# 3S-Trader: A Multi-LLM Framework for Adaptive Stock Scoring, Strategy, and Selection in Portfolio Optimization

Kefan Chen
The University of Adelaide, Australia
kefan.chen@student.adelaide.edu.au

Hussain Ahmad
The University of Adelaide, Australia
hussain.ahmad@adelaide.edu.au

Diksha Goel
CSIRO's Data61, Australia
diksha.goel@csiro.au

Claudia Szabo
The University of Adelaide, Australia
claudia.szabo@adelaide.edu.au

## Abstract

Large Language Models (LLMs) have recently gained popularity in stock trading for their ability to process multimodal financial data. However, most existing methods focus on single-stock trading and lack the capacity to reason over multiple candidates for portfolio construction. Moreover, they typically lack the flexibility to revise their strategies in response to market shifts, limiting their adaptability in real-world trading. To address these challenges, we propose **3S-Trader**, a training-free framework that incorporates scoring, strategy, and selection modules for stock portfolio construction. The scoring module summarizes each stock's recent signals into a concise report covering multiple scoring dimensions, enabling efficient comparison across candidates. The strategy module analyzes historical strategies and overall market conditions to iteratively generate an optimized selection strategy. Based on this strategy, the selection module identifies and assembles a portfolio by choosing stocks with higher scores in relevant dimensions. We evaluate our framework across four distinct stock universes, including the Dow Jones Industrial Average (DJIA) constituents and three sector-specific stock sets. Compared with existing multi-LLM frameworks and time-series-based baselines, 3S-Trader achieves the highest accumulated return of **131.83%** on DJIA constituents with a Sharpe ratio of **0.31** and Calmar ratio of **11.84**, while also delivering consistently strong results across other sectors.

## CCS Concepts

• **Applied computing → Economics**; • **Computing methodologies → Natural language processing**.

## Keywords

Stock Trading, Portfolio Management, Large Language Models, Self-Reflective Framework, Explainable AI

## 1 Introduction

In the stock market, a portfolio refers to the structured allocation of various stock assets [18, 34], carefully selected to achieve specific investment objectives. Compared to single-stock trading, portfolio-based strategies diversify exposure across multiple assets, mitigating the impact of individual stock volatility [41]. Constructing an effective portfolio requires analyzing diverse sources of market information, such as price trends, company fundamentals, and macroeconomic signals, to identify the most valuable assets for investment allocation [16, 32]. However, this task becomes increasingly challenging as the number of candidate stocks grows, especially when dealing with large volumes of heterogeneous, multimodal data that must be processed and compared in a coherent manner [29, 39].

With the rapid advancements in AI, particularly LLMs [3, 25], across domains such as cybersecurity [7, 15, 22, 42], software engineering [1, 21], and cloud computing [4–6, 27], the technology has also demonstrated significant benefits in the financial sector [45]. Equipped with powerful natural language understanding and reasoning capabilities [2, 11, 14], LLMs can analyze and summarize both financial texts [8] and indicators in a way similar to that of human experts and incorporate this information into their trading strategies [29, 32]. While the mainstream applications of LLMs in finance have primarily focused on single-stock operations, including price prediction, trend classification (up or down), and position adjustment [24]. Although such applications can offer useful signals, they are often insufficient for guiding actual stock selection for portfolio management. For instance, a model may predict high returns for a particular stock, but fail to account for associated volatility, making it a less desirable choice in practice. Similarly, strategies that recommend buying or selling individual stocks are not capable of allocating capital across multiple stocks. Furthermore, conventional trading models struggle to learn from past behavior and adapt their strategies accordingly. While reinforcement learning (RL) offers a way to guide output optimization [28], it often incurs substantial training costs and heavily relies on carefully crafted reward functions [17, 36]. However, due to the highly volatile nature of financial markets, it becomes challenging for RL-based models to learn robust and generalizable strategies [20, 23].

To tackle the challenges outlined above, we propose 3S-Trader, a training-free framework capable of constructing portfolios directly from candidate stocks' recent market signals. Moreover, it can iteratively refine its selection strategy by reflecting on past trading decisions, as illustrated in Figure 1. The framework consists of four
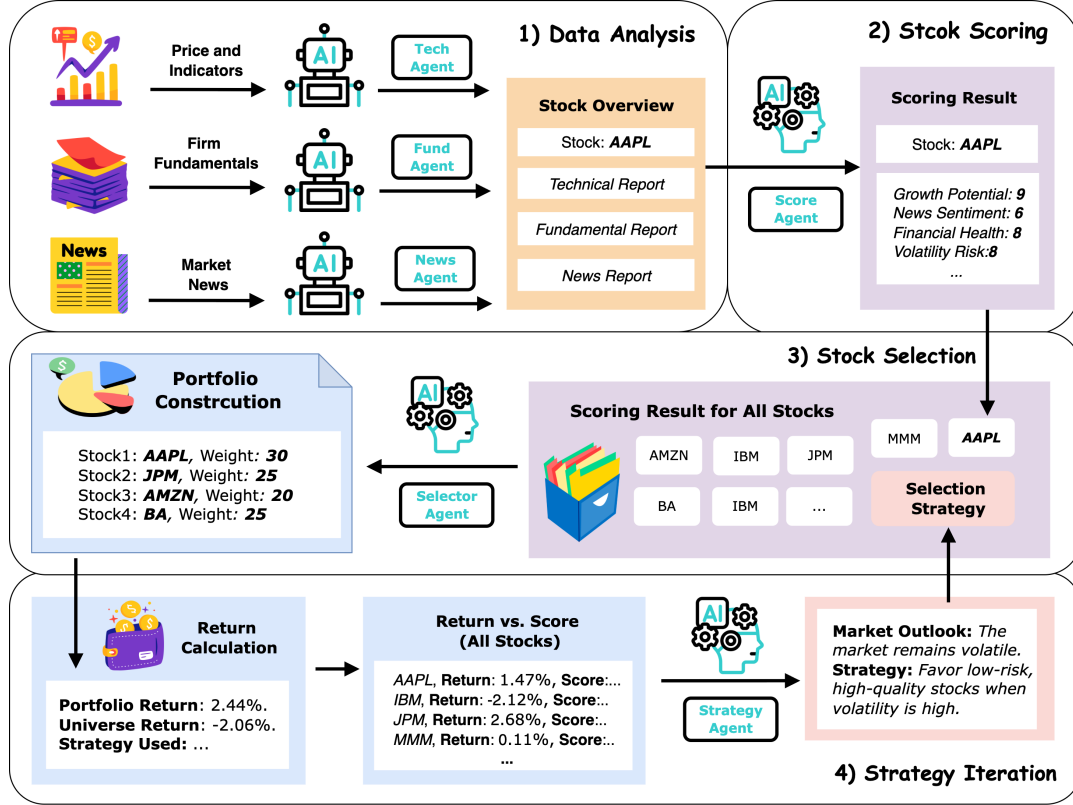
Figure 1: 3S-Trader Framework.

main stages. First, in the Data Analysis stage, three specialized LLM-based agents, the News Agent, Fundamental Agent, and Technical Agent, analyze recent market news, company fundamentals, and technical indicators for each stock to generate textual summaries. Next, the Score Agent evaluates each stock along multiple scoring dimensions such as growth potential, volatility risk, and news sentiment, producing stock-level scoring reports. Based on these reports and the current selection strategy (e.g., favoring financially healthy and low-volatility stocks), the Selector Agent selects a subset of stocks and allocates weights to construct the portfolio. Finally, after each trading round, the Strategy Agent analyzes the relationship between candidate scores and their realized returns to identify which types of stocks the market favors, and adjusts its strategy accordingly. Our major contributions are summarized as follows:

- We propose a training-free and easily deployable trading framework that directly constructs stock portfolios based on recent market information.
- We enhance traditional LLM-based trading pipelines by coupling strategy refinement with multi-dimensional scoring, thereby offering trading strategies a clear adjustment direction, without the need for supervision or rewards.
- We demonstrate the effectiveness and generality of our framework through experiments on stocks from multiple industry

sectors, showing its practical value for real-world investment scenarios.

## 2 Related Work

**Time-Series Models for Stock Prediction.** A common paradigm in traditional quantitative finance is to model stock trading as a time series forecasting problem. Auto-regressive models such as ARIMA and GARCH[10, 13] are widely adopted for modeling short-term trends and volatility in stock prices. More recently, deep learning techniques have been introduced to model complex temporal dependencies. Long Short-Term Memory (LSTM) networks and Temporal Convolutional Networks (TCNs) have shown promising results in capturing sequential patterns in financial time series [12, 19]. Transformer-based models like Informer and Stockformer further improve long-range forecasting ability by leveraging attention mechanisms [37, 47]. While these models achieve promising results in forecasting, they primarily focus on price predictions. As a result, they often overlook portfolio-level considerations such as risk control and cross-asset comparison.

**Multi-LLM Framework for Investment Guidance.** LLMs have been proposed to leverage their text processing capabilities in financial tasks such as sentiment analysis, financial question answering, and trend prediction. [9, 30, 40]. To achieve a complete pipeline covering the process from information processing to decision generation, frameworks that integrate LLMs with distinct

functionalities have been proposed. We categorize these frameworks into two types: *Single-Step* and *Reflective*. The *Single-Step* framework typically aggregates and restructures financial texts over a given period, after which one or more agents summarize the information. A designated output agent then directly relies on this summarized text to generate trading guidance [32, 44, 46]. Extending this design, the *Reflective* framework introduces an additional agent to analyze the model's recent behavior. The analysis is subsequently fed into the next trading cycle, thereby guiding its decision-making process [31, 33].

Building on these works, we upgrade the multi-LLM framework by introducing a multi-dimensional scoring mechanism that simplifies summarized information and enables effective cross-stock comparison. Furthermore, by recording historical strategy trajectories and analyzing the relationship between scores and returns, our framework provides more interpretable and iterative strategy refinement.

## 3 Preliminaries

### 3.1 Problem Formulation

We consider a weekly stock portfolio construction problem. At the beginning of each week $t$, a set of stocks $X = \{x_1, x_2, \ldots, x_n\}$ is available for selection. The goal is to construct a portfolio $\mathbf{w}_t = [w_t^{(1)}, w_t^{(2)}, \ldots, w_t^{(n)}]^\top$, where $w_t^{(i)} \in [0, 1]$ denotes the proportion of capital allocated to stock $x_i$. To allow for holding cash, we impose the constraint $\sum_{i=1}^n w_t^{(i)} \leq 1$.

Trades are executed in the following manner: on the first trading day of week $t$, capital is allocated to selected stocks according to $\mathbf{w}_t$ at the opening price; all positions are then liquidated at the closing price on the last trading day of the same week. For each stock $x_i \in X$, the weekly return $r_t^{(i)}$ is calculated as $r_t^{(i)} = (p_{\text{sell}}^{(i)} - p_{\text{buy}}^{(i)})/p_{\text{buy}}^{(i)}$, where $p_{\text{buy}}^{(i)}$ and $p_{\text{sell}}^{(i)}$ denote the buy and sell prices of stock $x_i$ in week $t$, respectively. Given the return vector $\mathbf{r}_t = [r_t^{(1)}, r_t^{(2)}, \ldots, r_t^{(n)}]^\top$, the total portfolio return is computed as $R_t = \mathbf{w}_t^\top \mathbf{r}_t$.

### 3.2 Data Collection

We collect three types of data for each stock: stock price and technical indicators, firm fundamentals, and market news. In this subsection, we describe how each type of data is processed for the portfolio construction task at week $t$.

*3.2.1 Price and Technical Indicators.* The raw price data consists of daily stock prices and trading volume. Based on these time series, we compute several technical indicators, including SMA, ATR, RSI, MACD, and Bollinger Bands. These indicators are commonly used in financial technical analysis [35, 38].

At week $t$, we extract the daily closing prices $\text{Price}_{i,d}$ and technical indicators $\text{Indicators}_{i,d}$ for stock $x_i$ over the preceding four calendar weeks $\mathcal{W}_{t-4:t-1}$, where $d$ indexes the days within that range. These values are concatenated into a plain text $\text{tech}_{i,t}$, denoted as the technical input of stock $x_i$ for week $t$:

$$\text{tech}_{i,t} = \text{ConcatText}(\{\text{Price}_{i,d}, \text{Indicators}_{i,d} \mid d \in \mathcal{W}_{t-4:t-1}\}) \quad (1)$$

*3.2.2 Market News.* The raw news data includes article titles, summaries, and their associated stock symbols. At week $t$, we collect

all news items related to stock $x_i$ from the previous week. Similar to Equation (1), we define the news input as:

$$\text{news}_{i,t} = \text{ConcatText}(\{\text{RawNews}_{i,d} \mid d \in \mathcal{W}_{t-1}\}) \quad (2)$$

Here, $\text{RawNews}_{i,d}$ denotes the news content related to stock $x_i$ published on day $d$, and $\mathcal{W}_{t-1}$ represents the calendar week immediately preceding week $t$.

*3.2.3 Firm Fundamentals.* The fundamental data consists of firm-specific earnings reports, balance sheets, and cash flow statements, which are updated on a quarterly basis. To enable trend analysis, we concatenate the fundamental records from the four most recent fiscal quarters released before week $t$. The resulting text serves as the fundamental input of stock $x_i$ at week $t$, denoted as:

$$\text{fund}_{i,t} = \text{ConcatText}(\{\text{RawFund}_{i,q} \mid q \in Q_{t_q-3:t_q}\}) \quad (3)$$

Here, $\text{RawFund}_{i,q}$ represents the fundamental data of stock $x_i$ reported in fiscal quarter $q$, and $Q_{t_q-3:t_q}$ refers to the four most recent quarters available before week $t$.

## 4 Methodology

In this section, we introduce the architecture of our proposed framework, 3S-Trader. As shown in Figure 1, the overall workflow is structured into four main steps. The framework incorporates six LLM-based agents that collaborate across these steps to support adaptive portfolio management. We implement all LLM-based agents using GPT-4o[1]. In the following subsections, we detail each step and explain how the agents interact to process information, refine strategies, and make investment decisions.

### 4.1 Market Analysis

In this part, we leverage three specialized agents: *News Agent*, *Fundamental Agent*, and *Technical Agent*, to analyze different types of input data. Specifically, at week $t$, for each stock $x_i$, based on its aggregated news text $\text{news}_{i,t}$ defined in Equation (2), the *News Agent* generates an analysis report $\alpha_{i,t}^{\text{news}}$ as follows:

$$\alpha_{i,t}^{\text{news}} = \text{agent}_{\text{news}}(\text{news}_{i,t}, \text{prompt}_{\text{news}}) \quad (4)$$

Similarly, we obtain $\alpha_{i,t}^{\text{tech}}$ and $\alpha_{i,t}^{\text{fund}}$ from the *Technical Agent* and *Fundamental Agent*, respectively, using corresponding inputs and prompts. These outputs are then concatenated to form the full data overview of stock $x_i$ at week $t$, denoted as $o_{i,t}$:

$$o_{i,t} = \text{ConcatText}(\alpha_{i,t}^{\text{news}}, \alpha_{i,t}^{\text{tech}}, \alpha_{i,t}^{\text{fund}}) \quad (5)$$

Example prompts used by the corresponding agents are illustrated in Figure 2.

### 4.2 Stock Scoring

In this part, the *Score Agent* performs multi-dimensional scoring based on each stock's recent market signals. It is instructed to assign a score from 1 to 10 for each dimension, reflecting the stock's relative strength in that aspect.

In designing the scoring dimensions, we focused on their relevance to the available input data. To ensure that each dimension is grounded in observable evidence, we derive them directly from

---

[1]https://openai.com/gpt-4o

---

**Prompt1: News Agent**

You are a financial news analysis agent. Your task is to filter and summarize recent news related to the stock *{stock code}*. The news content below includes summaries or full articles from the past week: *{raw news text}*

Please provide a concise and insightful weekly summary of the stock's recent news. Your output will be used to help a downstream stock selection agent make informed weekly investment decisions.

---

**Prompt2: Technical Agent**

You are a stock price analysis agent. Your task is to analyze the recent technical indicators and price data of the stock *{stock code}*. Below is the stock's daily technical indicator and prices from the past 4 weeks: *{raw technical text}*.

Please provide a technical analysis of the stock's recent performance. Your output will be used to help a downstream stock selection agent make informed weekly investment decisions.

---

**Prompt3: Fundamental Agent**

You are a stock fundamentals analysis agent. Your task is to analyze the recent financial performance of the stock *{stock code}* based on its past 4 quarterly reports. Below is the stock's recent financial data, including 4 quarters of: Income statements, Balance sheets, Cash flow statements. *{raw fundamental text}*.

Please provide a summary of the stock's fundamental trends. You may consider trends in revenue, profit, expenses, margins, cash flow, and balance sheet strength, as well as any notable improvements or warning signs.

---

**Figure 2: Example prompts used by News Agent, Technical Agent, and Fundamental Agent for market analysis.**

three primary sources: stock price movements, firm fundamentals, and market news. Based on these, we define the following six dimensions:

- **Financial Health**: Evaluates a company's current financial stability. A higher score reflects stronger fundamentals and lower short-term risk.
- **Growth Potential**: Assesses the company's future expansion capacity based on investment plans, and industry growth outlook. A higher score suggests stronger long-term earnings potential.
- **News Sentiment**: Reflects overall sentiment polarity extracted from recent news articles. A higher score implies more positive news coverage and investor perception.
- **News Impact**: Assesses the breadth and duration of news influence. Higher scores reflect more sustained impacts, e.g., from political events or industry-level shifts.
- **Price Momentum**: Captures recent upward or downward trends in stock price movement. A higher score reflects a stronger and more consistent upward price trend.
- **Volatility Risk**: Quantifies the level of recent price fluctuations, indicating risk exposure. A higher score represents higher volatility and less stable price behavior.

For stock selection at week $t$, we first obtain the data overview $o_{i,t}$ for stock $x_i$ as defined in Equation (5). This overview is then

---

**Prompt4: Score Agent**

You are an expert stock evaluation assistant. Tasked with assessing each stock using three input types: News summary, Fundamental analysis, and Recent price behavior.

From these inputs, evaluate the stock along six scoring dimensions. For each dimension: provide a score from 1 to 10, and give a brief justification (1–2 short sentences max).

Use only the information provided below. If anything is missing, score conservatively and state that in the reason.

**stock**: *{stock code}*
**News Summary**: *{news summary}*
**Fundamental Analysis**: *{fundamental analysis}*
**Price and Technical Analysis**: *{technical analysis}*
**Scoring Dimensions** (1–10):

1. Financial Health – based on profitability, debt, cash flow, etc.
2. Growth Potential– based on investment plans, innovation, and expansion prospects.
3. News Sentiment – overall tone of the news.
4. News Impact – the breadth and duration of news influence.
5. Price Momentum – recent trends, strength, and consistency
6. Volatility Risk – recent price stability or instability (higher = more risk)

---

**Figure 3: Example prompt used by the Score Agent for stock scoring.**

processed by the Score Agent to produce a textual scoring result, denoted as $s_{i,t}$. This scoring process can be formalized as:

$$s_{i,t} = \text{agent}_{\text{score}}(o_{i,t}, \text{prompt}_{\text{score}}) \qquad (6)$$

An example prompt used by the Score Agent is illustrated in Figure 3.

## 4.3 Stock Selection

In this part, the *Selector Agent* constructs the final portfolio based on two key inputs: the scoring results of all candidate stocks produced by the *Score Agent*, and the selection strategy provided by the *Strategy Agent*. For clarity, we will introduce the *Strategy Agent* later. Here, it suffices to note that the strategy is a textual description $\pi_t$ that specifies which types of stocks should be preferred, guiding the *Selector Agent* to prioritize stocks with higher scores in the relevant dimensions. For example, a strategy $\pi_t$ might be:

*Increase emphasis on financial health and reduce exposure to high-volatility stocks, as recent returns indicate stronger performance from fundamentally stable companies.*

The *Selector Agent* is tasked with selecting up to five stocks and assigning weights, allowing cash holding to avoid market downturns. Hence, the total portfolio weight may be less than 1. For the portfolio construction task at week $t$, the agent receives the scores $s_{i,t}$ for all candidate stocks $x_i \in \mathcal{X}$, along with the strategy $\pi_t$. The resulting portfolio can be expressed as:

$$\mathbf{w}_t = \text{agent}_{\text{select}}(\text{ConcatText}(\{s_{i,t}\}_{i=1}^{n}), \pi_t, \text{prompt}_{\text{select}}) \qquad (7)$$

where $\mathbf{w}_t = [w_t^{(1)}, w_t^{(2)}, \ldots, w_t^{(n)}]^{\top}$ is the portfolio weight vector, with $\sum_{i=1}^{n} w_t^{(i)} \leq 1$, and at most 5 elements of $\mathbf{w}_t$ are nonzero. An example prompt used by the Selector Agent is illustrated in Figure 4.
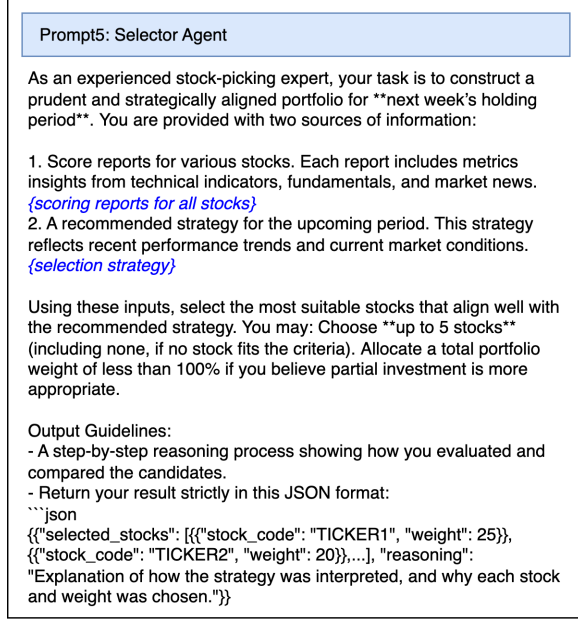
---

**Prompt5: Selector Agent**

As an experienced stock-picking expert, your task is to construct a prudent and strategically aligned portfolio for **next week's holding period**. You are provided with two sources of information:

1. Score reports for various stocks. Each report includes metrics insights from technical indicators, fundamentals, and market news. *{scoring reports for all stocks}*
2. A recommended strategy for the upcoming period. This strategy reflects recent performance trends and current market conditions. *{selection strategy}*

Using these inputs, select the most suitable stocks that align well with the recommended strategy. You may: Choose **up to 5 stocks** (including none, if no stock fits the criteria). Allocate a total portfolio weight of less than 100% if you believe partial investment is more appropriate.

Output Guidelines:
- A step-by-step reasoning process showing how you evaluated and compared the candidates.
- Return your result strictly in this JSON format:
```json
{{"selected_stocks": [{{"stock_code": "TICKER1", "weight": 25}},
{{"stock_code": "TICKER2", "weight": 20}},...], "reasoning":
"Explanation of how the strategy was interpreted, and why each stock and weight was chosen."}}
```

**Figure 4: Example prompt used by the Selector Agent for stock selection.**

## 4.4 Strategy Iteration

In this part, the *Strategy Agent* refines the current strategy based on the realized price changes of each stock and the historical trajectory of past strategies.

At the end of week $t$, following the trading rule defined in Section 3.1, we obtain the realized return $r_t^{(i)}$ for each stock $x_i$. To support strategic refinement, the *Strategy Agent* receives the realized returns $r_t^{(i)}$ of all candidate stocks and their corresponding scoring reports $s_{i,t}$. It is tasked with identifying shared characteristics among high- and low-performing stocks, and proposing updated strategies.

To ensure a stable and coherent strategy iteration process, we introduce the trajectory of past strategy updates into the input of the *Strategy Agent*. This design aims to prevent divergence from previously effective strategies, reduce unnecessary fluctuations in decision-making, and help the agent identify patterns that are indicative of long-term stable returns. The historical trajectory is denoted as:

$$\mathcal{H}t = \left\{ \text{ConcatText}\left(\pi_{t-k}, R_{t-k}^{\text{avg}}, R_{t-k}\right) \right\}_{k=1}^{K} \tag{8}$$

where $\pi_{t-k}$ denotes the strategy adopted in week $t-k$, $R_{t-k}^{\text{avg}}$ is the universe average return, and $R_{t-k}$ is the portfolio return. We set K = 10 in our experiments, enabling the review of the previous 10 weeks. Subsequently, the refined strategy for the next week can be defined as follows:

$$\pi_{t+1} = \text{agent}_{\text{strategy}}(\pi_t, \mathbf{w}_t, \mathbf{r}_t, \mathbf{s}_t, \mathcal{H}_t, \text{prompt}_{\text{strategy}}) \tag{9}$$

where $\mathbf{w}_t$ is the portfolio weight vector selected in week $t$, $\mathbf{r}_t = [r_t^{(1)}, r_t^{(2)}, \ldots, r_t^{(n)}]^\top$ is the realized return vector of all candidate stocks, and $\mathbf{s}_t = \{s_{i,t}\}_{i=1}^{n}$ denotes the set of scoring reports for all
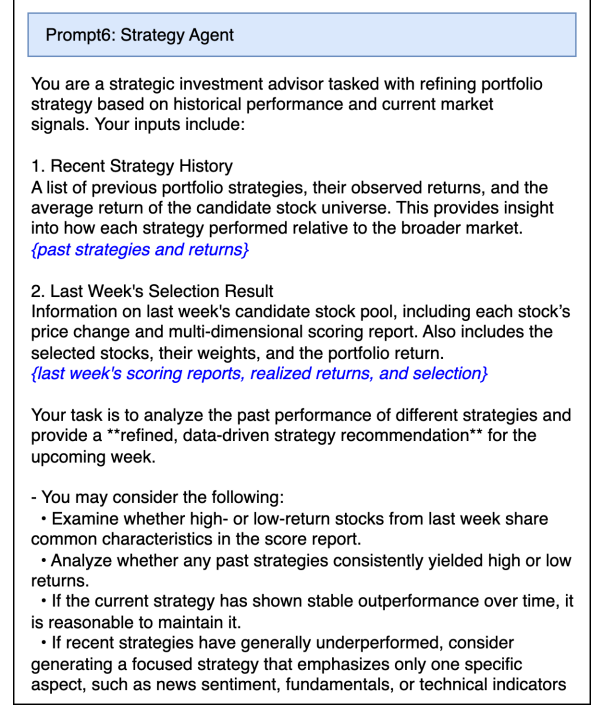
---

**Prompt6: Strategy Agent**

You are a strategic investment advisor tasked with refining portfolio strategy based on historical performance and current market signals. Your inputs include:

1. Recent Strategy History
A list of previous portfolio strategies, their observed returns, and the average return of the candidate stock universe. This provides insight into how each strategy performed relative to the broader market. *{past strategies and returns}*

2. Last Week's Selection Result
Information on last week's candidate stock pool, including each stock's price change and multi-dimensional scoring report. Also includes the selected stocks, their weights, and the portfolio return. *{last week's scoring reports, realized returns, and selection}*

Your task is to analyze the past performance of different strategies and provide a **refined, data-driven strategy recommendation** for the upcoming week.

- You may consider the following:
  • Examine whether high- or low-return stocks from last week share common characteristics in the score report.
  • Analyze whether any past strategies consistently yielded high or low returns.
  • If the current strategy has shown stable outperformance over time, it is reasonable to maintain it.
  • If recent strategies have generally underperformed, consider generating a focused strategy that emphasizes only one specific aspect, such as news sentiment, fundamentals, or technical indicators

**Figure 5: Example prompt used by the Strategy Agent for strategy iteration.**

stocks in week $t$. An example prompt used by the Strategy Agent is illustrated in Figure 5.

## 5 Experiments

In this section, we detail our experimental setup and evaluate the proposed framework by conducting portfolio construction across various stock universes. The performance of our framework is compared against several baseline models to demonstrate its effectiveness and applicability.

## 5.1 Experimental Settings

*5.1.1 Dataset Description.* We construct four distinct candidate stock universes for evaluation. For each universe, we independently conduct experiments to assess the model's performance under varying market conditions.

- **DJIA Constituents**: This set includes 30 large-cap U.S. companies from diverse sectors and serves as a representative benchmark of the broader market[2].
- **Technology Sector Stocks**: Comprising 44 constituent companies from the NASDAQ-100 Technology Sector Index[3].
- **Financial Sector Stocks**: A subset of 49 major companies with the highest weights in the SPDR Fund for the Financial Select Sector[4].

---

[2]https://en.wikipedia.org/wiki/Dow_Jones_Industrial_Average
[3]https://indexes.nasdaqomx.com/Index/Weighting/NDXT
[4]https://www.ssga.com/us/en/individual/etfs/the-financial-select-sector-spdr-fund-xlf

- **Healthcare Sector Stocks**: A selection of 46 top-weighted companies from the SPDR Fund for the Select Sector of Health Care[5].

Data are retrieved from Alpha Vantage [6], and preprocessed according to the methodology detailed in Section 3.2. The evaluation period spans from May 16, 2022 to May 27, 2024. For time-series baselines that require training, we provide historical price data covering the period from May 1, 2012 to May 15, 2022 as the training and validation set.

*5.1.2 Baselines.* We use the equal-weight strategy (1/N) as a market reference and compare our method against three representative categories of models.

- **Rule-based methods:** Including SMA (Simple Moving Average), MACD (Moving Average Convergence Divergence), and BOLL (Bollinger Bands) [38]. For each indicator, we construct a corresponding factor score and select the top 5 stocks each week, assigning 20% weight to each.
- **Deep learning prediction models:** We include the classical LSTM model [26] as well as two state-of-the-art transformer-based models: Informer [47] and Autoformer [43]. For each stock $i$ at week $t$, the response variable is the weekly return $r_{i,t}$. The input to the model includes a 4-week time window of technical features, combined with a stock-specific embedding to enable multi-stock prediction. All model hyperparameters are tuned on a validation set. At each week $t$, the top five stocks with the highest predicted returns are assigned an equal weight of 20% each.
- **Multi-LLM baselines:** We include both a single-step version and a reflective variant as described in Section 2. The single-step framework, implemented following by TradingAgent [46], directly generates portfolios from summarized information. The reflective version, following the design of CryptoTrade [31], further incorporates a *Reflection Agent* that analyzes previous portfolio returns to refine stock selection. For consistency and fair comparison, we make minor adjustments to the input and output structures of both implementations.

*5.1.3 Evaluation Metrics.* We use three metrics to evaluate portfolio performance: Accumulated Return (AR), Sharpe Ratio (SR), and Calmar Ratio (CR).

- **Accumulated Return (AR)**: AR is defined as the total compounded return over the evaluation period:

$$\text{AR} = \prod_{t=1}^{T}(1 + R_t) - 1 \tag{10}$$

where $R_t$ denotes the portfolio return at week $t$, and $T$ is the total number of evaluation weeks.

- **Sharpe Ratio (SR)**: SR measures the risk-adjusted return, capturing how efficiently a portfolio converts volatility into excess return. It is calculated as:

$$\text{SR} = \frac{\mathbb{E}[R_t]}{\sigma(R_t)} \tag{11}$$

assuming a zero risk-free rate. Here, $\mathbb{E}[R_t]$ is the mean return and $\sigma(R_t)$ is the standard deviation of returns.

- **Calmar Ratio (CR)**: CR evaluates return relative to the worst drawdown observed, providing an indication of how well a strategy balances profitability against downside risk. It is defined as:

$$\text{CR} = \frac{\text{AR}}{|\text{MDD}|} \tag{12}$$

where the maximum drawdown (MDD) is computed as:

$$\text{MDD} = \min_{t}\left(\frac{C_t - \max_{i \le t} C_i}{\max_{i \le t} C_i}\right) \tag{13}$$

and $C_t$ is the accumulated return up to week $t$.

## 5.2 Results and Analysis

We summarize the experimental results across different stock universes in a single table, as shown in Table 1. Below, we present our key findings.

*5.2.1 Overall Performance of 3S-Trader.* Our proposed method demonstrates stable performance across diverse market environments, consistently ranking among the top two models across nearly all evaluation metrics. This advantage is most evident in the mixed-sector DJIA constituents, where 3S-Trader achieves the highest accumulated return of 131.83%, significantly outperforming the second-best model. As illustrated in Figure 6, the upward trend of 3S-Trader is clearly visible. Moreover, 3S-Trader exhibits no obvious weaknesses across different performance metrics and sectors, indicating a high degree of robustness and stability in its returns.

*5.2.2 Comparison with Rule-based Models.* Our proposed method significantly outperforms traditional rule-based strategies across all sectors. In many cases, classical technical indicators such as SMA, MACD, and BOLL even underperform the simple 1/N benchmark. This underperformance can be attributed, in part, to the experimental design: to ensure fair comparison and interpretability, we use fixed-factor rule-based strategies with relatively low trading frequency. However, in real-world trading scenarios, the effectiveness of such strategies often depends on dynamic factor selection, regular parameter tuning, and more frequent position adjustments. These complexities are difficult to fully capture within our current experimental framework. By contrast, our LLM-based approach adapts more flexibly to changing market conditions and integrates diverse signals in a coherent and data-driven manner, which explains its consistent advantage over static rule-based methods.

*5.2.3 Comparison with Deep Learning Models.* Deep learning models demonstrate strong performance in terms of accumulated return, with each of the three sector-specific stock pools featuring a deep learning model as the top performer in this dimension. In particular, the classical sequence prediction model LSTM performs competitively against the more recent state-of-the-art architectures like Informer and Autoformer. In the technology sector, for instance, LSTM achieves an accumulated return of 193.39%, nearly double that of Informer (98.61%) and Autoformer (102.90%), highlighting its capacity to capture strong trends in high-momentum markets.
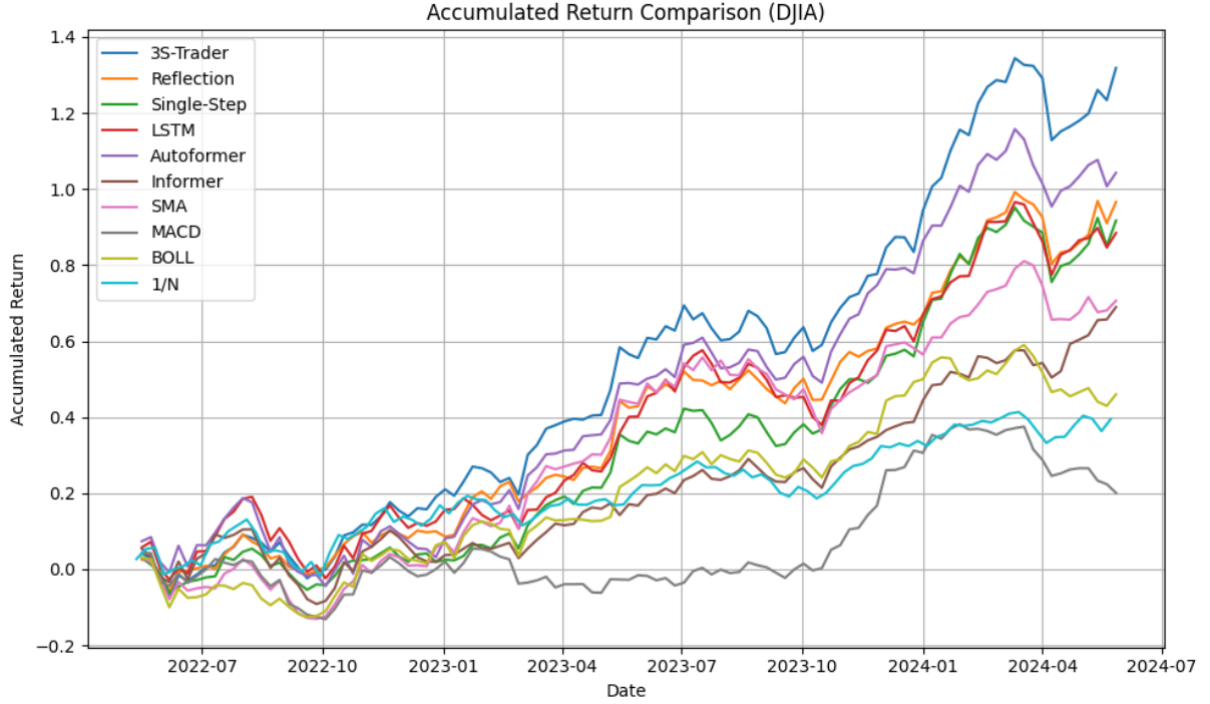
However, these models exhibit significant limitations in terms of risk control. This is reflected in both the Sharpe and Calmar Ratios,

**Table 1: Performance comparison across different models and stock universes. Red indicates the highest value in each column, and blue indicates the second highest.**

| Category | Approach | DJIA Constituents | | | Financial Sector | | | Healthcare Sector | | | Technology Sector | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | AR(%) | SR | CR | AR(%) | SR | CR | AR(%) | SR | CR | AR(%) | SR | CR |
| Benchmark | 1/N | 39.50 | 0.168 | 2.98 | 52.63 | 0.16 | 3.95 | 27.82 | 0.12 | 2.45 | 72.13 | 0.16 | 3.59 |
| Rule-based | SMA | 70.63 | 0.20 | 4.52 | 33.10 | 0.10 | 2.08 | 17.21 | 0.07 | 1.26 | 67.36 | 0.13 | 1.95 |
| | MACD | 46.00 | 0.15 | 3.02 | 61.16 | 0.15 | 3.42 | 20.58 | 0.08 | 1.29 | 116.29 | 0.19 | 4.30 |
| | BOLL | 20.09 | 0.08 | 1.20 | 25.09 | 0.06 | 1.60 | 6.82 | 0.04 | 0.51 | 96.98 | 0.16 | 5.16 |
| Deep Learning | LSTM | 88.42 | 0.23 | 4.91 | 61.24 | 0.15 | 3.54 | 59.78 | 0.15 | 3.26 | 193.39 | 0.21 | 5.81 |
| | Informer | 68.96 | 0.22 | 3.88 | 89.49 | 0.25 | 8.08 | 44.37 | 0.15 | 3.54 | 98.61 | 0.15 | 3.33 |
| | Autoformer | 104.26 | 0.24 | 5.35 | 75.92 | 0.17 | 4.50 | 30.44 | 0.11 | 2.30 | 102.90 | 0.14 | 2.97 |
| Multi-LLM | Single-Step | 91.69 | 0.27 | 8.94 | 34.92 | 0.15 | 3.10 | 31.40 | 0.12 | 3.11 | 152.84 | 0.25 | 9.04 |
| | Reflective | 96.61 | 0.28 | 9.90 | 44.86 | 0.17 | 4.56 | 43.73 | 0.18 | 6.51 | 128.27 | 0.21 | 5.41 |
| | **3S-Trader** | 131.83 | 0.31 | 11.84 | 84.93 | 0.21 | 7.57 | 51.41 | 0.17 | 3.82 | 183.29 | 0.27 | 11.81 |



**Figure 6: Comparison of accumulated returns for all methods on DJIA constituents.**

where their performance often falls short compared to LLM-based frameworks. The relatively lower risk-adjusted returns suggest that deep learning models may overfit to past trends or struggle to generalize under market volatility, leading to inconsistent or overly aggressive allocations. In contrast, multi-LLM systems show more balanced performance, with consistently higher Sharpe and Calmar Ratios. This indicates that while deep learning models are good at capturing trends, they often lack the ability to reason and handle complex information. Multi-LLM frameworks, by combining different data sources and using strategy texts, are better at making stable and informed decisions under uncertainty.

*5.2.4 Self-Refined Frameworks in Volatile Markets.* Both 3S-Trader and the *Reflective* framework are designed to refine their strategies by learning from past decisions. We focus on their performance in volatile markets such as healthcare and finance, as shown in Figure 7. In the healthcare sector, the performance gap between the two frameworks is not substantial, both demonstrate clear improvements compared to the *Single-Step* baseline. Notably, the *Reflective*
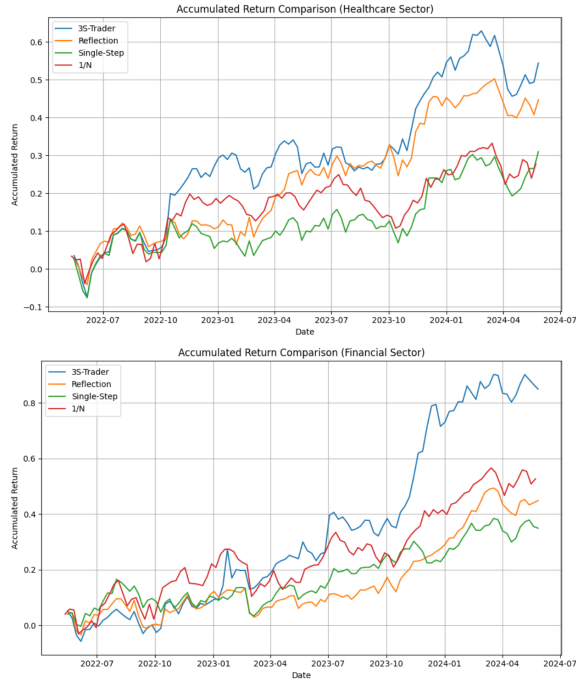
**Figure 7: Comparison of accumulated returns for 3S-Trader, Reflective and Single-Step Frameworks across Healthcare and Financial sectors.**

variant achieves the best stability in this sector, recording the highest SR of 0.18 and CR of 6.51.

However, in the financial sector, the performance gain from the reflection mechanism is less pronounced. The *Reflective* variant even underperforms both the *Single-Step* baseline and the market benchmark for most of the validation period, potentially due to over-adjustment or the absence of clear directional guidance during strategy updates. In contrast, 3S-Trader not only incorporates a reflective loop but also leverages a multi-dimensional scoring system that provides explicit and interpretable criteria for strategy refinement. This holistic design enables consistently superior and more stable performance across different market conditions.

## 6 Conclusion and Future Work

In this paper, we proposed 3S-Trader, a multi-LLM framework for portfolio construction that is capable of self-adjustment and adapts to diverse market conditions. The framework condenses recent market information into stock-level scoring reports and applies explicit, interpretable selection criteria to guide portfolio allocation. Compared with traditional rule-based and time-series forecasting models, our approach requires no model training or parameter tuning, yet achieves consistently strong performance across different stock universes, delivering competitive returns and robust stability.

For future work, several directions remain open. First, while the current scoring dimensions are designed based on domain expertise, a promising extension is to enable the automatic discovery and learning of scoring factors from data. Second, our experiments are

limited to backtesting; validating the framework in live trading environments will be essential to assess its real-world feasibility. Lastly, incorporating broader asset classes and exploring dynamic risk-control mechanisms could further enhance the generalizability and practical value of the proposed system.

## References

[1] Majid Abdulsatar, Hussain Ahmad, Diksha Goel, and Faheem Ullah. 2025. Towards deep learning enabled cybersecurity risk assessment for microservice architectures. *Cluster Computing* 28, 6 (2025), 350.

[2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[3] Hussain Ahmad and Diksha Goel. 2025. The future of ai: Exploring the potential of large concept models. *arXiv preprint arXiv:2501.05487* (2025).

[4] Hussain Ahmad, Christoph Treude, Markus Wagner, and Claudia Szabo. 2024. Smart HPA: A resource-efficient horizontal pod auto-scaler for microservice architectures. In *2024 IEEE 21st International Conference on Software Architecture (ICSA)*. IEEE, 46–57.

[5] Hussain Ahmad, Christoph Treude, Markus Wagner, and Claudia Szabo. 2025. Resilient Auto-Scaling of Microservice Architectures with Efficient Resource Management. *arXiv preprint arXiv:2506.05693* (2025).

[6] Hussain Ahmad, Christoph Treude, Markus Wagner, and Claudia Szabo. 2025. Towards resource-efficient reactive and proactive auto-scaling for microservice architectures. *Journal of Systems and Software* 225 (2025), 112390.

[7] Hussain Ahmad, Faheem Ullah, and Rehan Jafri. 2025. A survey on immersive cyber situational awareness systems. *Journal of Cybersecurity and Privacy* 5, 2 (2025), 33.

[8] Rabbia Ahmed, Sadaf Abdul Rauf, and Seemab Latif. 2024. Leveraging large language models and prompt settings for context-aware financial sentiment analysis. In *2024 5th International Conference on Advancements in Computational Sciences (ICACS)*. IEEE, 1–9.

[9] Dogu Araci. 2019. FinBERT: Financial Sentiment Analysis with Pre-trained Language Models. *CoRR* (2019).

[10] Tim Bollerslev. 1986. Generalized autoregressive conditional heteroskedasticity. *Journal of Econometrics* 31, 3 (1986), 307–327.

[11] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258* (2021).

[12] Anastasia Borovykh, Sander Bohte, and Cornelis W. Oosterlee. 2017. Conditional time series forecasting with convolutional neural networks. *arXiv preprint arXiv:1703.04691* (2017).

[13] George Edward Pelham Box and Gwilym Jenkins. 1990. *Time Series Analysis, Forecasting and Control.* Holden-Day, Inc., USA.

[14] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

[15] Shivansh Chopra, Hussain Ahmad, Diksha Goel, and Claudia Szabo. 2024. Chatnvd: Advancing cybersecurity vulnerability assessment with large language models. *arXiv preprint arXiv:2412.04756* (2024).

[16] Gregory Connor. 2000. Active portfolio management: a quantitative approach to providing superior returns and controlling risk.

[17] Yuhong Deng, Feng Bao, Yongcheng Kong, Zhiquan Ren, and Qionghai Dai. 2016. Deep direct reinforcement learning for financial signal representation and trading. *IEEE transactions on neural networks and learning systems* 28, 3 (2016), 653–664.

[18] Frank J Fabozzi, Petter N Kolm, Dessislava A Pachamanova, and Sergio M Focardi. 2007. *Robust Portfolio Optimization and Management.* John Wiley & Sons.

[19] Thomas Fischer and Christopher Krauss. 2018. Deep learning with long short-term memory networks for financial market predictions. In *European Journal of Operational Research*, Vol. 270. 654–669.

[20] Diksha Goel. 2023. Enhancing network resilience through machine learning-powered graph combinatorial optimization: Applications in cyber defense and information diffusion. *arXiv preprint arXiv:2310.10667* (2023).

[21] Diksha Goel, Hussain Ahmad, Ankit Kumar Jain, and Nikhil Kumar Goel. 2024. Machine learning driven smishing detection framework for mobile security. *arXiv preprint arXiv:2412.09641* (2024).

[22] Diksha Goel, Hussain Ahmad, Kristen Moore, and Mingyu Guo. 2025. Co-Evolutionary Defence of Active Directory Attack Graphs via GNN-Approximated Dynamic Programming. *arXiv preprint arXiv:2505.11710* (2025).

[23] Diksha Goel, Kristen Moore, Mingyu Guo, Derui Wang, Minjune Kim, and Seyit Camtepe. 2024. Optimizing cyber defense in dynamic active directories through

reinforcement learning. In *European Symposium on Research in Computer Security*. Springer, 332–352.

[24] Jingyi Gu, Junyi Ye, Guiling Wang, and Wenpeng Yin. 2024. Adaptive and explainable margin trading via large language models on portfolio management. In *Proceedings of the 5th ACM International Conference on AI in Finance*. 248–256.

[25] Mubin Ul Haque, Isuru Dharmadasa, Zarrin Tasnim Sworna, Roshan Namal Rajapakse, and Hussain Ahmad. 2022. " I think this is the most disruptive technology": Exploring Sentiments of ChatGPT Early Adopters using Twitter Data. *arXiv preprint arXiv:2212.05856* (2022).

[26] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.

[27] Raveen Kanishka Jayalath, Hussain Ahmad, Diksha Goel, Muhammad Shuja Syed, and Faheem Ullah. 2024. Microservice vulnerability analysis: A literature review with empirical insights. *IEEE Access* (2024).

[28] Zhengyao Jiang, Dixing Xu, and Jinjun Liang. 2017. A deep reinforcement learning framework for the financial portfolio management problem. *arXiv preprint arXiv:1706.10059* (2017).

[29] Alex Kim, Maximilian Muhn, and Valeri Nikolaev. 2024. Financial statement analysis with large language models. *arXiv preprint arXiv:2407.17866* (2024).

[30] Kelvin JL 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*. 4304–4315.

[31] Yuan Li, Bingqiao Luo, Qian Wang, Nuo Chen, Xu Liu, and Bingsheng He. 2024. CryptoTrade: A reflective LLM-based agent to guide zero-shot cryptocurrency trading. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 1094–1106.

[32] Xiao-Yang Liu, Guoxuan Wang, Hongyang Yang, and Daochen Zha. 2023. Fingpt: Democratizing internet-scale data for financial large language models. *arXiv preprint arXiv:2307.10485* (2023).

[33] Ziyang Liu, Meng Yu, Xueting Ren, Zhaoxin Yu, Zihan Zhang, Yuan Yang, Runzhi Zhang, Jiawei Zhou, Jing Liu, and Jie Tang. 2024. FinVision: A Multi-Agent Framework for Stock Market Prediction. *arXiv preprint arXiv:2405.18514* (2024).

[34] Harry Markowitz. 1952. Portfolio Selection. *The Journal of Finance* 7, 1 (1952), 77–91. http://www.jstor.org/stable/2975974

[35] Daniel Mebane. 2019. Technical Analysis Library in Python. Available at https://github.com/bukosabino/ta.

[36] Manale Mortaji, Azeddine Khiat, and Mohamed Benhouad. 2024. Reinforcement learning application in portfolio optimization: a comprehensive literature review. In *2024 International Conference on Intelligent Systems and Computer Vision (ISCV)*. IEEE, 1–6.

[37] Tashreef Muhammad, Anika Bintee Aftab, Muhammad Ibrahim, Md Mainul Ahsan, Maishameem Meherin Muhu, Shahidul Islam Khan, and Mohammad Shafiul Alam. 2023. Transformer-based deep learning model for stock price prediction: A case study on Bangladesh stock market. *International Journal of Computational Intelligence and Applications* 22, 03 (2023), 2350013.

[38] John J. Murphy. 1999. *Technical Analysis of the Financial Markets: A Comprehensive Guide to Trading Methods and Applications*. New York Institute of Finance.

[39] Ahmet Murat Ozbayoglu, Mehmet Ugur Gudelek, and Omer Berat Sezer. 2020. Deep learning for financial applications: A survey. *Applied soft computing* 93 (2020), 106384.

[40] Pragya Srivastava, Manuj Malik, Vivek Gupta, Tanuja Ganu, and Dan Roth. 2024. Evaluating LLMs' Mathematical Reasoning in Financial Document Question Answering. *arXiv preprint arXiv:2402.11194* (2024).

[41] Meir Statman. 1987. How many stocks make a diversified portfolio? *Journal of Financial and Quantitative Analysis* 22, 3 (1987), 353–363.

[42] Faheem Ullah, Xiaohan Ye, Uswa Fatima, Zahid Akhtar, Yuxi Wu, and Hussain Ahmad. 2025. What Skills Do Cyber Security Professionals Need? *arXiv preprint arXiv:2502.13658* (2025).

[43] Haixu Wu, Yifan Xu, Jiahui Wang, Guodong Long, Chengqi Zhang, and Lina Yao. 2021. Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. In *Advances in Neural Information Processing Systems (NeurIPS)*, Vol. 34. 22419–22430.

[44] Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao, Yilei Zhao, Xinyu Cai, et al. 2024. A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist. In *Proceedings of the 30th acm sigkdd conference on knowledge discovery and data mining*. 4314–4325.

[45] Yiyao Zhang, Hussain Ahmad, Diksha Goel, and Claudia Szabo. 2025. RegimeFolio: A Regime Aware ML System for Sectoral Portfolio Optimization in Dynamic Markets. *arXiv preprint arXiv:2510.14986* (2025).

[46] Yujia Zhao, Xiaozhong Liu, Yunkai Wang, Zihang Wang, Canwen Xu, Yixin Zhang, et al. 2023. TradingAgent: LLMs for Trading via Chain of Thought and Memory. *arXiv preprint arXiv:2311.09722* (2023).

[47] Haoyi Zhou, Shanghang Zhang, Jiehui Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang. 2021. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of AAAI*.