, Security, and Identification (ASID)*. IEEE, 74–78.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s Verify Step by Step. *arXiv preprint arXiv:2305.20050* (2023).
- Junmei Ma and Xinli Yu. 2013. Research on the Pricing of the Basket Credit Default Swap. *Mathematical Computation* 2, 4 (2013).
- Tejas Mankar, Tushar Hotchandani, Manish Madhwani, Akshay Chidrawar, and CS Lifna. 2018. Stock market prediction based on social sentiments using machine learning. In *2018 international conference on smart city and emerging technology (ICSCET)*. IEEE, 1–3.
- Romit Maulik, Arvind Mohan, Bethany Lusch, Sandeep Madireddy, Prasanna Balaprakash, and Daniel Livescu. 2020. Time-series learning of latent-space dynamics for reduced-order model closure. *Physica D: Nonlinear Phenomena* 405 (2020), 132368.
- Alexey Natekin and Alois Knoll. 2013. Gradient boosting machines, a tutorial. *Frontiers in neurorobotics* 7 (2013), 21.
- Nam Nguyen and Brian Quanz. 2021. Temporal latent auto-encoder: A method for probabilistic multivariate time series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 9117–9125.
- OpenAI. 2023a. GPT-4 Technical Report. *ArXiv abs/2303.08774* (2023).
- OpenAI. 2023b. GPT-4 Technical Report. *arXiv:2303.08774 [cs.CL]*
- Ping-Feng Pai and Chih-Sheng Lin. 2005. A hybrid ARIMA and support vector machines model in stock price forecasting. *Omega* 33, 6 (2005), 497–505.
- Andrew J Patton. 2012. A review of copula models for economic time series. *Journal of Multivariate Analysis* 110 (2012), 4–18.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277* (2023).
- Vasiliki Plerou, Parameswaran Gopikrishnan, Bernd Rosenow, Luís A Nunes Amaral, and H Eugene Stanley. 1999. Universal and nonuniversal properties of cross correlations in financial time series. *Physical review letters* 83, 7 (1999), 1471.
- Yuan Qi and Jing Xiao. 2018. Fintech: AI powers financial services to improve people’s lives. *Commun. ACM* 61, 11 (2018), 65–69.
- Sangeeta Rani and Geeta Sikka. 2012. Recent techniques of clustering of time series data: a survey. *International Journal of Computer Applications* 52, 15 (2012).
- Yu Rong, Yatao Bian, Tingyang Xu, Weiyang Xie, Ying Wei, Wenbing Huang, and Junzhou Huang. 2020. Self-supervised graph transformer on large-scale molecular data. *Advances in Neural Information Processing Systems* 33 (2020), 12559–12571.
- Omer Berat Sezer, Mehmet Ugur Gudelek, and Ahmet Murat Ozbayoglu. 2020. Financial time series forecasting with deep learning: A systematic literature review: 2005–2019. *Applied soft computing* 90 (2020), 106181.
- Li Shen, Yan Sun, Zhiyuan Yu, Liang Ding, Xinmei Tian, and Dacheng Tao. 2023. On Efficient Training of Large-Scale Deep Learning Models: A Literature Review. *arXiv preprint arXiv:2304.03589* (2023).
- Frank Smets and Raf Wouters. 2003. An estimated dynamic stochastic general equilibrium model of the euro area. *Journal of the European economic association* 1, 5 (2003), 1123–1175.
- Wen Song and Shigeru Fujimura. 2021. Capturing combination patterns of long-and short-term dependencies in multivariate time series forecasting. *Neuro-computing* 464 (2021), 72–82.
- James H Stock and Mark W Watson. 1999. Business cycle fluctuations in US macroeconomic time series. *Handbook of macroeconomics* 1 (1999), 3–64.
- Him Tang, Kai-Chun Chiu, and Lei Xu. 2003. Finite mixture of ARMA-GARCH model for stock price prediction. In *Proceedings of the Third International Workshop on Computational Intelligence in Economics and Finance (CIEF’2003)*, North Carolina, USA. 1112–1119.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An Instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023b. LLaMA: Open and Efficient Foundation Language Models. *ArXiv abs/2302.13971* (2023).
- Theodore B Trafalis and Huseyin Ince. 2000. Support vector machine for regression and applications to financial forecasting. In *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks. IJCNN 2000. Neural Computing: New Challenges and Perspectives for the New Millennium*, Vol. 6. IEEE, 348–353.
- Chih F Tsai and Sammy P Wang. 2009. Stock price forecasting by hybrid machine learning techniques. In *Proceedings of the international multiconference of engineers and computer scientists*, Vol. 1. 60.
- Manuel R Vargas, Beatriz SLP De Lima, and Alexandre G Evsukoff. 2017. Deep learning for stock market prediction from financial news articles. In *2017 IEEE international conference on computational intelligence and virtual environments for measurement systems and applications (CIVEMSA)*. IEEE, 60–65.
- Jar-Long Wang and Shu-Hui Chan. 2006. Stock market trading rule discovery using two-layer bias decision tree. *Expert Systems with Applications* 30, 4 (2006), 605–611.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903* (2022).
- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kambaradur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564* (2023).
- Haiqin Yang, Laiwan Chan, and Irwin King. 2002. Support vector machine regression for volatile stock market prediction. In *Intelligent Data Engineering and Automated Learning—IDEAL 2002: Third International Conference Manchester, UK, August 12–14, 2002 Proceedings 3*. Springer, 391–396.
- Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. 2021. Do transformers really perform badly for graph representation? *Advances in Neural Information Processing Systems* 34 (2021), 28877–28888.
- Paul D Yoo, Maria H Kim, and Tony Jan. 2005. Machine learning techniques and use of event information for stock market prediction: A survey and evaluation. In *International Conference on Computational Intelligence for Modelling, Control and Automation and International Conference on Intelligent Agents, Web Technologies and Internet Commerce (CIMCA-IAWTIC'06)*, Vol. 2. IEEE, 835–841.
- Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J Kim. 2019. Graph transformer networks. *Advances in neural information processing systems* 32 (2019).
- Jiawei Zhang, Haopeng Zhang, Congying Xia, and Li Sun. 2020. Graph-bert: Only attention is needed for learning graph representations. *arXiv preprint arXiv:2001.05140* (2020).
- Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923* (2023).
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223* (2023).
- Eric Zivot and Jiahui Wang. 2006. Vector autoregressive models for multivariate time series. *Modeling financial time series with S-PLUS®* (2006), 385–429.

A Appendix - Example Prompts and Responses

Prompt: Generate a short description for stock AAPL's company "Apple Inc". Also list general positive and negative factors that might impact the stock price; be brief and use keywords. Consider diverse general factors, such as macro economic situation (e.g. inflation, CPI growth), business factors (e.g. sales, investment, products), technology factors (e.g. innovation), and others. Use format Description: ..., Positive Factors: ..., Negative factors: ...

GPT-4 Response:

Description:

AAPL (Apple Inc.) is a world-renowned technology company based in Cupertino, California. Known for its sleek designs and innovative technologies, the company produces a wide range of products including the iPhone, iPad, Mac, Apple Watch, and services like iCloud, App Store, and more.

Positive Factors:

Macroeconomic: Economic growth, increased consumer spending

Business: Strong sales, expansion into new markets, positive earnings reports, strategic acquisitions

Technology: Innovation, new product launches, advancements in artificial intelligence and augmented reality

Brand: Strong global brand, loyal customer base

Legal/Regulatory: Favorable policy changes, successful patent applications

Negative Factors:

Macroeconomic: Inflation, economic downturn, trade restrictions

Business: Decreased sales, poor earnings reports, lack of new investments

Technology: Slow pace of innovation, delays in product launches

Brand: Damage to brand reputation, loss of customer loyalty

Legal/Regulatory: Antitrust investigations, unfavorable policy changes, patent disputes

Environmental: Failure to meet sustainability goals, negative environmental impact

Figure 1: An example of a stock's company profile consisting of the company description, the general positive/negative factors impacting the company's stock prices.

Prompt [template, those inside {} are placeholders]: Please summarize the following noisy but possible news data extracted from web page HTML, and extract keywords of the news. The news text can be very noisy due to it is HTML extraction. Give formatted answer such as Summary: ..., Keywords: ... The news is supposed to be for {symbol} stock. You may put 'N/A' if the noisy text does not have relevant information to extract.

News: {news}

GPT-4 Response:

Summary: Apple Inc is postponing perks for some corporate divisions and expanding a cost-cutting effort, joining its peers in streamlining operations. The company proposes to forego bonuses or promotions for April, with the payments occurring only in October. Apple is also restricting hiring and leaving additional positions open after employee departure. This is triggering fear of downsizing among employees.

Keywords: Apple, bonus, hiring, cost-cutting, promotions, downsizing, employees

Figure 2: An example of news summary & keywords extracted from one news article for a stock's company (AAPL in this case), generated by GPT-4. The original news is at <https://sports.yahoo.com/apple-joins-cost-cut-bandwagon-145845685.html>. The prompt in this example is a template. Texts inside [] are comments and are not really in the prompt when we submit it to the LLM, and we will fill stock symbol and the news into the placeholders enclosed by {}.

Summary: Apple Inc.'s stock is displaying rising relative strength, although shy of a key benchmark, and the company has been highlighted as a top-performing stock due to its sales of 1.3 billion iPhones and a significant \$19.4 billion profit last quarter. As the NASDAQ surges with Apple as a big cap, Apple expands cost-cutting measures, postpones perks, restricts hiring, and leaves positions open, instigating downsizing concerns. The overall stock market performance is boosted by a \$30 billion deposit infusion for struggling firms.

Keywords: Apple Inc., stock, relative strength, cost-cutting, downsizing, NASDAQ, \$30 billion deposit infusion, iPhones, profit.

Figure 3: An example of one week's meta summary & keywords condensed from all the company's summaries and keywords from the week.

Prompt [template, those inside { } are placeholders]:

Instruction: Forecast next week stock return (price change) for symbol, given the company profile, historical weekly news summary, keywords, and stock returns, and optionally the examples from other stocks of a similar company. The trend is represented by bins "D5+", "D5", "D4", "D3", "D2", "D1", "U1", "U2", "U3", "U4", "U5", "U5+", where "D5+" means price dropping more than 5%, D5 means price dropping between 4% and 5%, "D4" means price dropping between 3% and 4%, "U5+" means price rising more than 5%, "U5" means price rising between 4% and 5%, "D4" means price rising between 3% and 4%, etc.

Company Profile: {company_profile} [company description, general positive/negative factors as shown in Appendix Figure 1.]

Recent News: News are ordered from oldest news to latest news.

====8 Weeks ago====

{meta_news_summary_company} [meta news summary & keywords of the week for the company]

{meta_news_summary_macro} [meta news summary & keywords for the macro economy and finance status of the week] ...

====7 weeks ago====

...

====Last week====

...

Forecasting Examples: {few_shot_learning_examples_from_similar_stocks} [We use forecasting examples from top-3 similar stocks as the few-short learning examples. The example is structured similar to this prompt, with company profile and recent news, but without the instruction texts.]

Now predict what could be the next week's Summary, Keywords, and forecast the Stock Return. The predicted Summary/Keywords should explain the stock return forecasting. You should predict what could happen next week. Do not just summarize the history. The next week stock return need not be the same as the previous week. Use format Summary: ..., Keywords: ..., Stock Return: ...

GPT-4 Response [Forecasting week 05/01/2023 to 05/07/2023]:

Summary: Apple's stock is expected to see fluctuations in the week to come against a backdrop of mixed news. The company's announced expansion into India, coupled with the introduction of high-yield savings accounts, initially provided a boost to the stock. However, growing concerns surrounding reduced PC shipments and apprehensions about the company's upcoming earnings injected a degree of uncertainty. Despite these challenges, Apple's unwavering focus on cost control and the anticipation surrounding innovative offerings, such as the forthcoming mixed-reality headset, shored up investor confidence. Consequently, the stock maintained a general upward trend throughout the week, albeit punctuated by a few dips. The broader economic climate on Wall Street also had a role to play, with a recovering sentiment in Q1 towards high-tech companies like Microsoft and Alphabet. These companies released pivotal financial reports and ground-breaking products during the period, influencing the tech sector's performance, potentially impacting Apple's stock trajectory.

Keywords: AAPL, mixed-reality headset, India expansion, high-yield savings account, cost control strategies, Bernstein, Wall Street, economic recovery, tech sector performance

Stock Return: U1

Ground Truth [Summary from actual news from 05/01/2023 to 05/07/2023 and the actual return]:

Summary: Apple's stock saw significant gains after outperforming Wall Street's fiscal second quarter targets, driven by robust iPhone sales and strong services, despite posting a decrease in revenue for the second consecutive quarter and facing potential competition in the VR headset market from Reliance Jio's JioDive. The company also noted an improvement in supply and strength in emerging markets, further boosting investor confidence.

Keywords: Apple, stock, Wall Street, fiscal second quarter, iPhone, services sales, quarterly earnings, revenue, declining revenue, emerging markets, improved supply, Virtual Reality, competition, JioDive, Reliance Jio.

Stock Return: U3

Figure 4: The prompt structure for the experiments in this paper with LLMs. We also give an example of GPT-4 response to a concrete prompt constructed from information on and before 04/30/2023. We notice the cross-sequence information and macro-economy information are obviously considered in the LLM's reasoning. The stock return forecast U1 is for the next week from 05/01/2023 to 05/07/2023. The actual market performance for AAPL is U3. We note Apple 2023 Q2 earning call happened on May 04 beat expectation, which maybe the major contributor to the higher-than-forecasting gain in the week.

Prompt: {The same prompt as in Appendix Figure 4.}
Can you reason step by step before the finalized output?

GPT-4 Response [Forecasting week 05/01/2023 to 05/07/2023]:

Given the information provided, we can see that Apple (AAPL) has had a number of positive developments and a few potential concerns in the recent past.

Positive Developments:

- Recent sales and earnings have been strong, and the company appears to be exploring new markets such as India.
- The upcoming mixed-reality headset and the "Apple Pay Later" service show that Apple is still innovating and pushing into new product areas.
- The cost-cutting measures and the plan to use 100% recycled cobalt show that Apple is taking steps to improve its efficiency and sustainability, which could be attractive to certain investors.
- Winning the lawsuit against Britain's antitrust regulator and avoiding potential penalties.
- AAPL stock has been trending upwards recently, with multiple weeks of upward movement (U5, U4, U3, U2), and only a single day of slight downward movement (D1).

Potential Concerns:

- Apple is cutting costs, delaying bonuses, and slowing down hiring in certain areas, which could potentially affect employee morale and productivity.
- The negative sentiment from analyst Toni Sacconaghi about the overestimation of Q3 and Q4 revenues by Wall Street could potentially impact the stock price.

In addition to above factors, the general direction of the stock market appears to be on an upward trajectory, as evidenced by the recent increases in the Dow Jones, S&P 500, and Nasdaq. In recent weeks, the stock performances of high-tech giants such as Microsoft and Alphabet have also shown considerable strength. Furthermore, recent positive developments in the macroeconomic environment, including a surge in consumer spending, seem encouraging.

Now, we can generate a prediction for next week's Summary, Keywords, and Stock Return. However, it's important to note that while the information provided is used as a basis, this is a prediction and the actual outcome may be different.

Summary: Apple's mixed-reality headset previews have received positive feedback from industry experts, increasing investor confidence. However, apprehension about cost-cutting measures and fears of slower hiring might add some uncertainty. Wall Street anticipates a strong earnings report, boosting stock morale. Furthermore, Apple's expansion into India continues to show promise.

Keywords: Apple, mixed-reality headset, investor confidence, cost-cutting, hiring, Wall Street, earnings report, expansion, India.

Stock Return: U2. Given the positive trends and promising developments, the stock is likely to experience upward movement. However, cost-cutting concerns may moderate gains, leading to a prediction of "U2", suggesting a rise of 1% to 2% next week.

Figure 5: GPT-4 outputs its detailed reasoning steps if we simply add an instruction "Can you reason step by step before the finalized output?" to the end of the prompt in Appendix Figure 4. With detailed reasoning steps, GPT-4 captures a previously missed point "Wall Street anticipates a strong earnings report, boosting stock morale", and amends the stock return forecast as U2.

Summary: The stock market is likely to react favorably to Apple's continued focus on cost optimization and revenue generation through new product launches and market expansions. Despite facing regulatory challenges, the company's efforts towards environmental responsibility and diversification of financial offerings are likely to have a positive impact on its long-term prospects. As a result, we expect Apple's stock to see another strong week of growth, with a projected return of U5+.

Keywords: Cost-Cutting, New Product Launches, Expansion, Regulation, Supply Chain

Stock Return: U5+

Figure 6: Explainable forecasting by fine-tuned Open LLaMA 13B model.