

STOCKBENCH: Can LLM Agents Trade Stocks Profitably In Real-world Markets?

Anonymous ACL submission

Abstract

Large language models (LLMs) demonstrate strong potential as autonomous agents, with promising capabilities in reasoning, tool use, and sequential decision-making. While prior benchmarks have evaluated LLM agents in various domains, the financial domain remains underexplored, despite its significant economic value and complex reasoning requirements. Most existing financial benchmarks focus on static question-answering, failing to capture the dynamics of real-market trading. To address this gap, we introduce STOCKBENCH, a contamination-free benchmark designed to evaluate LLM agents in realistic, multi-month stock trading environments. Agents receive daily market signals—including prices, fundamentals, and news—and make sequential buy, sell, or hold decisions. Performance is measured using financial metrics such as cumulative return, maximum drawdown, and the Sortino ratio, capturing both profitability and risk management. We evaluate a wide range of state-of-the-art proprietary and open-source LLMs. Surprisingly, most models struggle to outperform the simple buy-and-hold baseline, while some models demonstrate the potential to achieve higher returns and stronger risk management. These findings highlight both the challenges and opportunities of LLM-based trading agents, showing that strong performance on static financial question-answering do not necessarily translate into effective trading behavior. We release STOCKBENCH as an open-source benchmark to enable future research on LLM-driven financial agents.

1 Introduction

Large language models (LLMs) have enabled extensive exploration of autonomous agents, demonstrating strong capabilities in reasoning, tool use, and long-horizon decision making (OpenAI, 2024; Anthropic, 2025a; DeepMind, 2025; Liu et al., 2024; Guo et al., 2025a; Meta-AI, 2025; Yang et al.,

2024a; Bai et al., 2025; OpenAI, 2025b). These agentic capabilities have been verified by benchmarks in various domains, such as software engineering (Jimenez et al., 2024; Yang et al., 2024b), scientific discovery (Bragg et al., 2025), and marketing (Chen et al., 2025; Barres et al., 2025). Evaluations using state-of-the-art LLMs such as GPT-5 (OpenAI, 2025a) and Claude-4 (Anthropic, 2025b), highlight the potential of LLM agents to support workflow automation and productivity gains. As LLMs continue to improve, their growing agentic capabilities increasingly push applications toward real-world deployment with tangible economic value.

Among various agent application domains, the financial domain stands out with its direct link to economic value and the high stakes of decision making (Wu et al., 2023; Lee et al., 2024; Nie et al., 2024). To rigorously evaluate the profitability and risk-management capabilities of LLM agents in financial settings, an ideal benchmark should satisfy three key principles: **(1) Realistic Market Interaction.** Agents should operate in a dynamic market environment, reacting to real-time price movements and news events. **(2) Continuous Decision Making.** Agents should make sequential trading decisions over extended horizons, reflecting the iterative nature of real investment strategies. **(3) Contamination-Free Data.** To ensure fair and reliable evaluation, agents must not have prior exposure to the test data during training, requiring careful data curation and strict temporal separation.

However, existing financial benchmarks for LLM agents largely focus on static question-answering tasks (Chen et al., 2021; Zhu et al., 2021; Yin et al., 2023), which are designed to test the financial knowledge of LLMs but fail to capture the dynamics of realistic trading scenarios. Although recent efforts like INVESTORBENCH (Li et al., 2025a) take a step towards simulating trading environments, they are limited to single-stock-trading

Benchmark	Market Simulation	Multi Month Horizon	Continuous Decision	Contamination Free	Direct Economic Value
FinQA (Chen et al., 2021)	✗	✗	✗	✗	✗
ConvFinQA (Chen et al., 2022)	✗	✗	✗	✗	✗
FLUE (Shah et al., 2022)	✗	✗	✗	✗	✗
FinEval (Guo et al., 2025b)	✗	✗	✗	✗	✗
CPA-QKA (Kuang et al., 2025)	✗	✗	✗	✗	✗
BizFinBench (Lu et al., 2025)	✗	✗	✗	✗	✗
Finance Agent Benchmark (Bigeard et al., 2025)	✓	✗	✓	✗	✗
INVESTORBENCH (Li et al., 2024)	✓	✓	✓	✗	✓
FinSearchComp (Hu et al., 2025)	✗	✓	✓	✗	✓
STOCKBENCH (Ours)	✓	✓	✓	✓	✓

Table 1: Comparison of STOCKBENCH with existing financial benchmarks.

and rely on historical data prior to 2021, raising concerns about data contamination and outdated market conditions.

To bridge this gap, we propose STOCKBENCH, an evolving benchmark that places LLM agents in realistic stock-trading environments, directly evaluating their profitability and risk-management capabilities. Specifically, STOCKBENCH is designed to be: **(1) Realistic.** Agents receive daily market signals including prices, company fundamentals, and news headlines, reflecting real-world trading conditions. **(2) Continuous.** Agents make sequential daily trading decisions (buy, sell, or hold) over a multi-month horizon, mirroring the iterative nature and long-term nature of investment strategies. **(3) Contamination-Free.** The benchmark is instantiated using recent market data from March 2025 to July 2025 and will be continuously updated to prevent overlap with the training corpora of LLMs. Performance is evaluated using key financial metrics such as cumulative return, maximum drawdown, and the Sortino ratio, providing a quantitative assessment of both profitability and risk control.

As a proof of concept, we evaluate a diverse set of LLM agents, including both proprietary models (*e.g.*, GPT-5 (OpenAI, 2025a), Claude-4 (Anthropic, 2025b)) and open-weight models (*e.g.*, Qwen3 (Yang et al., 2025), Kimi-K2 (Team et al., 2025), GLM-4.5 (Zeng et al., 2025)), alongside an equal-weight buy-and-hold baseline. Surprisingly, despite their strong performance on financial QA benchmarks, most LLM agents fail to outperform this simple baseline in terms of both cumulative return and risk-adjusted return. This finding suggests that success on static financial QA does not necessarily translate into effective trading strategies in dynamic market environments, underscoring a key

challenge for LLM-based financial agents.

The main contributions of this work are summarized as follows:

- We introduce STOCKBENCH, a novel benchmark for evaluating LLM agents in a realistic stock-trading environment, measuring both profitability and risk-management capabilities.
- We propose a comprehensive evaluation framework that incorporates realistic market dynamics, diverse input data, and multiple financial metrics to holistically assess agent performance.
- We conduct extensive experiments across a range of state-of-the-art LLMs, revealing their current limitations in achieving profitable trading strategies and underscoring the need for further methodological advances.
- We open-source the implementation of STOCKBENCH to facilitate reproducibility and to encourage community contributions, fostering further research on LLM-based financial agents.

2 STOCKBENCH

The construction of STOCKBENCH consists of two main building blocks. (1) A back-trading environment, which contains historical data necessary for stock-trading decision making. We simulate real-world stock trading using this back-trading setup. (2) An associated stock-trading agent workflow. This workflow allows us to evaluate LLM backbones as agents to engage in the back-trading environment. The overall framework of STOCKBENCH is demonstrated in Figure 1.

2.1 Back-Trading Environment

We design the back-trading environment to simulate realistic stock trading, where trading agents are exposed only to data available up to the time of

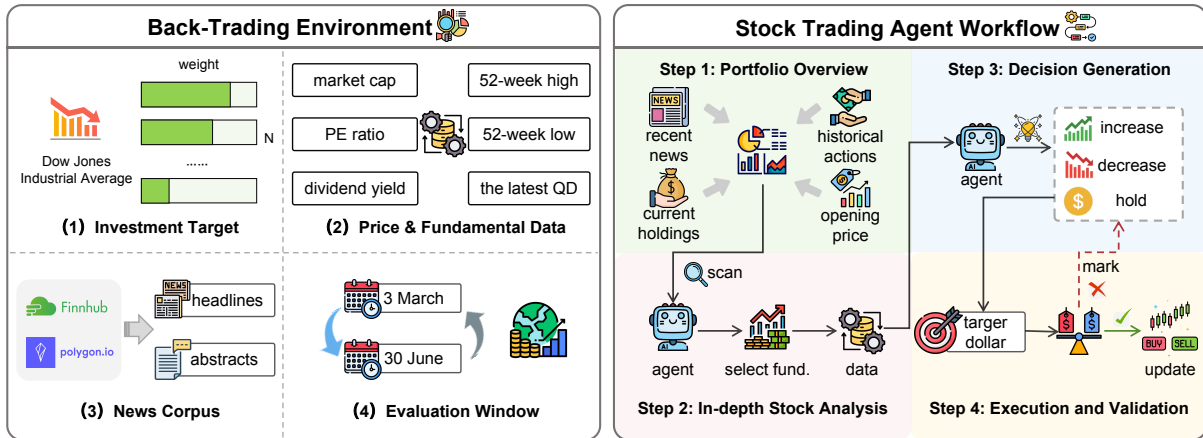


Figure 1: Overview of STOCKBENCH. The design of STOCKBENCH includes a back-trading benchmark dataset, and an associated workflow that converts backbone LLMs into agents.

each decision. To set up the environment, we identify three critical sources of information for trading decision making: (1) A bundle of investment targets, which defines the scope of the environment. We pre-define these investment targets to facilitate reproducibility of the evaluation on STOCKBENCH. (2) Historical market data, which includes both the prices and fundamental indicators. These enable the evaluated trading agents to perform quantitative analysis. (3) News corpora, which capture events that drive stock price fluctuations.

Investment Targets. The investment targets are a bundle of stocks that allow the trading agents to perform buy and sell operations. We manually select the investment targets to prevent potential outcome fluctuations caused by stock selection—*e.g.*, trading agents might otherwise happen to pick a stock driven by irrational market sentiment—thereby stabilizing the evaluation results.

To this end, we select 20 stocks from the Dow Jones Industrial Average (DJIA) with the highest weights as our investment targets. In particular, high-weighted DJIA stocks are representative of the global stock market and are less prone to short-term irrational sentiment-driven events. Constraining the trading action space to our selected investment targets mirrors real-world investor attention while keeping the dataset computationally tractable. Moreover, information about these well-known stocks is transparent and easy to collect, being readily accessible through web search engines. We show the distribution of the selected investment targets across different industries in Figure 2. Our selection covers technology, finance, and manufacturing, ensuring stock diversity.

Historical Market Data. We collect and preserve historical market data containing key quantitative information. For each stock, we use official opening prices with a concise set of fundamental metrics such as market capitalization, price-to-earnings (P/E) ratio, dividend yield, and trading range. These signals provide a reliable snapshot of company health and valuation, supporting informed decision making. We also retain the timestamps of the collected data to prevent leakage of future information to the agent.

News Corpora. We construct news corpora for stocks to enable stock-trading agents to interpret both sentiments and events in a manner that resembles how retail investors react to market narratives. For each stock, we collect news articles released within the previous 48 hours on a daily basis. These articles are retrieved using news-search API¹ with time restrictions. Since news analysis consumes substantial context length in backbone LLMs, we balance information coverage and computational cost by preserving the top five relevant news articles each time the search engine returns results.

We also carefully select the time window for collecting data in the back-trading environment. In principle, the evaluation window should satisfy two conditions: (1) the included stock information must not have been exposed to the evaluated stock-trading agents during their model training stages; and (2) the window should be sufficiently long to mitigate the impact of random noise that affects only short periods of time. To this end, we collect data spanning from March 3, 2025 to June

¹<https://finnhub.io/>

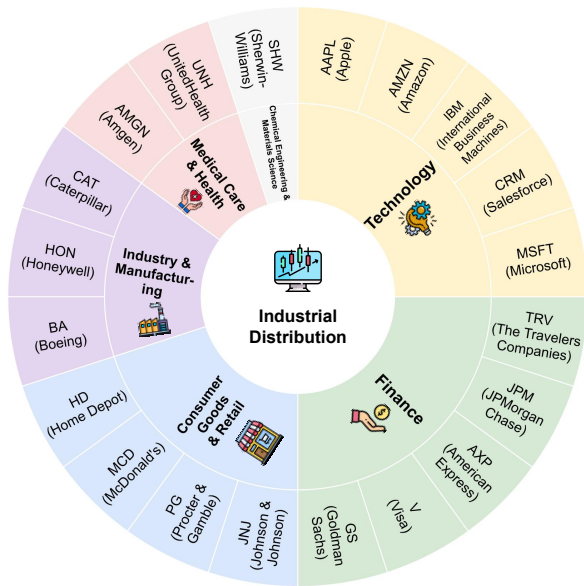


Figure 2: Industry distribution of selected stocks.

30, 2025, a four-month period that includes both volatility and trend reversals. This period also falls after the knowledge cutoff of mainstream LLMs, ensuring no data leakage. It is worth noting that we will continuously update the back-trading environment to avoid overlap with the training corpora of contemporary LLMs.

2.2 Stock-trading Agent Workflow

We provide a stock-trading agent workflow that enables backbone LLMs to interact with the back-trading environment as agents. The design of the workflow follows two goals. (1) Minimal workflow. We keep the workflow minimal, since overly complicated workflows introduce inductive biases that may favor certain backbone LLMs. (2) Realistic. We design the workflow to align with the iterative decision-making process of retail investors.

In particular, we follow previous frameworks (Zhang et al., 2020; Tsantekidis et al., 2017; Moody and Saffell, 2001; Deng et al., 2016) and organize the stock-trading workflow into four essential stages: portfolio overview, in-depth stock analysis, decision generation, and execution and validation.

Overall, the design prioritizes realism, fairness, and reproducibility, in line with earlier studies on benchmark construction for trading environments.

Step 1: Portfolio Overview. The agent first scans all available stocks in the market (the “investment target”), receiving relevant data for each stock. This includes recent news, current holdings

of the agent, historical actions, and the opening price. This step mirrors how a trader assesses the broader market and the overall status of each stock in their portfolio.

Step 2: In-Depth Stock Analysis. After the initial overview, the agent selects specific stocks for deeper analysis. For these selected stocks, the agent is provided with additional fundamental data such as market capitalization, P/E ratio, and dividend yield. This step simulates how a trader focuses on a subset of stocks identified in the initial overview, examining their financial health and other key metrics in greater depth.

Step 3: Decision Generation. With the enriched context, the agent generates decisions for each stock, choosing between three actions: (1) increase, (2) decrease, or (3) hold the position. These options ensure that the agent’s actions are clear, actionable, and executable within the constraints of a retail investor’s decision making process.

Step 4: Execution and Validation. Finally, the decisions are executed by converting dollar targets into share quantities based on the opening price. If the decisions of the agents exceed available liquidity, the system flags the issue and requires the agent to revise its decisions until they can be executed within available resources. Once validated, the new portfolio weights are locked, and the simulation advances to the next day.

2.3 Features of STOCKBENCH

We now discuss how the design of STOCKBENCH satisfies the following key principles:

Realistic Market Interaction. The design of the back-trading environment mimics real-world trading scenarios through three key elements: (1) a carefully selected bundle of investment targets, (2) reliable price and fundamental data, and (3) a concise yet timely news corpus. These elements ensure that the agent is exposed to information mirroring the complexities of real trading environments, avoiding unrealistic or overly expansive inputs.

Continuous Decision Making. In the workflow, the agent first performs a portfolio overview, then conducts in-depth stock analysis, and finally generates daily trading decisions (buy, sell, or hold) based on this analysis. These steps reflect the continuous decision-making process of retail investors,

Model	RT	DDN	Sortino	Rank
Kimi-K2	1.9	-11.8	0.0420	1
Qwen3-235B-Ins	2.4	-11.2	0.0299	2
GLM-4.5	2.3	-13.7	0.0295	3
Qwen3-235B-Think	2.5	-14.9	0.0309	4
OpenAI-O3	1.9	-13.2	0.0267	5
Qwen3-30B-Think	2.1	-13.5	0.0255	6
Claude-4-Sonnet	2.2	-14.2	0.0245	7
DeepSeek-V3.1	1.1	-14.1	0.0210	8
GPT-5	0.3	-13.1	0.0132	9
Qwen3-Coder	0.2	-13.9	0.0137	10
DeepSeek-V3	0.2	-14.1	0.0144	11
Passive Baseline	0.4	-15.2	0.0155	12
GPT-OSS-120B	-0.9	-14.0	0.0156	13
GPT-OSS-20B	-2.8	-14.4	-0.0069	14

Table 2: The performance of tested models over the evaluation period. The best performance in each metric is highlighted in bold. Models are ranked based on the z-score aggregation of all three metrics. RT stands for Final Return (%), DDN stands for Max Drawdown (%).

enabling the agent to adapt its strategies over time in response to market conditions.

Contamination-Free Data. We ensure that the agent has no prior exposure to the test data during its training. To achieve this, the benchmark is instantiated using recent market data, ensuring temporal separation and avoiding any overlap with the training corpora of contemporary LLMs.

3 Main Experiments

In this section, we present the experimental setup and results of evaluating various LLM agents within the STOCKBENCH trading workflow. We describe the trading environment, selected models, baseline strategy, and evaluation metrics. We then analyze performance outcomes, highlighting key insights into the capabilities of LLM agents in real-world financial markets.

3.1 Experiment Setup

We detail the experimental setup for evaluating LLM agents in the STOCKBENCH trading workflow. Specifically, we describe the trading environment, the models selected for benchmarking, the passive baseline, and the evaluation metrics used to assess performance.

Trading Environment. The top 20 DJIA stocks are selected as the investment targets, ensuring diverse representation across sectors. The evaluation period spans four months, from March 3 to June 30, 2025, covering 82 trading days and capturing

a range of market conditions. Each model starts with \$100,000 in cash and zero holdings, making daily trading decisions at market open. Key inputs include (1) the historical actions on held stocks over the past seven days, (2) up to five recent news articles from the previous 48 hours, and (3) for selected stocks, fundamental data such as market capitalization, P/E ratio, dividend yield, 52-week high/low, and recent quarterly dividends.

Models to Evaluate. We benchmark a diverse set of LLMs, including both open-weight models such as Qwen3 (Yang et al., 2025)², DeepSeek (Guo et al., 2025a; Liu et al., 2024), Kimi-K2 (Team et al., 2025), GLM-4.5 (Zeng et al., 2025) and GPT-OSS (OpenAI, 2024), as well as closed-source APIs like OpenAI’s O3 (OpenAI, 2025b) and Anthropic’s Claude-4-Sonnet (Anthropic, 2025b). This selection covers a range of architectures, sizes, and training methodologies to assess generality across different LLM designs. All models are equipped with 32,768 token context windows and decoded with official recommended settings to ensure their performance is optimized for the task. To hance a reliable result, each LLM agents would be run three times with different random seeds, and the average performance is reported.

Passive Baseline. As a reference point, we implement a passive equal-weight buy-and-hold strategy that allocates the initial capital equally across all selected stocks at the start of the evaluation period and holds these positions unchanged until the end. This naive allocation is a widely accepted benchmark in portfolio research, reflecting passive index tracking behavior and providing a robust lower bound against which more sophisticated active strategies can be compared (DeMiguel et al., 2009; Duchin and Levy, 2009).

Evaluation Metrics. We adopt three widely used measures in financial analysis:

Final Return. This metric measures overall profitability as the percentage change in portfolio value from the initial amount V_0 to the final amount V_T :

$$\text{Final Return} = \frac{V_T - V_0}{V_0} \quad (1)$$

It directly reflects the portfolio’s overall performance over the evaluation period and is a sim-

²Without special denote, the Qwen3 series in this papers refers to the 2507 variants

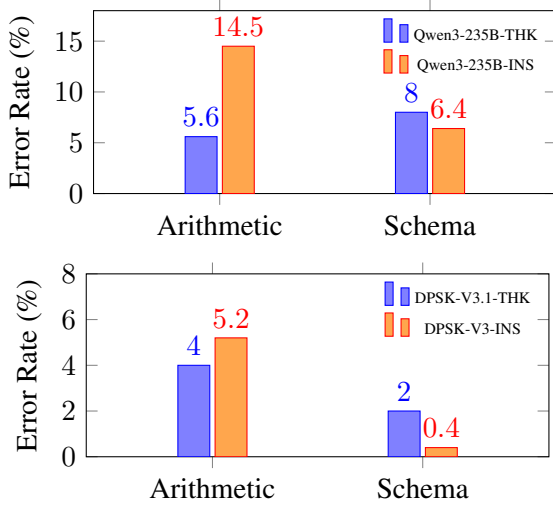


Figure 3: Error distribution (%) by type for Think vs Instruct models.

ple, widely used measure of investment profitability (Bodie et al., 2014).

Maximum Drawdown. The maximum drawdown quantifies the largest decline in portfolio value from its peak to its trough throughout the evaluation period, providing a measure for downside risk:

$$\text{Max Drawdown} = \min_{t \in [0, T]} \left(\frac{V_t - \max_{s \leq t} V_s}{\max_{s \leq t} V_s} \right) \quad (2)$$

It highlights the worst loss an investor could have faced and is commonly used to assess risk and volatility (Magdon-Ismail et al., 2004; Chekhlov et al., 2005).

Sortino Ratio. The Sortino ratio is a risk adjusted return metric that penalizes only downside volatility. It is defined as the excess return R_p divided by the downside deviation σ_d :

$$\text{Sortino Ratio} = \frac{R_p}{\sigma_d}, \quad \sigma_d = \sqrt{\frac{1}{N_d} \sum_{i=1}^{N_d} \min(R_i, 0)^2} \quad (3)$$

This metric is more appropriate than the Sharpe ratio when returns are asymmetric, as it focuses on negative volatility (Sortino and Van der Meer, 1991; Pedersen and Satchell, 2002).

After computing these metrics for each model, we derive a composite rank by leveraging the z-score of each metric, averaging them to produce a single performance score.

$$\text{Composite Rank} = \frac{z(\text{Ret}) - z(\text{DD}) + z(\text{SR})}{3} \quad (4)$$

Stocks	% Mean	% Std	CV
<i>Kimi-K2</i>			
5	-4.6	0.7	0.2
10	3.2	0.6	0.2
20	1.9	1.7	0.9
30	-0.5	1.2	2.2
<i>GPT-OSS-120B</i>			
5	-5.7	0.3	0.1
10	2.5	0.4	0.2
20	-0.4	3.9	10.2
30	-0.9	3.9	4.4

Table 3: Performance of representative models (Kimi-K2 and GPT-OSS-120B) across different investment target sizes. Results are reported as mean return (% Mean), standard deviation of returns (% Std), and coefficient of variation (CV).

Where Ret is Final Return, DD is Max Drawdown, and SR is Sortino Ratio. This approach balances profitability and risk, rewarding models that achieve high returns while effectively managing downside exposure.

3.2 Experiment Results

Table 2 presents the performance of all evaluated models over the four-month period without contamination. The results are reported across three key metrics—percentage return, maximum drawdown, and Sortino ratio—along with an overall ranking derived from a composite z-score of these metrics.

Here are the key observations: **(1) LLM agents can trade profitably in real-world markets.** Most tested models outperform the passive buy-and-hold baseline, which achieves a modest 0.4% return with a -15.2% drawdown and a Sortino ratio of 0.0155. Several agents deliver returns above 2%, with improved risk profiles. **(2) LLM agents can manage downside risk effectively.** All tested models achieve lower maximum drawdowns than the baseline, indicating that they can mitigate losses during market downturns. The best-performing agents limit drawdowns to around -11% to -14%, compared to the baseline’s -15.2%. **(3) Reasoning model does not guarantee better performance.** Although reasoning-tuned models such as Qwen3-235B-Think and Qwen3-30B-Think exhibit strong performance in tasks requiring complex reasoning, including math and coding (Yang et al., 2025), they do not consistently outperform instruction-

tuned counterparts in this trading task. For example, Qwen3-235B-Ins outperforms its reasoning-tuned version with a lower maximum drawdown (−11.2% vs. −14.9%). This suggests there is still a gap between reasoning ability and effective decision-making in dynamic, noisy environments like financial markets.

4 Analysis

4.1 Influence of Investment Target Size

To evaluate the impact of the investment target size on the agent’s performance, we study the effect of investment target size by running daily trading tasks on portfolios of 5, 10, 20, and 30 DJIA stocks, repeating each setting three times and measuring return variability. The results show that variability increases as the investment target expands. Specifically, as shown in Table 3, **(1) Scalability is inherently challenging.** All evaluated models exhibit performance degradation as the investment portfolio size increases, characterized by declining mean returns and rising return volatility. This indicates that scaling the number of tradable assets poses a non-trivial challenge for LLM agents. **(2) Model scale confers robustness.** The larger-scale model, Kimi-K2, demonstrates greater robustness to portfolio expansion, maintaining relatively stable risk-return profiles and achieving positive expected returns at moderate portfolio sizes (e.g., 10–20 stocks), whereas the smaller GPT-OSS-120B suffers from severe performance deterioration and excessive variability, suggesting that increased model capacity enhances generalization and stability in multi-asset decision-making contexts.

4.2 Influence of Error in the Trading Workflow

During trading, two common error types arise: **(1) Arithmetic Errors**, where agents miscalculate share quantities. **(2) Schema Errors**, where the LLM agent’s outputs violate the required JSON format. Figure 3 illustrates the frequency of these errors across thinking models and instruct models. The results show that thinking models make fewer arithmetic errors than instruct models, consistent with their stronger reasoning ability (Yu et al., 2025; Guo et al., 2025a; Yang et al., 2025). Yet, they also incur more schema errors, likely due to their tendency to generate overly complex outputs that deviate from the expected format (Fu et al., 2025; Li et al., 2025b).

Condition	Return (%)
Kimi-K2	1.9
w/o News	1.4
w/o News & Fund.	0.6
GPT-OSS-120B	−1.2
w/o News	−1.2
w/o News & Fund.	−3.4

Table 4: The cumulative return (CR, %) for Kimi-K2 and GPT-OSS-120B under three input settings: full input (Full), without news articles (w/o News), and without both news and fundamental data (w/o News & Fund.).

4.3 Ablation Study on Data Sources

In our workflow, LLM agents primarily rely on two main information sources: news and fundamental data. These two modalities provide complementary signals, with news capturing market sentiment and fundamentals grounding the model in key financial indicators. To better understand their respective contributions, we conduct an ablation study by progressively removing these inputs. As shown in Table 4, cumulative returns drop as news and then fundamentals are removed, confirming their importance. This behavior matches our expectation that both information sources play an important role in guiding trading decisions. The Kimi-K2 model remains relatively robust when only news is removed, but its performance deteriorates when both inputs are absent. In contrast, GPT-OSS-120B experiences a sharper decline, indicating that it relies more heavily on explicit signals provided by news and fundamentals. Overall, these results show that LLM agents effectively integrate textual and numerical inputs to guide trading decisions.

4.4 Impact of Evaluation Window

A good trading model should be able to adapt to changing market conditions over time. To investigate how the choice of evaluation window affects model rankings, we evaluate model robustness across two market time frames: a downturn period (January to April 2025) and an upturn period (May to August 2025). Figure 4 shows the models’ ranking shifts between these periods. Notably, we observe significant shifts in model rankings between the downturn and upturn periods. For example, GPT-OSS-120B moves from bottom-ranked in the downturn to top-ranked in the upturn, suggesting sensitivity to bullish markets, while Kimi-

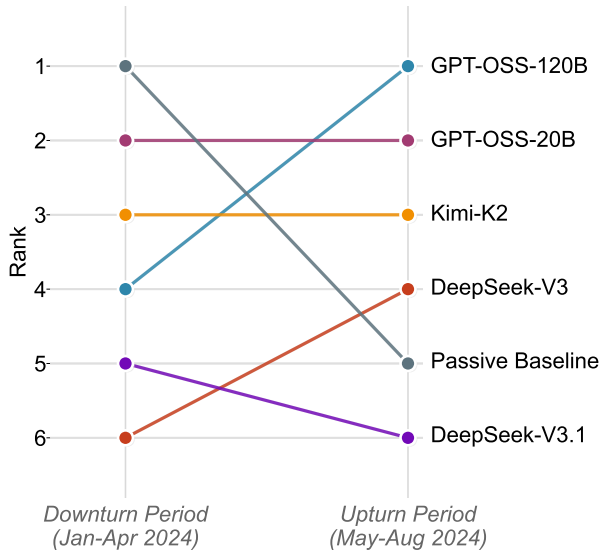


Figure 4: Model performance ranking based on the cumulative return, over two evaluation windows downturn (Jan-Apr 2024) and upturn (May-Aug 2024).

K2 remains relatively stable, indicating greater robustness to market fluctuations. This suggests that certain models may be better suited to specific market conditions, potentially due to their underlying architectures or training data. Notably, all LLM agents underperform the passive baseline during the downturn but outperform it in the upturn, highlighting a key weakness in bearish markets.

5 Related Work

5.1 LLM Agents and General Benchmarks

Large language models (LLMs) have evolved from text generators into autonomous agents capable of reasoning, planning, and interacting with environments (OpenAI, 2024; Anthropic, 2025a; DeepMind, 2025; Liu et al., 2024; Guo et al., 2025a; Meta-AI, 2025; Yang et al., 2024a; Bai et al., 2025; OpenAI, 2025b). Agentic behavior has been increasingly viewed as the next stage of LLM development, as it directly translates into real-world productivity (OpenAI, 2025a; Anthropic, 2025b). To evaluate these emerging capabilities, a range of agent benchmarks has been proposed across domains, including software engineering (SWE-Bench, SWE-Agent (Jimenez et al., 2024; Yang et al., 2024b)), scientific discovery (AstaBench (Bragg et al., 2025)), and commercial workflows (Xbench, Tau2Bench (Chen et al., 2025; Barres et al., 2025)). While these benchmarks demonstrate the potential of LLM agents for complex tasks, few existing works have exam-

ined domains where decisions translate into direct economic outcomes, such as financial trading.

5.2 Financial Agents and Benchmarks

Financial applications of LLMs have attracted growing attention due to their relevance to profitability, risk management, and high-stakes decision making (Wu et al., 2023; Lee et al., 2024; Nie et al., 2024). However, most existing benchmarks focus on static question-answering, such as FinQA, TAT-QA, and FinBench (Chen et al., 2021; Zhu et al., 2021; Yin et al., 2023). While useful for evaluating financial reasoning and domain knowledge, these benchmarks do not reflect the iterative, dynamic nature of real-world trading environments. Recent works like INVESTORBENCH (Li et al., 2025a) have begun to explore agent-based trading evaluation. However, it primarily considers single-stock settings and relies on historical data up to 2021, raising concerns about both scope and potential data contamination.

In contrast, STOCKBENCH is the first benchmark to embed LLM agents in realistic, multi-stock trading environments with continuously updated market data. By requiring LLM agents to make sequential trading decisions over extended horizons, STOCKBENCH directly evaluates both profitability and risk management capabilities. STOCKBENCH bridges the gap between static financial QA benchmarks and the practical challenges of real-world investment strategies, enabling a more faithful assessment of LLM-based financial agents.

6 Conclusion

We introduce STOCKBENCH, a benchmark for evaluating LLM agents in realistic stock-trading environments with dynamic markets and long-horizon decision making. Our experiments show that although current agents can be profitable, they rarely outperform simple baselines, revealing substantial room for improvement. We release STOCKBENCH to support future research on building more capable trading agents under complex market dynamics. We believe that STOCKBENCH will serve as a valuable resource for the research community, driving further advancements in the development of intelligent, autonomous financial agents capable of navigating complex market dynamics.

7 Limitations

STOCKBENCH evaluates LLM agents using daily trading over several months on a fixed set of large DJIA stocks. While this setup is realistic and avoids data leakage, it does not cover high-frequency trading, long-term market cycles, or a wide range of assets. Since agents can only trade once per day, intraday and event-driven strategies are not tested. The four-month window may also miss long bull or bear markets and rare but important market events, which can affect model rankings. In addition, we focus on large U.S. stocks and do not model trading costs, slippage, or liquidity limits, which makes trading easier than in real markets. We view these limitations as important directions for future work, including extending STOCKBENCH to longer horizons, higher-frequency trading, and more diverse asset classes to provide a more comprehensive assessment of LLM-based financial agents.

8 Ethical Statement

We strictly comply with all applicable financial regulations, data-protection laws, and academic ethical standards during the construction and use of STOCKBENCH. All market data (prices, fundamentals, and news) were collected through licensed data vendors or public APIs that explicitly allow research use; no non-public, insider, or personally identifiable information was accessed or stored. The benchmark is provided for academic and non-commercial research purposes only. Users are reminded that STOCKBENCH is not intended to offer, or serve as the basis for, any financial advice, trading recommendation, or commercial activity. Any trading strategy tested on STOCKBENCH carries inherent market risk; past performance recorded in the benchmark does not guarantee future returns.

References

Anthropic. 2025a. [Claude 3.7 Sonnet](#).

Anthropic. 2025b. Introducing claude 4. <https://www.anthropic.com/news/claude-4>. Accessed: 2025-09-14.

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. [Qwen2.5-VL technical report](#). *arXiv preprint arXiv:2502.13923*, abs/2502.13923.

Victor Barres and 1 others. 2025. [\$\tau^2\$ -bench: Evaluating conversational agents in a dual-control setting](#). *arXiv preprint arXiv:2506.07982*.

Antoine Bigeard, Langston Nashold, Rayan Krishnan, and Shirley Wu. 2025. [Finance agent benchmark: Benchmarking llms on real-world financial research tasks](#). *arXiv preprint arXiv:2508.00828*.

Zvi Bodie, Alex Kane, and Alan J. Marcus. 2014. *Investments*, 10th edition. McGraw-Hill Education.

Jonathan Bragg, Mike D’Arcy, Nishant Balepur, Dan Bareket, Bhavana Dalvi, Sergey Feldman, Dany Hadad, Jena D Hwang, Peter Jansen, Varsha Kishore, and 1 others. 2025. [Astabench: Rigorous benchmarking of ai agents with a scientific research suite](#). *arXiv preprint arXiv:2510.21652*.

Alexei Chekhlov, Stanislav Uryasev, and Michael Zabarankin. 2005. Drawdown measure in portfolio optimization. *International Journal of Theoretical and Applied Finance*, 8(1):13–58.

Kaiyuan Chen, Yixin Ren, Yang Liu, Xiaobo Hu, Haotong Tian, Fangfu Liu, Haoye Zhang, Hongzhang Liu, Yuan Gong, Chen Sun, Han Hou, Hui Yang, James Pan, Jianan Lou, Jiayi Mao, Jizheng Liu, Jinpeng Li, Kangyi Liu, Kenkun Liu, and 13 others. 2025. [xbench: Tracking agents productivity scaling with profession-aligned real-world evaluations](#). *arXiv preprint arXiv:2506.13651*, abs/2506.13651.

Zhiyu Chen, Wenhui Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. 2021. [Finqa: A dataset of numerical reasoning over financial data](#). *Proceedings of EMNLP 2021*.

Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. 2022. [Convinqa: Exploring the chain of numerical reasoning in conversational finance question answering](#). *arXiv preprint arXiv:2210.67890*.

Google DeepMind. 2025. [Gemini 2.5](#).

Victor DeMiguel, Lorenzo Garlappi, and Raman Uppal. 2009. Optimal versus naive diversification: How inefficient is the 1/n portfolio strategy? *Review of Financial Studies*, 22(5):1915–1953.

Yue Deng, Feng Bao, Youyong 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):653–664.

Ran Duchin and Haim Levy. 2009. Markowitz versus the talmudic portfolio diversification strategies. *Journal of Portfolio Management*, 35(2):71–84.

Tingchen Fu, Jiawei Gu, Yafu Li, Xiaoye Qu, and Yu Cheng. 2025. [Scaling reasoning, losing control: Evaluating instruction following in large reasoning models](#). *arXiv preprint arXiv:2505.14810*, abs/2505.14810.

- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. DeepSeek-R1: Incentivizing reasoning capability in LLMs via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Xin Guo, Haotian Xia, Zhaowei Liu, Hanyang Cao, Zhi Yang, Zhiqiang Liu, Sizhe Wang, Jinyi Niu, Chuqi Wang, Yanhui Wang, and 1 others. 2025b. **Fineval: A chinese financial domain knowledge evaluation benchmark for large language models**. *arXiv preprint arXiv:2504.06258*.
- Liang Hu, Jianpeng Jiao, Jiashuo Liu, Yanle Ren, Zhoufutu Wen, Kaiyuan Zhang, Xuanliang Zhang, Xi-ang Gao, Tianci He, Fei Hu, Yali Liao, Zaiyuan Wang, Chenghao Yang, Qianyu Yang, Mingren Yin, Zhiyuan Zeng, Ge Zhang, Xinyi Zhang, Xiying Zhao, and 4 others. 2025. **Finsearchcomp: Towards a realistic, expert-level evaluation of financial search and reasoning**. *arXiv preprint arXiv:2509.13160*.
- Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2024. **SWE-bench: Can language models resolve real-world github issues?** *arXiv preprint arXiv:2310.06770*.
- Ziyan Kuang, Feiyu Zhu, Maowei Jiang, Yanzhao Lai, Zelin Wang, Zhitong Wang, Meikang Qiu, Jiajia Huang, Min Peng, Qianqian Xie, and Sophia Ananiadou. 2025. From scores to skills: A cognitive diagnosis framework for evaluating financial large language models. *arXiv preprint arXiv:2508.13491*.
- Jean Lee, Nicholas Stevens, Soyeon Caren Han, and Minseok Song. 2024. **A survey of large language models in finance (finllms)**. *arXiv preprint arXiv:2402.02315*, abs/2402.02315.
- Haohang Li, Yupeng Cao, Yangyang Yu, Shashidhar Reddy Javaji, Zhiyang Deng, Yueru He, Yuechen Jiang, Zining Zhu, K. P. Subbalakshmi, Jimin Huang, Lingfei Qian, Xueqing Peng, Qianqian Xie, and Jordan W. Suchow. 2024. **Investorbench: A benchmark for financial decision-making tasks with llm-based agent**. *arXiv preprint arXiv:2412.18174*.
- Haohang Li, Yupeng Cao, Yangyang Yu, Shashidhar Reddy Javaji, Zhiyang Deng, Yueru He, Yuechen Jiang, Zining Zhu, K.p. Subbalakshmi, Jimin Huang, Lingfei Qian, Xueqing Peng, Jordan W. Suchow, and Qianqian Xie. 2025a. **INVESTORBENCH: A benchmark for financial decision-making tasks with LLM-based agent**. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2509–2525.
- Xiaomin Li, Zhou Yu, Zhiwei Zhang, Xupeng Chen, Ziji Zhang, Yingying Zhuang, Narayanan Sadagopan, and Anurag Beniwal. 2025b. **When thinking fails: The pitfalls of reasoning for instruction-following in llms**. *arXiv preprint arXiv:2505.11423*, abs/2505.11423.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. **DeepSeek-V3 technical report**. *arXiv preprint arXiv:2412.19437*, abs/2412.19437.
- Guilong Lu, Xuntao Guo, Rongjunchen Zhang, Wenqiao Zhu, and Ji Liu. 2025. **Bizfinbench: A business-driven real-world financial benchmark for evaluating llms**. *arXiv preprint arXiv:2505.19457*.
- Malik Magdon-Ismail, Amir F. Atiya, Anurag Pratap, and Yaser S. Abu-Mostafa. 2004. On the maximum drawdown of a brownian motion. *Journal of Applied Probability*, 41(1):147–161.
- Meta-AI. 2025. **The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation**.
- John Moody and Matthew Saffell. 2001. Learning to trade via direct reinforcement. *IEEE Transactions on Neural Networks*, 12(4):875–889.
- Yuqi Nie, Yaxuan Kong, Xiaowen Dong, John M. Mulvey, H. Vincent Poor, Qingsong Wen, and Stefan Zohren. 2024. **A survey of large language models for financial applications: Progress, prospects and challenges**. *arXiv preprint arXiv:2406.11903*, abs/2406.11903.
- OpenAI. 2024. **Hello GPT-4o**.
- OpenAI. 2024. Introducing gpt-oss. <https://openai.com/index/introducing-gpt-oss/>. Accessed: 2025-09-14.
- OpenAI. 2025a. Gpt-5 is here. <https://openai.com/gpt-5/>. Accessed: 2025-09-14.
- OpenAI. 2025b. **Introducing openai o3 and o4-mini**.
- Carsten S. Pedersen and Stephen Satchell. 2002. On the foundation of performance measures under asymmetric returns. *Quantitative Finance*, 2(3):217–223.
- Raj Shah, Kunal Chawla, Dheeraj Eidnani, Agam Shah, Wendi Du, Sudheer Chava, Natraj Raman, Charese Smiley, Jiaao Chen, and Diyi Yang. 2022. When flue meets flang: Benchmarks and large pretrained language model for financial domain. *arXiv preprint arXiv:2210.12345*.
- Frank A. Sortino and Robert Van der Meer. 1991. Downside risk. *Journal of Portfolio Management*, 17(4):27–31.
- Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, and 1 others. 2025. **Kimi k2: Open agentic intelligence**. *arXiv preprint arXiv:2507.20534*.
- Avraam Tsantekidis, Nikolaos Passalis, Anastasios Tefas, Juho Kannianen, Moncef Gabbouj, and Alexandros Iosifidis. 2017. Forecasting stock prices from the limit order book using convolutional neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 29(9):2856–2870.

Shijie Wu, Ozan Irsoy, Steven Lu, and 1 others. 2023. [Bloomberggpt: A large language model for finance](#). *arXiv preprint arXiv:2303.17564*. Accessed: 2025-09-14.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024a. [Qwen2.5 technical report](#). *arXiv preprint arXiv:2412.15115*, abs/2412.15115.

John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan, and Ofir Press. 2024b. [SWE-agent: Agent-computer interfaces enable automated software engineering](#). *arXiv preprint arXiv:2405.15793*.

Yuwei Yin, Yazheng Yang, Jian Yang, and Qi Liu. 2023. [Finpt: Financial risk prediction with profile tuning on pretrained foundation models](#). *arXiv preprint arXiv:2308.00065*, abs/2308.00065.

Qiyong Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, and 1 others. 2025. [Dapo: An open-source llm reinforcement learning system at scale](#). *arXiv preprint arXiv:2503.14476*, abs/2503.14476.

Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, and 1 others. 2025. [Glm-4.5: Agentic, reasoning, and coding \(arc\) foundation models](#). *arXiv preprint arXiv:2508.06471*, abs/2508.06471.

Yue Zhang, Stefan Zohren, and Stephen Roberts. 2020. Deep reinforcement learning for trading. *The Journal of Financial Data Science*, 2(2):25–40.

Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. [TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3277–3287. Association for Computational Linguistics.

A Prevent Data Leakage

In this study, we minimize the risk of data leakage by carefully planning and evaluating the time frame. When testing large language models (LLMs) in the financial field, a potential concern is that during the training process, the model will learn a lot of

past financial knowledge, which may lead to the model’s performance being artificially exaggerated. For instance, when asking GPT-5 (without using the search function), we found that the model could accurately predict the stock trend of AAPL in 2021, and the model’s response was consistent with the facts.

This discovery indicates that if the evaluation time is relatively early, the model may have obtained future information that could not have been reasonably acquired at the time of evaluation. In view of this, we have decided to limit the data used for evaluation to a more recent time frame, thereby minimizing the possibility of such “data leakage” and ensuring that the model is tested more fairly. By focusing on a narrow evaluation time window, we aim to simulate real-world scenarios where agents can only make trading decisions based on the publicly available information at the time of each decision.

This approach conforms to the best practices of financial model evaluation, ensuring that the evaluation results truly reflect the predictive and decision-making capabilities of LLM agents without being disturbed by the unintentional availability of future data

B Model Return Variance

Table 5: Model Return Variance Across Different Models. This table presents the variance of model returns for various LLMs.

Rank	Model	Var ($\times 10^{-4}$)
1	<i>DeepSeek-V3</i>	0.074
2	<i>DeepSeek-V3.1</i>	0.203
3	<i>GPT-5</i>	0.210
4	<i>Claude-4-Sonnet</i>	0.153
5	<i>GLM-4.5</i>	0.099
6	<i>Qwen3-30B-Think</i>	0.115
7	<i>Qwen3-235B-Think</i>	0.321
8	<i>Qwen3-235B-Ins</i>	0.281
9	<i>Qwen3-4B-Ins</i>	1.382
10	<i>GPT-OSS-20B</i>	1.337
11	<i>Qwen3-Coder</i>	1.655
12	<i>Openai-O3</i>	3.250
13	<i>Kimi-K2</i>	1.866
14	<i>GPT-OSS-120B</i>	10.19

In this section, we analyze the return variances of different models. Models with higher return variances may exhibit more unpredictable behaviors,

895 which is undesirable in many real-world applica-
896 tions, especially in high-risk environments such as
897 financial decision-making.

898 We ranked several large language models
899 (LLMs) based on their return variances, as shown
900 in table 5. In the evaluated model, *DeepSeek-V3*
901 exhibited the smallest performance fluctuation, in-
902 dicated high stability. In contrast, *GPT-OSS-120B*
903 exhibits the highest return variance, indicating a
904 volatility in its performance.

905 **C The Use of Large Language Models**

906 We use LLMs for two purposes. (1) Code imple-
907 mentation. When implementing the code for this
908 paper, including data gathering and experiment im-
909 plementation, we use LLMs in the form of `copilot`
910 to complete code snippets. The architecture design
911 is conducted by human researchers. (2) Proofread-
912 ing. To fix grammar issues, we use LLMs as a
913 writing tools to refine the draft.

914 We would like to highlight that LLMs are not
915 responsible for creativity tasks during conducting
916 the research of this paper, including but not limited
917 to: ideation, experiment design, paper organizing.