



Learning to Generate Explainable Stock Predictions using Self-Reflective Large Language Models

Kelvin J.L. Koa

National University of Singapore
kelvin.koa@u.nus.edu

Ritchie Ng

Eastspring Investments, Singapore
ritchie.ng@easpring.com

Yunshan Ma*

National University of Singapore
yunshan.ma@u.nus.edu

Tat-Seng Chua

National University of Singapore
dcscts@nus.edu.sg

ABSTRACT

Explaining stock predictions is generally a difficult task for traditional non-generative deep learning models, where explanations are limited to visualizing the attention weights on important texts. Today, Large Language Models (LLMs) present a solution to this problem, given their known capabilities to generate human-readable explanations for their decision-making process. However, the task of stock prediction remains challenging for LLMs, as it requires the ability to weigh the varying impacts of chaotic social texts on stock prices. The problem gets progressively harder with the introduction of the explanation component, which requires LLMs to explain *verbally* why certain factors are more important than the others. On the other hand, to fine-tune LLMs for such a task, one would need expert-annotated samples of explanation for every stock movement in the training set, which is expensive and impractical to scale.

To tackle these issues, we propose our Summarize-Explain-Predict (SEP) framework, which utilizes a verbal self-reflective agent and Proximal Policy Optimization (PPO) that allow a LLM teach itself how to generate explainable stock predictions, in a fully autonomous manner. The reflective agent learns how to explain past stock movements through a self-reasoning process, while the PPO trainer trains the model to generate the most likely explanations given the input texts at test-time. The training samples for the PPO trainer are also the responses generated during the reflective process, which eliminates the need for human annotators. Using our SEP framework, we fine-tune a specialized LLM that can outperform both traditional deep-learning and LLM methods in prediction accuracy and Matthews correlation coefficient, for the stock classification task. To justify the generalization capability of our framework, we further test it on the portfolio construction task, and demonstrate its effectiveness through various portfolio metrics. Our code can be accessed through <https://github.com/koa-fin/sep>.

CCS CONCEPTS

• Information systems → Web mining; • Applied computing → Forecasting; Economics.

*Corresponding author.



This work is licensed under a Creative Commons Attribution International 4.0 License.

WWW '24, May 13–17, 2024, Singapore, Singapore
© 2024 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0171-9/24/05.
<https://doi.org/10.1145/3589334.3645611>

KEYWORDS

Stock Prediction, Large Language Models, Explainable AI

ACM Reference Format:

Kelvin J.L. Koa, Yunshan Ma, Ritchie Ng, and Tat-Seng Chua. 2024. Learning to Generate Explainable Stock Predictions using Self-Reflective Large Language Models. In *Proceedings of the ACM Web Conference 2024 (WWW '24)*, May 13–17, 2024, Singapore, Singapore. ACM, Singapore, 12 pages. <https://doi.org/10.1145/3589334.3645611>

1 INTRODUCTION

The Efficient Market Hypothesis (EMH) states that in financial markets, stock prices reflect all available information [20], and should only react to new information. Through mining and analysing external data sources, the goal of investors is to quickly understand the impact of new information on the market, in order to anticipate future stock price movements [27]. However, analyzing the impact of these data on the stock market is a huge undertaking and imposes a heavy workload on financial experts, due to the large volume of information available [24]. Because of this, many have explored the use of deep-learning techniques [23, 37, 42] for stock prediction.

1. Task Prompt:

Given a list of facts, estimate their overall impact on the price movement of AAPL stock. Give your response in this format:
(1) Price Movement, which should be either Positive or Negative.
(2) Explanation, which should be in a single, short paragraph.
Facts: [...]

2. LLM Response:

Price Movement: Mixed
Explanation: The given facts provide a mixed sentiment regarding AAPL stock. On one hand, there are positive indicators such as the stock being at an all-time high, Apple Watch outselling the Swiss watch industry, and Apple's acquisition of NeXT. On the other hand, there are negative aspects including the potential impact of the coronavirus on AirPods orders, France's fine on Apple for slowed iPhones, and the stock price falling after a downgrade. Other facts provide neutral information about stock repurchases, trademark appeals, and the growth of Google's Canadian workforce. Overall, the sentiment is mixed, with both positive and negative factors influencing the perception of AAPL stock.

Figure 1: While LLMs can classify the sentiment of individual texts (highlighted in orange), they are not trained to weigh between them to produce an aggregate prediction (highlighted in blue). An improved response by our fine-tuned LLM will be presented in the results. [...] refers to truncated text.

However, due to their complex and quantitative nature, traditional deep-learning methods in stock prediction are black box models and do not address the explainability of their predictions

[41]. This reduces their usability in practical applications, as users might not be able to trust [4] the results to invest their capital. Even among works that deal with explainable stock predictions [6, 41], the "explanations" are often simply defined as the specific texts that caused the price movement, which are usually obtained by analyzing learnt attention weights [15, 56]. For example, these models could analyze a series of texts regarding Apple stock and determine that its *Positive* prediction is attributed to the text "*Apple reported revenue of \$90.1 billion, beating expectations*". However, these models do not go beyond that to explain *why* these texts caused the stock movement, and require the user to make their own inference.

Today, the emergence of Large Language Models (LLMs) has presented a solution to this problem. Recent surveys [69, 79] have shown that LLMs possess both strong Natural-Language Understanding capabilities, which allow them to perform tasks like text summarization [54] and text classification [43] in a few-shot manner; and strong Natural-Language Generation capabilities, which let them generate human-readable explanations for their own decision-making process [44, 63]. Currently, works that utilize LLMs for stock prediction [10, 76] are few, and use limited techniques such as pre-trained LLMs or instruction tuning. Our work seeks to fill this gap by designing a reinforcement learning (RL) framework which can fine-tune a LLM to generate explanations for stock prediction.

To tackle the explainable stock prediction task using LLMs, we can identify two main challenges. Firstly, it is well-established in past stock prediction literature that social texts are *chaotic*, where the influence of different texts on stock prices can be highly diverse [33, 67]. For example, breaking news such as surprise earnings announcements or crisis events often have a visible impact on the stock price, while unsubstantiated opinions or vague remarks usually cause little to no change [60]. This requires a prediction model to have the ability to weigh the varying impacts of new information [21], and arrive at a maximum-likelihood prediction [28]. Typically, this involves training a regression-based neural network, and is not a known capability of LLMs (see Figure 1). Secondly, the problem becomes progressively harder with the introduction of the explanation component, as it requires the LLM to explain *verbally* why certain information are more important than others. However, to train a LLM for this task using RL [32, 52], one would need good and bad samples [38, 44] of explanations for each price movement in the training set. This requires substantial amount of labour by financial experts, which is expensive and impractical to scale.

To deal with the above-mentioned problems, we propose our Summarize-Explain-Predict (SEP) framework, which utilizes a self-reflective agent [59] and Proximal Policy Optimization (PPO) [57] to let a LLM teach itself how to make explainable stock predictions in a fully autonomous manner (see Figure 2). Firstly, the Summarize module utilizes the strong summarization capabilities of LLMs [54] to convert large volumes of text input data into point-form summaries of factual information. Secondly, in the Explain module, a reflective agent teaches itself to generate correct stock predictions and explain their reasoning [63] given a sequence of summarized facts, via an iterative, verbal self-reflective process [49, 59]. The iterative process additionally allows us to obtain a series of *correct* and *incorrect* predictions with annotated explanations through its past mistakes, which can be used as fine-tuning samples without human-in-the-loop. Lastly, in the Predict module, a specialized LLM

is fine-tuned [32, 52] via PPO training [57] using its own self-taught responses, in order to generate the most likely stock predictions and explanations, given the input texts from an unseen test set.

To demonstrate the effectiveness of the SEP framework, we validate through experimental results that our model is able to outperform both traditional deep-learning and LLM methods in terms of its prediction accuracy and Matthews correlation coefficient (MCC) for the binary stock classification task. We also analyze some responses from the fine-tuned LLM qualitatively, to show how it is better able to understand and weigh the impacts of different information within the input texts. Additionally, to justify the generalization capability of the framework, we test it on the portfolio construction task, by generating explainable weights for a stock portfolio. We also demonstrate the effectiveness of this method through portfolio metrics, such as its profitability and Sharpe Ratio.

The main contributions of this paper are summarized as:

- We investigate the limitations of teaching LLMs to weigh information in multiple texts for stock prediction in an explainable manner, without expert-annotated explanation samples.
- We propose a solution that utilizes a self-reflective agent and PPO techniques, that can allow a LLM to teach itself how to make explainable stock predictions in a fully autonomous manner.
- We validate the effectiveness of SEP through experimental results on tweet data, and show that the fine-tuned LLM is able to provide improvements in both the prediction performance and the quality of its explanations. We further demonstrate the generalizability of the framework by fine-tuning a LLM to generate quantitative weights for multiple stocks, to tackle the portfolio task.

2 RELATED WORKS

In this section, we trace the progress of textual analysis techniques in stock prediction works, and also explore some pioneering works that utilized Large Language Models (LLMs) in the financial domain.

Text Analysis in Stock Prediction. Early text analysis works in stock prediction first studied the effectiveness of using different textual representations of news, such as Bag of Words, Noun Phrases, and Named Entities, in Support Vector Machines (SVM) [58]. These "shallow" features were later replaced in favor of structured information, where events in the form of (*Actor, Action, Object*) tuples were used as inputs for deep neural networks [17, 18].

Later works would define the challenges in text analysis more clearly, which was attributed to the chaotic and diverse nature of text data [33]. This led to the popular use of attention-based models to capture the "most important" information in texts directly from pre-trained text embeddings [15]. Some other notable works include the use of Variational Auto-Encoders (VAEs) to model the latent factors in market information [67], and Transformer models [71].

Most recent works have moved away from improving text analysis methods, and opted instead to enhance the current models with additional forms of information, such as the vocal features from audio data [70] or cross-stock impacts from company relational graphs [42, 56]. In contrast, our work return to purely text-based models, to isolate the effects of text information on stock movements.

Large Language Models in Finance. Out of the existing works that utilize LLMs on general financial tasks, the most well-known one is BloombergGPT [65], which trained a 50B parameters LLM using their existing large financial text corpus. Their model was

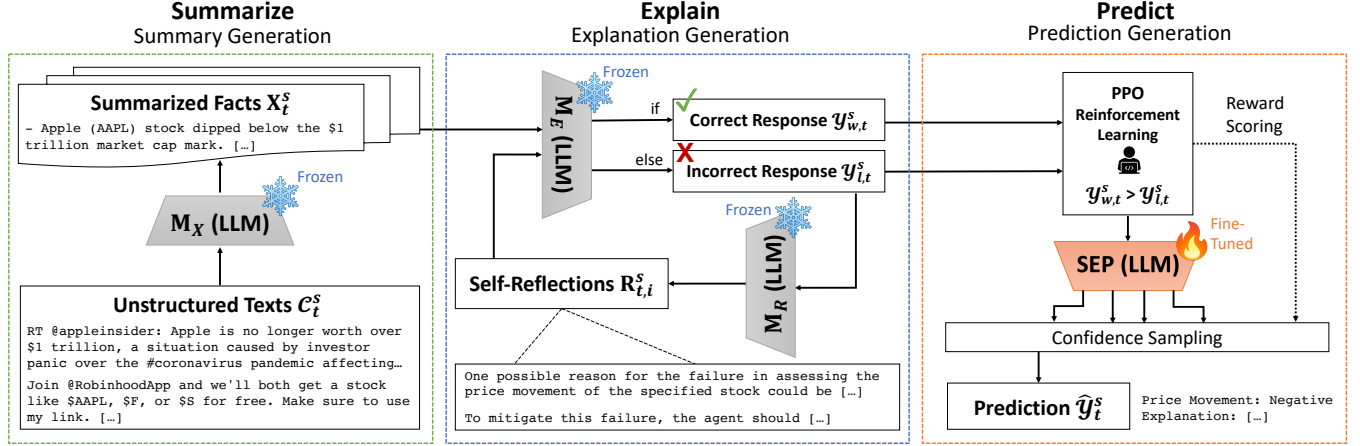


Figure 2: Overall framework of our proposed SEP method, which consists of three components: Summarize, Explain and Predict.

evaluated on several downstream tasks such as sentiment analysis and named-entity recognition (NER), with optimistic results. Along this direction, some works have also attempted to fine-tune their own financial LLM, which include FinMA [66] and FinGPT [68].

Other works explored the use of existing LLMs such as ChatGPT to perform specialized downstream tasks, such as stock sentiment prediction from news headlines [47], and classification of Federal announcements [31]. These early works focused on analyzing *individual* texts, as opposed to a sequence of texts. More recent works have explored the use of LLMs to make stock predictions using sequences of stock-related texts, via instruction-tuning [76] or pre-trained models enhanced with relational graphs [10]. We build on these works by implementing an additional verbal self-reflective agent to learn how to generate better explanations, and a PPO trainer to fine-tune a more specialized LLM for stock predictions.

3 METHODOLOGY

In this section, we first define the task and data for explainable stock prediction. We then present the proposed SEP framework, which was illustrated in Figure 2. There are three main components: (1) a Summarize module, which generates a summary of factual information from the unstructured text inputs; (2) an Explain module, which generates explanations for its stock predictions and refines them through an iterative self-reflective process; and (3) a Predict module, which generates confidence-based predictions after fine-tuning a LLM using its self-generated annotated samples.

3.1 Preliminaries

3.1.1 Problem Formulation. Given a stock $s \in \mathcal{S} = \{s_i\}_{i=1}^O$ and its associated text corpora for the past T days $\{C_{t-T}^s, \dots, C_{t-2}^s, C_{t-1}^s\}$, we aim to generate a stock prediction for the next trading day \hat{y}_t^s , which consists of a binary price movement $\hat{y}_t^s \in \{0, 1\}$ and a human-readable explanation \hat{e}_t^s . Each corpus contains a variable number of unstructured texts $C_t^s = \{c_{t,n}^s\}_{n=1}^{N_t^s}$, where $c_{t,n}^s$ is a single text, and $N_t^s = |C_t^s|$ is the number of texts for the stock s on day t .

3.1.2 Data Collection and Clustering. In this work, we construct a new dataset by following the data collection methodology used for the ACL18 StockNet dataset [67], which is a popular benchmark used in many stock prediction works [22, 41, 56]. The duration of

the original dataset ranges from year 2014–2016, and we collect an updated version for year 2020–2022. Since the previous work, the number of industries have expanded, and the number of tweets have also increased exponentially. We collect data for the top 5 stocks in the 11 industries, giving us a total of 55 stocks. The price data is collected from Yahoo Finance¹, while the tweet data is collected using the Twitter API². Additionally, given the large volume of tweets for each day, we utilize a clustering pipeline via BERTopic [29] to identify the representative tweets for each day. These tweets would be used as the text inputs for all models. More details on the dataset and clustering pipeline can be found in Appendix A.

3.2 Summary Generation

The goal of the Summary module is to generate summarized information from the unstructured input texts. Current LLMs are known for their summarization ability, which surpass even humans [54]. Given that a sequence of raw texts from T days would exceed the character limit, even for 16K-context LLMs, we first prompt a LLM to generate point-form summaries of factual information from the given texts [76] for each day. The prompt takes in two variable inputs: the specified stock s , and the unstructured text inputs C_t^s for each day t . The LLM M_X then generates a summary of facts X_t^s that can be learnt from the input texts, which can include specific information for stock s and other related news in its industry, e.g., "Big Tech stocks, including Apple (AAPL), Google, Amazon, and Facebook, beat earnings expectations." This can be formulated as:

$$X_t^s = M_X(s, C_t^s). \quad (1)$$

Within the prompt, we also provide two in-context examples [75] that were composed from selected cases in the dataset. The examples for all prompts in this work can be found in Appendix B.

3.3 Explanation Generation

The goal of the Explain module is two-fold: While the key aim of the module is to generate clear explanations for stock predictions, the generated explanations also serve as a reasoning step [63] for the LLM to do self-reflection to improve its own predictions [59]. In the following subsections, we discuss the initial prompt design and the subsequent self-reflective process for the module.

¹<https://finance.yahoo.com/>

²<https://developer.twitter.com/>

3.3.1 Explanation Prompting. The prompt for the Explain module contains two variable inputs: the specified stock s , and a *sequence* of extracted information that was generated from the previous module. Given these inputs, the LLM M_E then generate the response \mathcal{Y}_t^s , which should contain the next-day price movement y_t^s , and the annotated explanation e_t^s , i.e., $\mathcal{Y}_t^s = (y_t^s, e_t^s)$. We formalize this as:

$$\mathcal{Y}_t^s = M_E(s, X_{t-T}^s, \dots, X_{t-2}^s, X_{t-1}^s). \quad (2)$$

Similar to the previous summarization prompt, we select two cases from the dataset and manually compose the response trajectories to use as few-shot exemplars [75]. Additionally, the two example cases chosen have specifically one Positive and one Negative movement label, in order to avoid any majority label bias [80]. The prompt trajectories are designed in a fashion similar to ReAct [73], albeit in a singular, prediction-explanation step.

3.3.2 Self-Reflective Process. Current LLMs are not trained to generate stock predictions, which could cause incorrectly-generated annotated examples in the previous step. To tackle this, we deploy the LLM as an autonomous agent that can iteratively improve on its past responses, through a verbal self-reflection loop (see Figure 3). The loop is first seeded with the response from the previous step, i.e., $\mathcal{Y}_{t,0}^s = \mathcal{Y}_t^s$, which is taken to be the initial iteration $i = 0$.

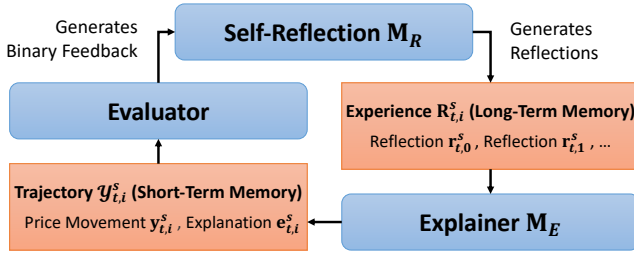


Figure 3: Diagram of the self-reflective process.

From the generated price movement $y_{t,i}^s$, we can obtain a binary feedback by evaluating its alignment with the ground truth. For the incorrect samples, we then prompt a LLM M_R to generate a verbal feedback $r_{t,i}^s$ for each iteration i , given its previous inputs and outputs, which we refer to as its short-term memory [59]. The feedback should explain clearly where it went wrong in its previous reasoning $e_{t,i}^s$, and also come up with a high-level plan to mitigate this failure for the next iteration. The overall formalization is:

$$r_{t,i}^s = M_R(s, X_{t-T}^s, \dots, X_{t-2}^s, X_{t-1}^s, \mathcal{Y}_{t,i}^s). \quad (3)$$

For every iteration, each reflection $r_{t,i}^s$ represent a lesson that the LLM learnt from its failures, which is added to its experiences, or long-term memory [59]. We represent this as a set of reflections, $R_{t,i}^s = [r_{t,0}^s, r_{t,1}^s, \dots, r_{t,i}^s]$. The reflections, together with the original inputs, are fed again into LLM M_E to generate the price movement and explanation for the next iteration. The formalization is:

$$\mathcal{Y}_{t,i+1}^s = M_E(s, X_{t-T}^s, \dots, X_{t-2}^s, X_{t-1}^s, R_{t,i}^s). \quad (4)$$

The prompt and response examples can be found in Appendix B.

Through this process, we are then able to obtain pairs of correct and incorrect responses, for each successful reflection. We define these as $\mathcal{Y}_{w,t}^s = (y_{t,\tilde{i}}^s, e_{t,\tilde{i}}^s)$ and $\mathcal{Y}_{l,t}^s = (y_{t,\tilde{i}-1}^s, e_{t,\tilde{i}-1}^s)$ respectively, where \tilde{i} refers to the iteration in which the reflective process resulted in the LLM M_E generating the correct stock movement.

3.4 Prediction Generation

The goal of the Predict module is to fine-tune a LLM to generate good stock predictions and explanations for the unseen test period. In this section, we discuss the overall fine-tuning process of the model and the subsequent inference procedure at test-time.

3.4.1 Model Fine-Tuning. Following previous works that tackles Reinforcement Learning from Human Feedback (RLHF) [52, 61], we utilize a similar three-step process to fine-tune a LLM. Instead of human feedback, we use the binary evaluations from the reflections to choose the "better" response during training (see Figure 4).

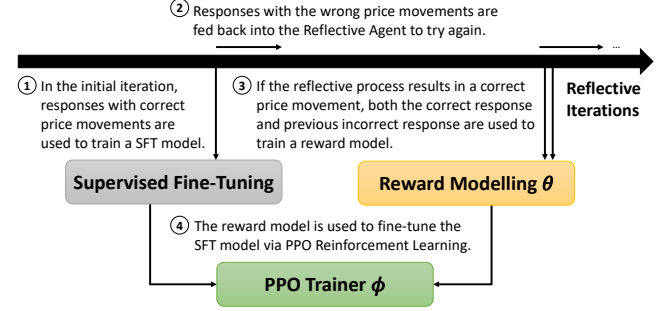


Figure 4: Diagram of the fine-tuning process.

In the first step, we collect the demonstration data, which are taken from the correct predictions in the initial iteration $\mathcal{Y}_{t,0}^s$. These samples do not have corresponding "wrong" responses, as they were taken from the initial prompt. The samples are used to train a supervised policy π^{SFT} using Supervised Fine-Tuning (SFT).

In the second step, we collect the comparison data \mathcal{D} , which contains pairwise correct and incorrect responses $\mathcal{Y}_{w,t}^s, \mathcal{Y}_{l,t}^s$ for each structured input X_t^s , taken from the successful reflection iterations. These are used to train a reward model r_θ , which learns to give higher reward scores to the correct responses. Specifically, we train the model to minimize the following cross-entropy loss [61]:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(X, \mathcal{Y}_w, \mathcal{Y}_l, s, t) \sim \mathcal{D}} \left[\log \left(\sigma \left(r_\theta(X_t^s, \mathcal{Y}_{w,t}^s) - r_\theta(X_t^s, \mathcal{Y}_{l,t}^s) \right) \right) \right]. \quad (5)$$

In the third step, we use the reward model to optimize the trained policy using PPO [57]. We first initialize the model with the supervised policy π^{SFT} , and use it to generate predictions $\hat{\mathcal{Y}}_t^s$ for randomly selected samples X_t^s from the overall dataset $\mathcal{D}_{\pi_\phi^{RL}}$. Next, the reward model r_θ is used to generate a reward for each response. We then try to optimize a PPO model π_ϕ^{RL} by maximizing the overall reward. This is achieved by minimizing the following loss objective:

$$\mathcal{L}(\phi) = -\mathbb{E}_{(X, \mathcal{Y}, s, t) \sim \mathcal{D}_{\pi_\phi^{RL}}} \left[r_\theta(X_t^s, \hat{\mathcal{Y}}_t^s) - \beta \log \frac{\pi_\phi^{RL}(\hat{\mathcal{Y}}_t^s | X_t^s)}{\pi^{SFT}(\hat{\mathcal{Y}}_t^s | X_t^s)} \right]. \quad (6)$$

We note that the objective includes an additional term that penalizes the KL divergence between the trained policy π_ϕ^{RL} and the supervised policy π^{SFT} [35], which is used to deter the policy from collapsing into a single mode [61], and prevent it from generating responses that are too different from those of the original reference model π^{SFT} [74]. The term is controlled by a hyper-parameter β .

3.4.2 Confidence-based Sampling. During inference, the unstructured input texts C_t^s are first summarized using a pre-trained LLM.

We then use the trained policy π_{ϕ}^{RL} to generate the next-day predictions \hat{Y}_t^s from the summarized facts X_t^s . For generating predictions, we use a best-of- n sampler, where we generate n responses and use the scores from reward model r_{θ} to select the best response [74].

4 EXPERIMENT

We evaluate the performance of SEP on our collected dataset. Our work aims to answer the following three research questions:

- **RQ1:** How does the SEP model perform against traditional deep-learning and other LLM methods in the stock prediction task, in both its classification accuracy and quality of explanations?
- **RQ2:** How does each proposed component (*i.e.*, Summarize, Explain, Predict) help to improve the performance of the SEP model?
- **RQ3:** Is the SEP framework sufficiently generalizable to other finance-related tasks, such as explainable portfolio construction?

4.1 Experimental Settings

4.1.1 Baselines. To demonstrate the effectiveness of our SEP-trained model, we compare it against baselines from both traditional deep-learning models and fine-tuned Large Language Models (LLMs).

Deep Learning Models:

- **VAE+Attention** [67]: In this model, a Variational Auto-encoder (VAE) [36] is used to model the latent market factors within texts. News-level [33] and temporal [19] attention are used to weigh texts on their salience in the corpus and across the input period. Texts are represented on the word level using GloVe [53].
- **GRU+Attention** [56]: This model utilize a hierarchical attention model using Gated Recurrent Networks (GRU) [55] with multiple stages of attention layers [2, 72] to capture the corpus-level and day-level importance of each text. The texts are encoded on the sentence level using the Universal Sentence Encoder [8].
- **Transformer** [70]: This model uses stacked transformer encoders to perform multi-headed self-attention on the token- and sentence-level, before decoding with multiple feed-forward layers [71]. For preprocessing, the texts are encoded on the token level using the Whole Word Masking BERT (WWM-BERT) [16].

Large Language Models:

- **GPT-3.5-turbo** [52]: We provide the same prompts to a GPT-3.5-turbo-16k LLM for comparison. ChatGPT has previously been explored in other stock sentiment prediction works [47, 76].
- **Vicuna-7b-v1.5** [11]: Similarly, we provide the same prompts to a Vicuna-7b-v1.5-16k LLM. This is also the model used for fine-tuning in our work, and serves as a base model for comparison.
- **FinGPT-Forecaster** [68]: This is an instruction-tuned LLM model by FinGPT, which can take in a series of market news to make stock predictions. This is the most recent model for our task.

For the deep-learning methods, we keep only the text-processing components for an equivalent comparison. The inputs for all models are the unstructured representative tweets C_t^s . Following the previous works that deals with the binary stock classification task [18, 22, 67], we use the prediction accuracy and Matthews Correlation Coefficient (MCC) as our evaluation metrics. For all LLM results, any predictions that are made in the wrong format, or are "Neutral" or "Mixed", will be considered as an incorrect prediction.

Additionally, a key feature of the SEP framework is the Summarize module, which extracts key information from unstructured tweets for the LLM to base its predictions on. However, there are

some days when there are no useful information to be found in the tweets. In such cases, there can still be significant price movements, which could be due to external factors such as stock price stochasticity [37] or daily interest rates fluctuations [1]. For the LLM experiments, we report both the results before and after removing such cases. In practice, this could be seen as a benefit of LLMs, as it is able to actively tell that it has not enough information to make a prediction, and investors could choose to either look for more information to analyze or not invest their capital for the day.

4.1.2 Implementation Details. For the Summarize and Explain components, we evaluate two different models for generating the responses. We use OpenAI GPT-3.5-turbo-16k for the top 1 stock in each industry, and Vicuna-13b-v1.5-16k for the remaining stocks. Both are set to a temperature of zero. The input length is $T = 5$.

For training the prediction model, we use Vicuna-7b-v1.5-16k. The LLM is trained using *trl*, which supports transformer reinforcement learning with PPO trainer³. For the supervised fine-tuning, we run two epochs with a learning rate of 3×10^{-4} . For the reward model tuning, we run one epoch with a learning rate of 2×10^{-4} . For the RL learning with PPO, we run four epochs with a learning rate of 1.4×10^{-5} . All components are trained using 4-bit quantized low-rank adapters (LoRA) [32] with a setting of $r = 8$. At inference, we set $n = 4$ for n -shot sampling, where the temperature of the model is set at 0.7. The best response, based on reward scoring, will be used as the selected output for the final comparisons.

4.2 Performance Comparison (RQ1)

In this section, we evaluate both the prediction and explanation responses generated by our SEP model, through quantitative and qualitative comparisons against the relevant baselines.

4.2.1 Prediction Performance. Table 1 reports the quantitative results on the stock prediction task. On the prediction accuracy, we observe that the SEP model fine-tuned on the GPT-generated explanations (Table 1, left) was able to obtain the best results, achieving an improvement of 2.4% over the strongest baseline (GRU+Att) using all texts. On the other hand, the SEP model fine-tuned on explanations generated by Vicuna-v1.5 (Table 1, right) under-performed the baselines in terms of accuracy. A possible reason for this is that the Vicuna-generated explanations used for training the model are prone to hallucinations, which could negatively impact the reasoning ability of the SEP model (see Figure 5). The poorer performance of GPT-3.5, a pre-trained LLM, is largely attributed to its inability to make decisive predictions from mixed sentiments. The instruction-tuned FinGPT-Forecaster is able to improve on this by guiding the LLM towards trained responses, which are in the correct format. Finally, our SEP model produces the best accuracy, likely due to its additional self-reflective process and reinforcement learning.

For this task, a more telling metric is the Matthews Correlation Coefficient (MCC), which takes into account the ratios of True and False Positives and Negatives of the predictions [12, 13]. Given that not all stock movements are necessarily caused by the provided texts, the accuracy results might not be fully indicative of the model's natural language processing capabilities, as it includes some random guesses on the non-informative texts. After filtering for informative texts only, we can see increases in the MCC ratio, possibly from having less random guesses in the prediction results.

³<https://huggingface.co/docs/trl>

Table 1: Performance comparisons in accuracy and MCC of our SEP model against baselines. The best results are boldfaced.

Models		Top 1 Stock, GPT-3.5				Remaining Stocks, Vicuna			
		All Texts		Informative Texts		All Texts		Informative Texts	
		Accuracy	MCC	Accuracy	MCC	Accuracy	MCC	Accuracy	MCC
Deep-Learning Models	VAE+Att	49.96	0.0046	-	-	49.83	0.0070	-	-
	GRU+Att	50.15	0.0125	-	-	50.77	0.0189	-	-
	Transformer	50.06	0.0089	-	-	50.17	0.0135	-	-
Large Language Models	GPT-3.5	20.80	0.0094	29.35	0.0298	17.57	0.0027	22.99	0.0052
	Vicuna	40.85	0.0114	45.29	0.0368	39.66	0.0115	43.30	0.0301
	FinGPT	47.61	0.0158	51.56	0.0384	45.76	0.0161	46.12	0.0379
	SEP (Ours)	51.38	0.0302	54.35	0.0993	47.59	0.0203	50.57	0.0508

On the MCC metric, our SEP model was able to outperform all models under all settings, which showcase the true ability of the model to understand the impacts of natural language texts on stock movements, after accounting for the random guesses. Under the all-texts setting, we are able to outperform the strongest deep-learning baseline (GRU+Att) by 0.0177 for the GPT-3.5-based model, and 0.0014 for the Vicuna-based model. After filtering for informative texts only, our fine-tuned SEP model is also able to outperform the strongest LLM baseline, FinGPT-Forecaster, by 0.0609 and 0.0129 for the GPT-3.5 and Vicuna-based SEP models respectively.

[illegible]

Figure 5: An example of a hallucinated response from Vicuna. [...] refers to truncated text, which are all repeated text.

4.2.2 Explanation Performance. In addition to generating better predictions, the natural advantage of using LLMs over traditional deep-learning methods is their capability to generate explanations for their predictions. Here, we compare the generated explanations qualitatively between the pre-trained LLMs and our SEP model.

After SEP fine-tuning, we can observe two main improvements. The first deals with the ability to decisively weigh between news information to make a stock movement prediction. While pre-trained LLMs are known to be able to classify the sentiment of individual texts [47, 77], they typically do not try to weigh between these sentiments and make a decisive stock prediction, even if specifically requested by the prompt (see Figure 1). This is generally an easier task to tackle, which is similar to fine-tuning an expert LLM [30], albeit ours is trained without human experts-in-the-loop. Figure 6 shows an example of how our SEP model can learn how to make a decisive stock prediction after the fine-tuning process.

The second improvement deals with the ability to generate better-quality explanations. This is a more difficult task for LLMs, as it requires them to not only understand the meaning of natural language texts, but also to correctly reason out their overall impact on the stock price movement. Through the SEP framework, our LLM first learns to reason out the correct explanations via self-reflection and teach them to the PPO model, which learns to determine heuristically what is the most probable explanation at test-time. For this comparison, we came up with a set of metrics for explanation quality, and use GPT-4 to rate each response from 1 to 7 for the samples

Before Fine-Tuning

Price Movement: Mixed

Explanation: The given facts provide a mixed sentiment regarding AAPL stock. On one hand, there are positive indicators such as the stock being at an all-time high, Apple Watch outselling the Swiss watch industry, and Apple's acquisition of NeXT. On the other hand, there are negative aspects including the potential impact of the coronavirus on AirPods orders, France's fine on Apple for slowed iPhones, and the stock price falling after a downgrade. Other facts provide neutral information about stock repurchases, trademark appeals, and the growth of Google's Canadian workforce. Overall, the sentiment is mixed, with both positive and negative factors influencing the perception of AAPL stock.

↓

After Fine-Tuning

Price Movement: Negative

Explanation: The given facts regarding AAPL stock indicate a mixed sentiment with predominantly negative factors. The delay in production and potential impact on sales due to the coronavirus outbreak raise concerns about Apple's financial performance. The restart of a key China plant suggests a potential recovery, but its overall effect is uncertain. Additionally, the mention of a fine imposed on Apple for slowing down old phones adds a negative aspect. The fluctuations in the stock, discussions about market cap and potential earnings, as well as options trading, provide some neutral factors that do not significantly impact the overall sentiment. However, considering the predominantly negative factors mentioned, the estimate leans towards a negative sentiment for AAPL stock.

Figure 6: An example of SEP learning to make a decisive, aggregate stock prediction. After fine-tuning, the SEP model is able to make a correct Negative prediction based on the predominantly negative events contained within the texts.

in the top-1 stocks. The average score for each metric is reported in Table 2. Explanations of the metrics can be found in Appendix C.

From Table 2, we can make the following observations:

- The highest scores come from more generalizable metrics, such as *Consistency with Information*. Some metrics are sometimes not observable given that no such information is in the texts, which will lower their scores. However, it is fair to compare the relative scores across the LLMs, as their input texts are the same.
- All LLMs give good-quality explanations even if the prediction is wrong, as they can naturally understand the input texts and make reasonable comparisons, but have done so incorrectly. Thus, prediction accuracy should still be the first metric to look at.
- Our SEP model was able to score the highest for all metrics. We note that the model was not trained on these metrics, but was only provided a reward based on correct binary predictions. Through its self-reflection and reinforcement learning, the SEP framework was able to intrinsically teach the model to better compare these factors, in order to generate better predictions.

Table 2: Comparison of explanation quality from our SEP model against baselines. The best results are boldfaced.

Metric	GPT-3.5	Vicuna	SEP (Ours)
Relevance to Stock Movement	5.407	5.396	5.449
Financial Metrics	2.957	3.146	3.334
Global & Industry Factors	3.180	3.576	3.700
Company Developments	3.905	4.066	4.224
Temporal Awareness	3.951	4.066	4.170
Balance of Positive & Negative	4.030	4.084	4.224
Contextual Understanding	4.012	4.098	4.193
Clarity & Coherence	6.271	6.325	6.439
Consistency with Information	5.575	5.652	6.006
Sensitivity to Updates	4.112	4.172	4.362

4.3 Ablation Study (RQ2)

In this section, we evaluate the efficiency of each of the three proposed components: the Summarize, Explain and Predict modules.

4.3.1 Summarize Module. The Summarize module reduces the noise and length of the input texts by extracting only the important, factual information. For the ablation study, we compare against the performance of using non-summarized, raw social texts in our trained SEP model. To keep the input lengths within the LLM's token limit, we try two ways of using the raw texts: Taking the 30 most shared texts and randomly sampling 30 texts, for each day.

Table 3: Comparison of SEP with non-summarized and summarized input texts. The best results are boldfaced.

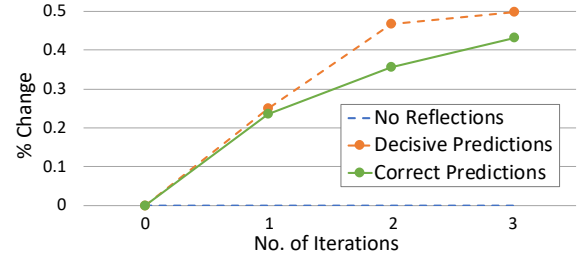
	Top 1 Stock (GPT-3.5)		Remaining Stocks (Vicuna)	
	Accuracy	MCC	Accuracy	MCC
Non-Summ. (Random 30)	50.75	0.0208	40.81	-0.0037
Non-Summ. (Top 30)	50.81	0.0219	41.27	0.0023
Summarized (All Texts)	51.38	0.0302	47.59	0.0203
Summarized (Informative)	54.35	0.0993	50.57	0.0508

From Table 3, we can make the following observations:

- Using the most shared texts is better than randomly sampling.
- The model trained on Vicuna-generated responses fare much worse without summarizing than the GPT-3.5-trained model. This could be attributed to the texts causing more hallucination (Table 5), given their chaotic content (*e.g.*, emojis, spam, *etc.*).
- The summarized texts provide better results. One possible reason here could be due to having information from more than 30 texts. However, it would also show that the summarization process did not lose any important information that would cause degradation.
- Finally, removing the non-informative texts, which is only possible with the Summarize module, provides the best results.

4.3.2 Explain Module. In the SEP model, we have observed two improvements: 1) the ability to make *decisive* stock predictions from mixed sentiments; and 2) the ability to make *correct* stock predictions (*i.e.*, better prediction accuracy). In order to fine-tune the LLM to produce these predictions and explanations, the Explain module must first try to generate correctly-annotated samples through binary feedback and self-reflection. To demonstrate its effectiveness,

we plot the percentage change in the number of generated decisive and correct predictions after each of its reflective iteration.

**Figure 7: Percentage change in number of decisive and correct explanation samples over the self-reflective process.**

From Figure 7, we see that with multiple self-reflective iterations, the model generates more and more decisive and correct annotated samples, to be used for fine-tuning. We also observe that there is a greater number of decisive samples produced given that it is an easier task, which starts to slow down as more samples become non-Neutral. Overall, the number of decisive samples grew by 49.8% while the correct samples grew by 43.2% after 3 iterations, which highlights the effectiveness of the Explain module in generating annotated explanation samples, without the help of human experts.

4.3.3 Predict Module. For the Predict module, we conduct an ablation study over different variants of the model. We remove one additional component for each variant, *i.e.*, no n -shot sampling at inference [SEP (1-shot)]; no PPO reinforcement learning [SEP (no PPO)]; and no explanations [SEP (binary)], which is simply instruction-tuning the LLM to make binary up/down predictions. We make the comparisons on the top-1 stock from each industry.

Table 4: Ablation study. The best results are boldfaced.

	All Texts		Informative Texts	
	Accuracy	MCC	Accuracy	MCC
SEP (binary)	40.84	-0.0042	42.75	0.0295
SEP (no PPO)	44.08	0.0144	45.29	0.0368
SEP (1-shot)	50.08	0.0270	52.54	0.0715
SEP (Ours)	51.38	0.0302	54.35	0.0993

From Table 4, we can make the following observations:

- The addition of the explanation component during the instruction-tuning process, *i.e.*, from SEP (binary) to SEP (no PPO), gives the model an average improvement of 6.9%. It is likely that by tuning the LLM to generate explanations, we are able to elicit a reasoning process from the LLM [63] when generating stock movement predictions, resulting in better prediction accuracy.
- The instruction-tuned variant, *i.e.*, SEP (no PPO), shows very similar results to the base model that it is tuned on (*i.e.*, the Vicuna model in Table 1). It is possible that the instruction tuning process has no impact on the SEP model given that the samples, taken before the reflective iterations (*i.e.*, Step 1 in Figure 4), are "easy" samples that the base model could already handle. We also note that supervised-tuned models have been seen to produce little to even negative improvements in previous literature [61].
- The largest improvement comes from the PPO reinforcement learning, *i.e.*, from SEP (no PPO) to SEP (1-shot), with an average improvement of 14.8%. This highlights the ability of the

PPO trainer in teaching the LLM to generate stock predictions more effectively. Additionally, the n -shot sampling weighs between n generated samples using the learnt reward model to select the best output. The average improvement of this variant *i.e.*, 3.0% from **SEP (1-shot)** to **SEP (Ours)**, further reinforces the usefulness of the reward model trained during the PPO process.

4.4 Portfolio Optimization (RQ3)

From our results, we have observed that the SEP framework is able to teach an LLM to weigh the impact of information within the input texts in a binary manner. We further explore its generalization capability by using it to fine-tune a LLM to weigh between information within its own generated explanations quantitatively, in order to generate portfolio weights for the stock portfolio task.

For the portfolio task, we follow the same method as above to fine-tune a LLM. Here, the input information are now all the generated explanations for the basket of stocks for each day. For this experimental task, we filter only the stocks with positive predictions, in order to reduce the number of stocks the LLM have to weigh, and to prevent negative weights (hence setting a no short-sales constraint [37]). We then prompt the LLM to generate portfolio weights given the outlook for each given stock (see Figure 8). An example of the prompt and response can be found in Appendix B.

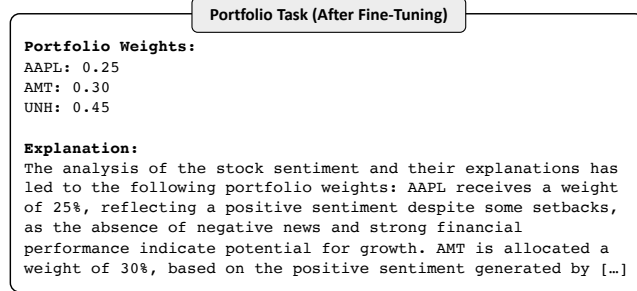


Figure 8: An example response for the portfolio task after fine-tuning with SEP. [...] refers to truncated text.

As there is no binary feedback for this task, in each self-reflective iteration, we provide the reflective LLM with the overall profits based on the provided portfolio weights, and prompt it to reflect on how it can improve itself to obtain higher profits. The reflections are then used to generate an updated set of portfolio weights. Finally, we feed both sets of generated weights into a PPO trainer, where the one with higher profits is used as the "better" response.

We compare the performances of portfolios generated by three different LLMs: GPT-3.5-turbo, Vicuna, and our fine-tuned SEP model. We also include three baselines: the 1/N portfolio, where all 11 stocks in the basket are bought at equal weights [14]; the S&P500 stock market index; and Positive-Only, where only the predicted positive stocks are bought at equal weights. The latter can also be seen as evaluating the results of the original stock prediction LLM in a practical setting, without the portfolio weighing prompts.

We evaluate the portfolio performance using four metrics: the overall gain, which simply sums up the gains for each day; the cumulative gain, which is the final gain after re-investing any additional profits or losses over the evaluation period; the standard deviation of the profits; and the annualized Sharpe Ratio [46].

Table 5 reports the portfolio results. From the table, we observe:

Table 5: Portfolio results comparison. The best results are boldfaced. The Sharpe Ratio values are annualized.

Approach	Overall	Cumulative	Std. Dev.	Sharpe
1/N	-0.0330	-0.0502	1.613e-2	-0.225
Market Index	0.0180	0.0003	1.533e-2	0.123
Positive-Only	0.1243	0.1065	1.911e-2	0.807
GPT-3.5	0.1497	0.1353	1.893e-2	0.980
Vicuna	0.1541	0.1447	1.731e-2	1.104
SEP (Ours)	0.1661	0.1569	1.792e-2	1.150

- The **Positive-Only** portfolio, *i.e.*, evenly buying the stocks that are predicted to be Positive, already showcases good performance. This highlights the capability of our original stock prediction model to produce good trading signals in a practical setting.
- For the standard deviation results, we note that the top 2 portfolio methods, *i.e.*, 1/N and **Market Index**, contains more number of stocks, which allow them to spread out the stock price fluctuations more evenly. However, their Sharpe Ratios are still lower than the other models, which shows a lower reward-to-risk ratio.
- The pre-trained LLM models, *i.e.*, **GPT-3.5** and **Vicuna**, already shows better performance than the **Positive-Only** portfolio in most metrics, which shows the capabilities of using LLMs to weigh between information to produce portfolio weights.
- Our **SEP** model was able to outperform all other methods in most portfolio metrics, and achieve comparable performance in its standard deviation, which showcases the effectiveness of our SEP framework. In addition to the shown metrics, we also re-emphasize the ability of the LLM-based models to *explain* the generated portfolio weights, which further adds to the interpretability and trustability of their results for practitioners.

5 CONCLUSION

In this work, we explored the explainable stock prediction task, which was largely difficult to solve before generative models. For this task, we highlighted two challenges: the limitations of current LLMs in weighing varied market factors to make aggregate stock predictions, and the lack of annotated training samples for fine-tuning LLMs to make explanations. To tackle these challenges, we proposed our SEP framework, which utilizes a verbal self-reflective agent and PPO techniques to let a LLM teach itself how to generate stock explanations in a fully autonomous manner. Through experimental results, we validated that our SEP model is able to outperform both traditional deep-learning and LLM methods in the accuracy of the predictions and quality of the generated explanations. Furthermore, we also demonstrated the generalizability of the SEP framework by fine-tuning a model for the portfolio task, and show that it can outperform other methods on portfolio metrics.

6 ACKNOWLEDGEMENT

This research is supported by the National Research Foundation, Singapore under its Industry Alignment Fund – Pre-Positioning (IAF-PP) Funding Initiative, by the National Research Foundation, Singapore through the National Cybersecurity R&D Lab at the National University of Singapore under its National Cybersecurity R&D Programme (Award No. NCR25-NCL P3-0001). Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore.

REFERENCES

- [1] Md Mahmudul Alam and Gazi Uddin. 2009. Relationship between interest rate and stock price: empirical evidence from developed and developing countries. *International Journal of Business and Management (ISSN 1833-3850)* 4, 3 (2009), 43–51.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural Machine Translation by Jointly Learning to Align and Translate. In *ICLR*.
- [3] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023* (2023).
- [4] Or Biran and Kathleen R. McKeown. 2017. Human-Centric Justification of Machine Learning Predictions. In *IJCAI*. ijcai.org, 1461–1467.
- [5] Ricardo J. G. B. Campello, Davoud Moulavi, and Jörg Sander. 2013. Density-Based Clustering Based on Hierarchical Density Estimates. In *PAKDD (2) (Lecture Notes in Computer Science, Vol. 7819)*. Springer, 160–172.
- [6] Salvatore M. Carta, Sergio Consoli, Luca Piras, Alessandro Sebastian Podda, and Diego Reforgiato Recupero. 2021. Explainable Machine Learning Exploiting News and Domain-Specific Lexicon for Stock Market Forecasting. *IEEE Access* 9 (2021), 30193–30205.
- [7] Stephen Casper, Jason Lin, Joe Kwon, Gatlen Culp, and Dylan Hadfield-Menell. 2023. Explore, Establish, Exploit: Red Teaming Language Models from Scratch. *arXiv preprint arXiv:2306.09442* (2023).
- [8] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175* (2018).
- [9] Canyu Chen and Kai Shu. 2023. Combating misinformation in the age of llms: Opportunities and challenges. *arXiv preprint arXiv:2311.05656* (2023).
- [10] Zihan Chen, Lei Nico Zheng, Cheng Lu, Jialu Yuan, and Di Zhu. 2023. ChatGPT Informed Graph Neural Network for Stock Movement Prediction. *arXiv preprint arXiv:2306.03763* (2023).
- [11] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90% chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023) (2023).
- [12] Davide Chicco and Giuseppe Jurman. 2020. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC genomics* 21, 1 (2020), 1–13.
- [13] Davide Chicco, Niklas Tötsch, and Giuseppe Jurman. 2021. The Matthews correlation coefficient (MCC) is more reliable than balanced accuracy, bookmaker informedness, and markedness in two-class confusion matrix evaluation. *BioData Min.* 14, 1 (2021), 13.
- [14] Victor DeMiguel, Lorenzo Garlappi, and Raman Uppal. 2009. Optimal versus naive diversification: How inefficient is the 1/N portfolio strategy? *The review of Financial studies* 22, 5 (2009), 1915–1953.
- [15] Shumin Deng, Ningyu Zhang, Wen Zhang, Jiaoyan Chen, Jeff Z. Pan, and Huajun Chen. 2019. Knowledge-Driven Stock Trend Prediction and Explanation via Temporal Convolutional Network. In *WWW (Companion Volume)*. ACM, 678–685.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT (1)*. Association for Computational Linguistics, 4171–4186.
- [17] Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. 2014. Using Structured Events to Predict Stock Price Movement: An Empirical Investigation. In *EMNLP*. ACL, 1415–1425.
- [18] Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. 2015. Deep Learning for Event-Driven Stock Prediction. In *IJCAI*. AAAI Press, 2327–2333.
- [19] Yujuan Ding, Yunshan Ma, Lizi Liao, Wai Keung Wong, and Tat-Seng Chua. 2022. Leveraging Multiple Relations for Fashion Trend Forecasting Based on Social Media. *IEEE Trans. Multim.* 24 (2022), 2287–2299.
- [20] Eugene F Fama. 1970. Efficient capital markets: A review of theory and empirical work. *The journal of Finance* 25, 2 (1970), 383–417.
- [21] Eugene F Fama and Kenneth R French. 2015. A five-factor asset pricing model. *Journal of financial economics* 116, 1 (2015), 1–22.
- [22] Fuli Feng, Huimin Chen, Xiangnan He, Ji Ding, Maosong Sun, and Tat-Seng Chua. 2019. Enhancing Stock Movement Prediction with Adversarial Training. In *IJCAI*. ijcai.org, 5843–5849.
- [23] Fuli Feng, Xiangnan He, Xiang Wang, Cheng Luo, Yiqun Liu, and Tat-Seng Chua. 2019. Temporal Relational Ranking for Stock Prediction. *ACM Trans. Inf. Syst.* 37, 2 (2019), 27:1–27:30.
- [24] Fuli Feng, Moxin Li, Cheng Luo, Ritchie Ng, and Tat-Seng Chua. 2021. Hybrid Learning to Rank for Financial Event Ranking. In *SIGIR*. ACM, 233–243.
- [25] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2023. Bias and fairness in large language models: A survey. *arXiv preprint arXiv:2309.00770* (2023).
- [26] Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple Contrastive Learning of Sentence Embeddings. In *EMNLP (1)*. Association for Computational Linguistics, 6894–6910.
- [27] Gyoza Gidofalvi and Charles Elkan. 2001. Using news articles to predict stock price movements. *Department of computer science and engineering, university of california, san diego* 17 (2001).
- [28] Stefano Giglio, Bryan Kelly, and Dacheng Xiu. 2022. Factor models, machine learning, and asset pricing. *Annual Review of Financial Economics* 14 (2022), 337–368.
- [29] Maarten Grootendorst. 2022. BERTopic: Neural topic modeling with a class-based TF-IDF procedure. *arXiv preprint arXiv:2203.05794* (2022).
- [30] Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *arXiv preprint arXiv:2301.07597* (2023).
- [31] Anne Lundgaard Hansen and Sophia Kazinnik. 2023. Can ChatGPT Decipher FedSpeak? Available at SSRN (2023).
- [32] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *ICLR*. OpenReview.net.
- [33] Ziniu Hu, Weiqing Liu, Jiang Bian, Xuanzhe Liu, and Tie-Yan Liu. 2018. Listening to Chaotic Whispers: A Deep Learning Framework for News-oriented Stock Trend Prediction. In *WSDM*. ACM, 261–269.
- [34] Sima Jannati, Alok Kumar, Alexandra Niessen-Ruenzi, and Justin Wolfers. 2023. In-group bias in financial markets. Available at SSRN 2884218 (2023).
- [35] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard. 2019. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456* (2019).
- [36] Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *ICLR*.
- [37] Kelvin J. L. Koa, Yunshan Ma, Ritchie Ng, and Tat-Seng Chua. 2023. Diffusion Variational Autoencoder for Tackling Stochasticity in Multi-Step Regression Stock Price Prediction. In *CIKM*. ACM, 1087–1096.
- [38] Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. 2023. RLAI: Scaling Reinforcement Learning from Human Feedback with AI Feedback. *arXiv preprint arXiv:2309.00267* (2023).
- [39] Jure Leskovec, Anand Rajaraman, and Jeffrey D. Ullman. 2014. *Mining of Massive Datasets, 2nd Ed.* Cambridge University Press.
- [40] Haochen Li, Maria Polukarov, and Carmine Ventre. 2023. Detecting Financial Market Manipulation with Statistical Physics Tools. In *ICAIF*. ACM, 1.
- [41] Shuqi Li, Weiheng Liao, Yuhang Chen, and Rui Yan. 2023. PEN: Prediction-Explanation Network to Forecast Stock Price Movement with Better Explainability. In *AAAI*. AAAI Press, 5187–5194.
- [42] Wei Li, Ruihan Bao, Keiko Harimoto, Deli Chen, Jingjing Xu, and Qi Su. 2020. Modeling the Stock Relation with Graph Network for Overnight Stock Movement Prediction. In *IJCAI*. ijcai.org, 4541–4547.
- [43] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110* (2022).
- [44] Hao Liu, Carmelo Sferrazza, and Pieter Abbeel. 2023. Chain of hindsight aligns language models with feedback. *arXiv preprint arXiv:2302.02676* 3 (2023).
- [45] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [46] Andrew W Lo. 2002. The statistics of Sharpe ratios. *Financial analysts journal* 58, 4 (2002), 36–52.
- [47] Alejandro Lopez-Lira and Yuehua Tang. 2023. Can chatgpt forecast stock price movements? return predictability and large language models. *arXiv preprint arXiv:2304.07619* (2023).
- [48] Yunshan Ma, Chenchen Ye, Zijian Wu, Xiang Wang, Yixin Cao, Liang Pang, and Tat-Seng Chua. 2023. Structured, Complex and Time-complete Temporal Event Forecasting. *arXiv preprint arXiv:2312.01052* (2023).
- [49] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651* (2023).
- [50] Leland McInnes, John Healy, and James Melville. 2018. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426* (2018).
- [51] Alexandra Niessen-Ruenzi and Stefan Ruenzi. 2019. Sex Matters: Gender Bias in the Mutual Fund Industry. *Manag. Sci.* 65, 7 (2019), 3001–3025.
- [52] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda

- Askill, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *NeurIPS*.
- [53] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global Vectors for Word Representation. In *EMNLP. ACL*, 1532–1543.
- [54] Xiao Pu, Mingqi Gao, and Xiaojun Wan. 2023. Summarization is (Almost) Dead. *arXiv preprint arXiv:2309.09558* (2023).
- [55] Yao Qin, Dongjin Song, Haifeng Chen, Wei Cheng, Guofei Jiang, and Garrison W. Cottrell. 2017. A Dual-Stage Attention-Based Recurrent Neural Network for Time Series Prediction. In *IJCAI. ijcai.org*, 2627–2633.
- [56] Ramit Sawhney, Shivam Agarwal, Arnav Wadhwa, and Rajiv Ratn Shah. 2020. Deep Attentive Learning for Stock Movement Prediction From Social Media Text and Company Correlations. In *EMNLP (1)*. Association for Computational Linguistics, 8415–8426.
- [57] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [58] Robert P. Schumaker and Hsinchun Chen. 2009. Textual analysis of stock market prediction using breaking financial news: The AZFin text system. *ACM Trans. Inf. Syst.* 27, 2 (2009), 12:1–12:19.
- [59] Noah Shinn, Federico Cassano, Beck Labash, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: Language Agents with Verbal Reinforcement Learning. *arXiv:2303.11366 [cs.AI]*
- [60] Timm O Sprenger and Isabell M Welp. 2011. News or noise? The stock market reaction to different types of company-specific news events. *The Stock Market Reaction to Different Types of Company-Specific News Events* (2011).
- [61] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F. Christiano. 2020. Learning to summarize with human feedback. In *NeurIPS*.
- [62] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How does llm safety training fail? *arXiv preprint arXiv:2307.02483* (2023).
- [63] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *NeurIPS*.
- [64] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023. NExT-GPT: Any-to-Any Multimodal LLM. *arXiv preprint arXiv:2309.05519* (2023).
- [65] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhakaran Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564* (2023).
- [66] Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. 2023. PIXIU: A Large Language Model, Instruction Data and Evaluation Benchmark for Finance. *arXiv preprint arXiv:2306.05443* (2023).
- [67] Yumo Xu and Shay B. Cohen. 2018. Stock Movement Prediction from Tweets and Historical Prices. In *ACL (1)*. Association for Computational Linguistics, 1970–1979.
- [68] Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. 2023. FinGPT: Open-Source Financial Large Language Models. *arXiv preprint arXiv:2306.06031* (2023).
- [69] Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *arXiv preprint arXiv:2304.13712* (2023).
- [70] Linyi Yang, Jiazheng Li, Ruihai Dong, Yue Zhang, and Barry Smyth. 2022. NumHTML: Numeric-Oriented Hierarchical Transformer Model for Multi-Task Financial Forecasting. In *AAAI*. AAAI Press, 11604–11612.
- [71] Linyi Yang, Tin Lok James Ng, Barry Smyth, and Ruihai Dong. 2020. HTML: Hierarchical Transformer-based Multi-task Learning for Volatility Prediction. In *WWW. ACM / IW3C2*, 441–451.
- [72] Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alexander J. Smola, and Edward H. Hovy. 2016. Hierarchical Attention Networks for Document Classification. In *HLT-NAACL*. The Association for Computational Linguistics, 1480–1489.
- [73] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *ICLR*. OpenReview.net.
- [74] Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh Murthy, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, et al. 2023. Retroformer: Retrospective large language agents with policy gradient optimization. *arXiv preprint arXiv:2308.02151* (2023).
- [75] Seonghyeon Ye, Hyeonbin Hwang, Sohee Yang, Hyeonung Yun, Yireun Kim, and Minjoon Seo. 2023. In-context instruction learning. *arXiv preprint arXiv:2302.14691* (2023).
- [76] Xinli Yu, Zheng Chen, Yuan Ling, Shujing Dong, Zongyi Liu, and Yanbin Lu. 2023. Temporal Data Meets LLM—Explainable Financial Time Series Forecasting. *arXiv preprint arXiv:2306.11025* (2023).
- [77] Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023. Sentiment Analysis in the Era of Large Language Models: A Reality Check. *arXiv preprint arXiv:2305.15005* (2023).
- [78] Xuan Zhang and Wei Gao. 2023. Towards llm-based fact verification on news claims with a hierarchical step-by-step prompting method. *arXiv preprint*

arXiv:2310.00305 (2023).

- [79] 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).
- [80] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate Before Use: Improving Few-shot Performance of Language Models. In *ICML (Proceedings of Machine Learning Research, Vol. 139)*. PMLR, 12697–12706.

A DATASET AND CLUSTERING PIPELINE

In this section, we include additional details on the statistics of the collected dataset and the overall clustering pipeline.

A.1 Dataset

In this work, we construct a new dataset by following the data collection methodology used for the **ACL18** StockNet dataset [67], updated for the year 2020–2022. Since the previous work, the number of tweets have increased exponentially (see Table 6). To keep the texts within a reasonable length, we first employ a clustering pipeline to obtain the most representative tweets for each day.

A.2 Clustering Pipeline

Following previous works that perform clustering on full-length documents for LLM inputs [48], we make use of the BERTopic [29] pipeline for clustering: First, we generate embeddings for the tweets using a pre-trained language model RoBERTa [45], which have also been fine-tuned using the SimCSE [26] framework. Next, UMAP [50] was used for dimensionality reduction of the embeddings, and HDBSCAN [5] was used to cluster them into semantically similar groups. Finally, we use a class-based TF-IDF procedure [29, 39] to rank and extract the most representative tweet for each cluster.

For the hyper-parameters, we set the number of neighbors for UMAP dimensionality reduction as 15. For HDBSCAN clustering, the minimum cluster size is set to 10. The statistics of the tweet data before and after clustering can be found in Table 6.

Table 6: Statistics of tweets before and after clustering.

	Avg. tweets	Avg. tokens	Max tweets	Max tokens
Before Clustering	469	27,951	46,569	1,911,495
After Clustering	16	1,068	1,599	63,392

In total, the dataset consists of tweets for 757 trading days. The overall number of samples used is 29,997, which is split in a train-test ratio of 8:2. Within the training set, 10% of the generated explanation samples are used for validation during fine-tuning.

B FULL PROMPT EXAMPLES

In this section, we provide some examples of the prompts used in SEP and the responses. Examples for four tasks are shown:

- Figure 9 shows an example for the summarization task, where summarized factual information is generated from the chaotic input tweets. In the example, we can see that tweets that contain useless information, such as unsubstantiated comments, are ignored by the LLM. Additionally, the facts extracted from the tweets are also summarized in a concise and succinct manner.
- Figure 10 shows a successful example for the explanation task. In the example, we can see that while there are some positive news, there are more recent and impactful negative facts which caused a negative price movement. The example showcases the ability of the LLM to weigh between these factors effectively, and generate the correct price movement with a reasonable explanation.

- Figure 11 shows an example for the reflection task. In the example, the incorrect previous response is fed into the LLM to generate a reflection, which consists of what went wrong and a plan on how to mitigate this problem. The reflection tells the LLM to further consider the positive earnings, overall market for big tech companies, and the long-term strategic initiatives, which allowed it to obtain a correct prediction in the next iteration.
- Figure 12 shows an example for the portfolio task. Given the self-predicted explanations for all positive stocks for each day, the LLM further weigh between their outlook to recommend the amount of each stock to purchase. In the example, we can see the LLM gave more weight to factors such as digital transformation, which could signify potential future growth for the company.

C EVALUATION OF EXPLANATION QUALITY

The explanation of the metrics used in Table 2 are given below. These metrics are manually curated by us through the assistance of ChatGPT. There are currently limited works on evaluating generated stock explanations, given that it is a relatively new application.

- **Relevance to Stock Movement:**
 - Does the explanation focus on factors directly related to the stock’s movement?
- **Financial Metrics:**
 - Does the explanation include relevant financial metrics (e.g., earnings estimates, market cap)?
 - Are these metrics explained in the context of their impact on stock performance?
- **Global & Industry Factors:**
 - Does the explanation consider broader economic conditions or industry trends that may impact the stock?
 - Is there an understanding of how global events could influence the stock’s performance?
- **Company Developments:**
 - Are specific developments related to the company discussed?
 - Is there an understanding of how these developments might influence the stock?
- **Temporal Awareness:**
 - Does the explanation consider the timing of events and developments?
 - Is there an acknowledgment of the temporal dynamics of the stock market?
- **Balance of Positive & Negative:**
 - Is there an attempt to balance positive and negative factors?
 - Does the explanation recognize mitigating factors that could counteract positive or negative sentiments?
- **Contextual Understanding:**
 - Does the explanation demonstrate a nuanced understanding of the context in which the news is presented?
 - Is there an awareness of the complexities and uncertainties in predicting stock movements?
- **Clarity & Coherence:**
 - Is the explanation clear and easy to understand?
 - Does it present a coherent narrative that connects various factors logically?
- **Consistency with Information:**
 - Is the information presented consistent with known facts and events?

- Are there any inaccuracies or contradictions in the explanation?

- **Sensitivity to Updates:**

- Does the explanation show sensitivity to the possibility of changing circumstances or new information that could affect the stock?

D FUTURE WORK

There are some directions that can be explored in future works. Firstly, we address the possibility of cumulative errors in the SEP framework. At each stage, poorly generated summaries or explanations could lead to poorer responses in the next step. In practice, it is possible for experts to vet through the responses before using them, which would be an easier task than generating them manually. However, more can be done to increase the robustness of the generated responses and reduce the need for human-in-the-loop. Secondly, using additional data sources, such as knowledge graphs [31] or audio features [70], could increase the quality of the predictions. At the same time, such works would also help to explore the multi-modal capabilities of the most recent LLM upgrades [3, 64]. Finally, as this is a relatively new task, there are currently limited works on evaluating the generated stock explanations. Further studies can be done to improve the metrics created in this work.

E ETHICAL USE OF DATA

For this research, we have utilized datasets derived from publicly available sources, and no human annotators were involved in the data collection process. Rights pertaining to the data used, such as text data, remain the sole property of the original rights holders. This study is intended exclusively for academic purposes only.

There are potential ethical and social implications of using LLMs for stock prediction. We list some of them here and suggest possible ways to mitigate them when deploying our model for practical use:

- *Risk of Manipulation.* Market manipulation has always been a problem in stock markets [40]. Using LLMs for stock prediction can increase this risk, given their known vulnerabilities such as jailbreak prompting [62] and model red-teaming [7]. To mitigate these, there should be measures to scrutinize user inputs to the model before processing them. Access to the LLM’s internal knowledge base should be restricted to authorized users only.
- *Misinformation.* While the point of explainable forecasting is to generate trustable results, LLMs can also be leveraged to generate deceptive misinformation [9]. Measures should be taken to verify the correctness of facts before utilizing them in the explanations, either by automated verification [78] or human-in-the-loop.
- *Prediction Bias.* It is known that LLMs tend to inherit stereotypes or existing biases due to the internet-based data they are trained on [25]. As stock prediction with LLMs is relatively new, it is unknown whether existing investor biases [34, 51] will also carry over into the LLMs’ generated responses. Some mitigation strategies include removing biased responses, verifying all information are factual, before training the LLM with reinforcement learning.

In general, the most effective mitigation strategy is to include human-in-the-loop to anticipate and mitigate various potential risks. While LLMs can assist humans in labor-intensive tasks such as processing large volume of texts and analyzing their stock market impacts, it is not able to replace the need for human oversight.

Summarize Prompt

Given a list of tweets, summarize all key facts regarding AAPL stock.

Here are some examples: [...]

Tweets:
 RT @CalebGregory304: Markets will be down big again tomorrow! I'm on the hunt for companies with strong balance sheets. I think I'll be b..."
 RT @appleinsider: Apple is no longer worth over \$1 trillion, a situation caused by investor panic over the #coronavirus pandemic affecting...
 RT @Firefight9221: So much chop intraday and false signals. Had the \$SPY Drop at the open then pop in the afternoon, while \$AAPL was fading..."
 \$AAPL Trump Administration Gives Apple More Tariff Relief
 Trade officials approved the company's request to remove tariffs on the Apple Watch. \n OAKLAND, Calif. - U.S. trade offici...
<https://t.co/zwmxc0IXot>
 Join @RobinhoodApp and we'll both get a stock like \$AAPL, \$F, or \$S for free. Make sure to use my link. <https://t.co/4d6pAckBI2>
 RT @OphirGottlieb: \$AAPL iPhone sales no longer limited to two per customer
 Mar 23 Monday, \$AAPL NEW LOW IS \$212.63 and hit Bearish targets \$229.24 📉 \$221.25 📉
 [...]

Facts:
 Apple (AAPL) stock dipped below the \$1 trillion market cap mark.
 Apple requested tariff relief for the Apple Watch and received approval.
 The iPhone sales limit per customer has been lifted by Apple.
 Apple stock reached a new low of \$212.63 on March 23.
 [...]

Figure 9: An example prompt and its response for generating summarized facts. Underlined words denote the end of the input prompt text. [...] refers to truncated text.

Reflect Prompt

You are an advanced reasoning agent that can improve based on self reflection. You will be given a previous reasoning trial in which you were given access to a list of facts to assess their overall impact on the price movement of GOOG stock. You were unsuccessful in tackling the task because you gave the wrong price movement. In a few sentences, diagnose a possible reason for failure and devise a new, concise high-level plan that aims to mitigate the same failure.

Previous trial:
 Facts: [...]
 Price Movement: [...]
 Explanation: [...]

Reflection: The failure in the previous reasoning trial was due to an overemphasis on the negative impacts of the COVID-19 pandemic on Alphabet's core search and advertising revenue and insufficient consideration of the positive aspects of the earnings report.

New Plan: To mitigate this failure, I will focus on a more balanced assessment of the facts, considering both positive and negative aspects of the earnings report and the pandemic's impact on Alphabet's business. I will also consider the broader market context and the performance of other major tech companies to provide a more comprehensive analysis of the sentiment surrounding GOOG stock.

1. Analyze the earnings report and the pandemic's impact on Alphabet's core search and advertising revenue, as well as other business segments.
2. Compare Alphabet's performance with that of other major tech companies, such as Amazon, Apple, and Facebook, to gauge the overall market sentiment for big tech stocks.
3. Evaluate the company's strategic initiatives, such as Google Cloud partnerships and investments in Waymo and Verily [...]

Figure 11: An example prompt and its response for generating reflections. Underlined words denote the end of the input prompt text. [...] refers to truncated text.

Explain Prompt

Given a list of facts, estimate their overall impact on the price movement of AAPL stock. Give your response in this format:
 (1) Price Movement, which should be either Positive or Negative.
 (2) Explanation, which should be in a single, short paragraph.

Here are some examples: [...]

Facts:
 2020-03-22
 Apple's stock price experienced a fast crash due to the COVID-19 pandemic.
 Apple donated millions of masks to healthcare professionals in the US and Europe.
 2020-03-23
 [...]

Price Movement: Negative

Explanation: The COVID-19 pandemic has had a significant impact on Apple's stock performance, causing a decline in its stock price. Wedbush lowered their earnings estimates for Apple, reflecting concerns about the company's financial performance. Additionally, Apple's stock dipped below the \$1 trillion market cap mark and reached a new low on March 23rd. The disruptions caused by the pandemic have also resulted in lowered sales expectations for Apple in the first half of 2020. While Apple has made efforts to contribute to relief efforts by donating masks to healthcare professionals, these positive actions have not been sufficient to counteract the negative sentiment surrounding the stock. Overall, the combination of the decline in stock performance, lowered earnings estimates, and reduced sales expectations contribute to the negative sentiment surrounding AAPL stock.

Figure 10: An example prompt and its response for generating stock explanations. Underlined words denote the end of the input prompt text. [...] refers to truncated text.

Portfolio Prompt

Given a list of stock price movements and their explanations, weigh the amount of each stock to purchase for a stock portfolio. Give your response in this format:
 (1) Portfolio Weights, which should sum up to 1.
 (2) Explanation, which should be in a single paragraph.

Here are some examples: [...]

Summary:
 WMT: Positive
 Explanation: [...]

UNH: [...]

UPS: [...]

Portfolio Weights:
 WMT: 0.20
 UNH: 0.30
 UPS: 0.50

Explanation: The portfolio weights for WMT, UNH, and UPS are determined based on the provided stock price movements and their explanations. WMT receives a weight of 20% due to the positive sentiment reflected in the tweets, which mention strong quarterly reports and sustainability initiatives. UNH is allocated a weight of 30% based on the positive sentiment from the increased dividend payout and strong financial performance. UPS is given the highest weight of 50% due to the positive sentiment from the digital transformation appointment, labor stability, and absence of negative tweets or information regarding the company's stock performance. The allocation of weights aims to create a balanced and diversified portfolio while considering the positive sentiment of each stock.

Figure 12: An example prompt and its response for generating portfolio weights. Underlined words denote the end of the input prompt text. [...] refers to truncated text.