

# LiveTradeBench: Seeking Real-World Alpha with Large Language Models

Haofei Yu, Fenghai Li, Jiaxuan You

University of Illinois, Urbana-Champaign

 [trade-bench.live](https://trade-bench.live)

 [github.com/ulab-uiuc/live-trade-bench](https://github.com/ulab-uiuc/live-trade-bench)

Large language models (LLMs) achieve strong performance across benchmarks—from knowledge quizzes and math reasoning to web-agent tasks—but these tests occur in static settings, lacking real dynamics and uncertainty. Consequently, they evaluate isolated reasoning or problem solving rather than decision-making under uncertainty. To address this, we introduce LiveTradeBench, a live trading environment for evaluating LLM agents in realistic and evolving markets. LiveTradeBench follows three design principles: (i) Live data streaming of market prices and news, eliminating dependence on offline backtesting and preventing information leakage while capturing real-time uncertainty; (ii) a portfolio-management abstraction that extends control from single-asset actions to multi-asset allocation, integrating risk management and cross-asset reasoning; and (iii) multi-market evaluation across structurally distinct environments—U.S. stocks and Polymarket prediction markets—differing in volatility, liquidity, and information flow. At each step, an agent observes prices, news, and its portfolio, then outputs percentage allocations that balance risk and return. Using LiveTradeBench, we run 50-day live evaluations of 21 LLMs across families. Results show that (1) high LMArena scores do not imply superior trading outcomes; (2) models display distinct portfolio styles reflecting risk appetite and reasoning dynamics; and (3) some LLMs effectively leverage live signals to adapt decisions. These findings expose a gap between static evaluation and real-world competence, motivating benchmarks that test sequential decision making and consistency under live uncertainty.

## 1. Introduction

Large language models (LLMs) have achieved near-saturation performance on diverse benchmarks—such as knowledge quizzes (Hendrycks et al., 2020; Phan et al., 2025; Rein et al., 2024), math reasoning tests (Cobbe et al., 2021; Contributors, 2023; Quan et al., 2025), and instruction-following tasks (Jiang et al., 2023; Pyatkin et al., 2025; Zhou et al., 2023a). However, these benchmarks are static, evaluating models on fixed inputs with single-turn reasoning. High scores on such tests do not necessarily reflect real-world intelligence, where agents must perceive, act, and adapt through feedback over time.

To move beyond static evaluation, recent work has introduced interactive environments that allow LLMs to perform sequential actions and observe feedback (Jimenez et al., 2023; Zhou et al., 2023c). Examples include web and computer-use agents (He et al., 2024; Koh et al., 2024a; Xie et al., 2024; Zhou et al., 2023b), which operate in discrete and deterministic environments—each action produces a predictable transition defined by backend logic. These environments test perception and reasoning, but remain fully controllable and support tree-based searching (Aksitov et al., 2023; Koh et al., 2024b; Putta et al., 2024). In contrast, trading environments represent continuous and autonomous systems. The world evolves independently of the agent, and actions only adjust the agent’s internal portfolio state

**Table 1: Comparison of LiveTradeBench with existing trading benchmarks.** We compare our work with others through four dimensions: (1) *sequential decision* for whether its current trading actions rely on the previous actions; (2) *portfolio management* for whether its task is multi-asset portfolio management; (3) *live trading* for whether the evaluation belongs to backtest with historical market data or live test with real-time streaming data; (4) *multi-market evaluation* for whether it includes markets beyond the stock market.

Benchmark	Sequential Decision	Portfolio Management	Live Trading	Multi-market Evaluation
FinQA (Chen et al., 2021)	✗	✗	✗	✗
ConvFinQA (Chen et al., 2022)	✗	✗	✗	✗
FLUE (Shah et al., 2022)	✗	✗	✗	✗
FinEval (Zhang et al., 2023)	✗	✗	✗	✗
BizFinBench (Bigeard et al., 2025)	✗	✗	✗	✗
FinAgentBench (Bigeard et al., 2025)	✗	✗	✗	✗
FinSearchComp (Hu et al., 2025)	✓	✗	✗	✗
INVESTORBENCH (Li et al., 2024)	✓	✗	✗	✗
StockBench (Chen et al., 2025)	✓	✓	✗	✗
DeepFund (Li et al., 2025b)	✓	✓	✓	✗
<b>LiveTradeBench (ours)</b>	✓	✓	✓	✓

rather than directly determining future observations. Feedback is delayed and noisy, emphasizing adaptation over control. This difference in environment structure—from deterministic systems to dynamic processes—defines a deeper frontier for evaluating LLM agents’ ability to reason and act in open-ended, real-world settings (Garrido-Merchán et al., 2024; Li et al., 2024).

Despite this importance, current applications for building LLM-based trading agents remain oversimplified and disconnected from live market dynamics. Specifically, (1) most evaluation frameworks rely on offline backtesting, which is prone to information leakage and fails to capture the uncertainty, volatility, and feedback of real-world environments (Li et al., 2025a,b,c; Papadakis et al., 2025); and (2) most trading agents model trading as low-level local actions (*e.g.*, buy/sell/hold) on a single asset, neglecting higher-level reasoning and planning across multiple assets (Briola et al., 2021; Han et al., 2023a,b; Ma et al., 2025). This naturally raises a broader question: *How can we effectively evaluate the trading ability of LLM-based agents under realistic market conditions at low cost?*

To answer this question, we introduce LiveTradeBench, a live trading environment designed to address both limitations above. (1) LiveTradeBench streams live market data, financial news, and social signals, eliminating the dependence on offline backtesting and thereby avoiding information leakage from the root while capturing real-world uncertainty and feedback. (2) It adopts the portfolio management abstraction, framing trading as a strategic allocation process that integrates risk management, temporal reasoning, and decision consistency across multiple assets (Gu et al., 2024; Kou et al., 2024; Yu et al., 2024). At each step, the environment exposes dynamic observations—market conditions, contextual signals, and the agent’s historical decisions—and the LLM must output an updated portfolio allocation that balances risk and return over time. By combining live data streaming with portfolio-level reasoning, LiveTradeBench offers a realistic, end-to-end platform to evaluate the true trading competence of LLM-based agents under evolving market dynamics.

Using this benchmark, we conduct two types of live trading evaluations: stock market (U.S.

stocks) trading and prediction market (Polymarket<sup>\*</sup>) betting. We compare 21 mainstream LLMs across multiple model families and capability tiers. Our analysis yields three key findings: (1) State-of-the-art models in LMArena (Chiang et al., 2024) do not exhibit state-of-the-art trading performance—high benchmark scores in general reasoning do not translate to superior trading outcomes; (2) LLMs display distinct portfolio management styles, differing in their risk appetite, asset selection patterns, and allocation dynamics; and (3) LLMs can effectively leverage real-time market and news signals to make more informed and adaptive trading decisions. Together, these results reveal a disconnect between conventional LLM evaluation and real-world financial competence, motivating the development of more adaptive and robust portfolio management agents.

## 2. Related Work

**Evaluation of trading agents** The evaluation of LLM-based trading agents generally relies on three types of environments or benchmarks. (1) Backtesting with historical market data is the mainstream approach (Li et al., 2025c; Tang et al., 2025; Tian et al., 2025; Xiao et al., 2024a). However, such evaluations often suffer from information leakage (Li et al., 2025a,d) and poor generalization across longer or multi-regime market periods (Gao et al., 2024a; Jiang and Zhou, 2025). To address these issues, several studies propose data contamination audits, entity anonymization, and temporal de-biasing protocols for fairer backtesting evaluation (He and Xu, 2024; Wu and Zhang, 2025). (2) Market simulators provide an alternative by constructing synthetic or self-designed trading environments (Chen et al., 2023; Emmanoulopoulos et al., 2025; Lopez-Lira, 2025; Papadakis et al., 2025; Zhang et al., 2024a). Yet, these simulators serve mainly as testbeds for behavioral analysis rather than producing realistic trading actions aligned with actual market dynamics. (3) Live evaluation with real-time data represents an emerging direction. While widely explored in other domains such as question answering (Kasai et al., 2022; Nie et al., 2025) and coding (Liang and Zhang, 2024), this approach remains largely unexplored in trading (Li et al., 2025b). Our work focuses on this live evaluation paradigm, which we argue offers the most faithful and future-proof assessment of LLM trading intelligence.

**Action space design for trading agents** The design of trading tasks varies substantially across objectives, which can be formalized through differences in the *action space* of trading agents. In stock markets, most LLM-based systems adopt a *single-asset trading* formulation, where actions are discrete decisions such as buy, hold, or sell (Gao et al., 2024b; Li et al., 2023; Ma et al., 2025; Zhang et al., 2024b, 2025). While intuitive, this setup overlooks cross-asset dependencies and realistic portfolio interactions. Other approaches focus on *alpha prediction*, producing continuous vectors of alpha signals that represent expected excess returns or relative performance across assets (Heinrich et al., 2021; Islam, 2025; Sun et al., 2024; Zhang et al., 2020). However, these signals describe predictions rather than directly executable trading actions. In betting markets, agents often output probability estimates for mutually exclusive outcomes (e.g., “Yes” vs. “No”) (DeHaven et al., 2024; Jumadinova and Dasgupta, 2011; Koning and Zijm, 2022), which can be interpreted as implicit portfolio positions in complementary assets. We unify these perspectives under a *portfolio management* framework, where the agent outputs allocation ratios across multiple assets or outcomes (Lucarelli and Borrotti, 2020; Sun et al., 2021; Ye et al., 2020). This formulation generalizes discrete trading, alpha prediction, and probabilistic betting within a single continuous decision space that naturally emphasizes risk–return trade-offs and inter-asset correlations.

---

<sup>\*</sup><https://polymarket.com/>

**Framework for LLM-based trading agents** Various frameworks leverage LLMs to build trading agents in different styles. One line of research focuses on fine-tuning a single LLM with reinforcement learning (RL) to enhance decision-making and trading performance (Koa et al., 2024; Wang et al., 2025; Xiong et al., 2025; Zha and Liu, 2025; Zhang et al., 2025). Another line explores multi-agent systems, where agents collaborate or compete through role differentiation to simulate realistic market dynamics (Li and Zhao, 2025; Xu et al., 2024; Zhang and Wang, 2025). In addition, capabilities such as tool use (e.g., API calls, data collectors) (Islam, 2025; Papadakis et al., 2025), self-reflection (Koa et al., 2024), and memory (Li et al., 2024, 2023; Yu et al., 2023) have been recognized as key components for improving trading intelligence. To provide a controlled yet extensible setup, we adopt a React-style (Yao et al., 2022) framework equipped with tool use and memory as our agent configuration.

### 3. Building Live Trading Environment for Portfolio Management

#### 3.1. Definition of Portfolio Management

**Problem definition** We formulate the portfolio management task as a partially observable Markov decision process (POMDP)  $\mathcal{E} = \langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \Omega \rangle$ , where  $\mathcal{S}$  is the latent market state space,  $\mathcal{A}$  the action space,  $\mathcal{O}$  the observation space,  $\mathcal{T}$  the transition dynamics, and  $\Omega$  the observation emission function. At each timestep  $t$ , the environment is in a latent state  $\mathbf{s}_t \in \mathcal{S}$ , which encapsulates the true market condition, including asset fundamentals, volatility, liquidity, and other unobserved factors. The agent receives a partial observation

$$\mathbf{o}_t = (\mathbf{q}_t, \mathbf{p}_t, \mathbf{c}_t) = \Omega(\mathbf{s}_t), \quad (1)$$

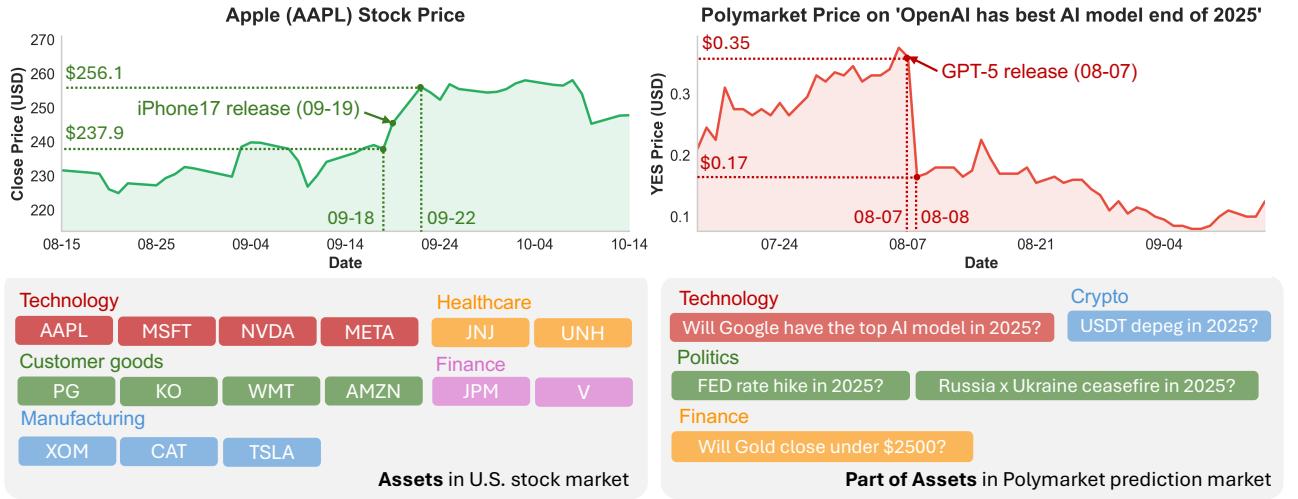
where  $\mathbf{q}_t \in \mathbb{R}^N$  denotes the current asset holdings (including cash),  $\mathbf{p}_t \in \mathbb{R}^N$  the observable market prices, and  $\mathbf{c}_t$  contextual signals such as news, sentiment, or macro indicators. The total portfolio value is computed as  $v_t = \mathbf{q}_t^\top \mathbf{p}_t$ . Conditioned on the observation history  $\mathbf{o}_{\leq t}$ , the agent produces an action  $\mathbf{a}_t \in \mathcal{A}$  representing a target allocation vector, subject to  $\sum_i a_t^{(i)} = 1$ .

**State transition function** The environment transition captures the joint evolution of the market and the agent’s portfolio. It consists of two coupled processes: an *exogenous* market-state evolution, governed by real-world dynamics and observable as  $(\mathbf{p}_t, \mathbf{c}_t) \rightarrow (\mathbf{p}_{t+1}, \mathbf{c}_{t+1})$ , and an *endogenous* portfolio adjustment induced by the agent’s allocation decision  $\mathbf{a}_t$  under the new market state, leading to  $\mathbf{q}_t \rightarrow \mathbf{q}_{t+1}$ . Concretely, after executing  $\mathbf{a}_t$ , the market evolves according to  $\mathcal{T}$ , producing new prices  $\mathbf{p}_{t+1}$  and contextual signals  $\mathbf{c}_{t+1}$ . The portfolio is revalued and rebalanced under the new prices as

$$v_{t+1}^- = \mathbf{q}_t^\top \mathbf{p}_{t+1}, \quad \mathbf{q}_{t+1} = v_{t+1}^- \frac{\mathbf{a}_t}{\mathbf{p}_{t+1}}, \quad v_{t+1} = \mathbf{q}_{t+1}^\top \mathbf{p}_{t+1} = v_{t+1}^- \quad (2)$$

where division is element-wise. The next observation  $\mathbf{o}_{t+1}$  is emitted by  $\Omega$ , capturing the updated portfolio state and market context. Since we focus on highly liquid assets such as Nvidia (NVDA) stocks and trending Polymarket markets to make up the portfolio, the agent’s individual trades exert negligible influence on real-world prices. This assumption justifies modeling the simulated transition function  $\mathcal{T}$  as closely aligned with the real world.

**From allocation decisions to executable trading actions** At each timestep  $t$ , after outputting the allocation decision  $\mathbf{a}_t$ , the agent can update its holdings from  $\mathbf{q}_t$  to  $\mathbf{q}_{t+1}$  through executable trading actions (BUY/SELL/HOLD). The executed trade vector is defined as  $\Delta \mathbf{q}_t = \mathbf{q}_{t+1} - \mathbf{q}_t$ , where a



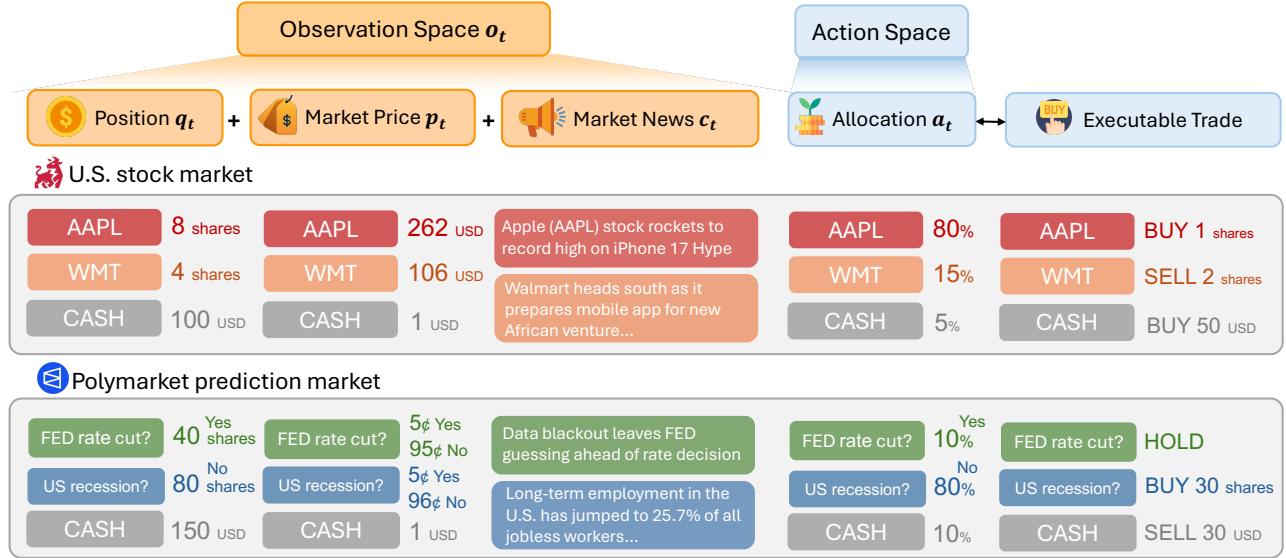
**Fig. 1: Market selection in LiveTradeBench.** The top panels show AAPL in the U.S. stock market (left) and the contract “OpenAI has the best AI model by the end of 2025” in the Polymarket prediction market (right). In prediction markets, the price directly reflects the probability of a given outcome. Both markets respond to news and historical price trends, but Polymarket exhibits sharper fluctuations, faster reactions, and higher sensitivity to external signals. The bottom panels display representative assets across various domains, including technology, finance, cryptocurrency, manufacturing, and politics.

positive  $\Delta q_t^{(i)} > 0$  indicates buying  $\Delta q_t^{(i)}$  shares of asset  $i$ , a negative  $\Delta q_t^{(i)} < 0$  indicates selling  $|\Delta q_t^{(i)}|$  shares, and  $\Delta q_t^{(i)} = 0$  corresponds to holding the current position. Once these trades are executed, the portfolio transitions to the new holdings  $\mathbf{q}_{t+1}$ , and the total portfolio value  $v_{t+1}$  is updated according to Eq. 2. This formulation provides a direct mapping from the high-level allocation action space to explicit buy, sell, and hold operations, without modeling low-level order execution mechanics.

### 3.2. Market Selection

We evaluate agents in two complementary environments (stock market and prediction market) designed to test their generalization across both structured and information-driven regimes. This dual setup enables a comprehensive assessment of whether agents can perform consistently across markets that differ in structure, information flow, and reaction speed. Importantly, these two markets demand distinct strategies and reasoning perspectives for profitability: the stock market rewards long-horizon analysis and disciplined diversification, whereas the prediction market requires short-horizon adaptation and event-driven belief updating.

**U.S. stock market** The U.S. stock market represents a mature, institutionally regulated system where asset prices evolve smoothly, exhibit strong cross-sector correlations, and reflect aggregated fundamentals and macroeconomic signals over time. Effective portfolio management in this environment requires capturing long-term dependencies, modeling hidden correlations, and maintaining diversified risk exposure. We construct a representative portfolio of 15 equities and ETFs spanning major U.S. sectors to ensure diverse responses to external information and macroeconomic shifts. The portfolio includes technology stocks—Apple (AAPL), Microsoft (MSFT), NVIDIA (NVDA), and Meta Platforms (META); financial stocks—JPMorgan Chase (JPM) and Visa (V); energy and industrial stocks—Exxon



**Fig. 2: Observation and action space for LiveTradeBench.** We illustrate examples from both the U.S. stock market and the Polymarket prediction market to demonstrate the observation and action spaces. The observation space consists of three components: the agent’s position, market prices, and relevant news context. The action space represents the portfolio allocation decisions generated by the agent, which can be directly translated into executable trading actions.

Mobil (XOM), Caterpillar (CAT), and Tesla (TSLA); consumer goods stocks—Procter & Gamble (PG), Coca-Cola (KO), Amazon (AMZN), and Walmart (WMT); and healthcare stocks—Johnson & Johnson (JNJ) and UnitedHealth Group (UNH). In addition, a *cash asset* with a constant unit price and zero return rate is included to represent risk-free capital allocation. This composition provides balanced exposure across key sectors in a highly liquid and regulated financial environment, and we collect real-time stock prices as the data source for evaluation.

**Polymarket prediction market** In contrast, the Polymarket prediction market is decentralized, sentiment-driven, and characterized by loosely coupled contracts that respond sharply and asynchronously to real-time information. These markets often move more abruptly and less coherently than stocks, reflecting shifts in collective belief rather than fundamentals. As a result, effective portfolio management here demands rapid adaptation, event-driven reasoning, and sensitivity to evolving narratives. We continuously track ten active binary prediction markets from *Polymarket*, focusing on betting markets related to politics, crypto, technology and finance—such as “Fed rate hike in 2025?”, “Tether insolvent in 2025?”, “U.S. recession in 2025?”, and “USDT depeg in 2025?”. We hypothesize that prediction markets and stock markets respond to the same information with different latency and magnitude: stock markets integrate signals gradually through institutional consensus, while prediction markets react instantly and often overshoot due to speculative sentiment. Together, the two environments—structured financial markets and decentralized prediction markets—offer complementary testbeds for evaluating agents under both stability and uncertainty.

### 3.3. Observation Space

At each timestep, the agent receives an observation  $\mathbf{o}_t = (\mathbf{q}_t, \mathbf{p}_t, \mathbf{c}_t)$  that encapsulates the current market condition, external context, and portfolio status. This observation serves as the primary input for the agent’s decision-making process, integrating quantitative market signals, position information, and qualitative contextual cues. Based on these dynamic inputs, the agent determines its next allocation action, adapting to evolving market trends and external developments. Details on the data collection process of the observation space are available in Appendix §S1.

**Position  $\mathbf{q}_t$**  The position observation  $\mathbf{q}_t$  represents the agent’s current holdings across all assets, including cash. Each component  $q_t^{(i)}$  is a continuous, non-negative value indicating the number of units (or fraction thereof) of asset  $i$  currently held in the portfolio. This formulation differs from discrete buy/hold/sell signals used in traditional trading formulations and instead provides a fine-grained representation of continuous capital allocation. The non-negativity constraint ensures that the agent cannot take short positions, reflecting realistic market restrictions and emphasizing capital distribution among long-only assets.

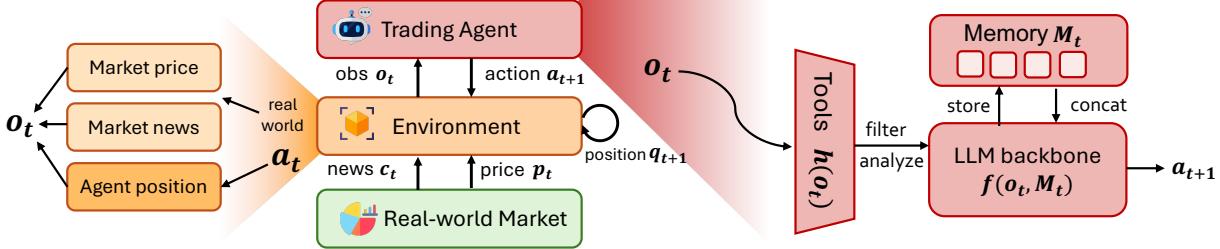
**Market price  $\mathbf{p}_t$**  The price observation includes the latest asset prices and corresponding timestamps for all instruments in the portfolio. For the stock market,  $\mathbf{p}_t$  contains closing prices of the 15 selected equities; for the prediction market, it includes the real-time trading prices of 10 trending *Polymarket* markets. These values serve as the direct basis for portfolio valuation and allocation updates, enabling the agent to track how asset values evolve over time.

**Market context  $\mathbf{c}_t$**  The contextual observation  $\mathbf{c}_t$  mainly provides real-time market news. We collect recent articles from Google News using asset- and topic-specific keywords (e.g., “Federal Reserve,” “inflation,” “NVIDIA stock”). Such information reflects short-term market sentiment, investor attention, and macro-level signals that often precede price movements. To enable the model to reason about these factors, we include the textual summaries of these news items directly in the prompt, allowing LLM-based agents to incorporate qualitative context—such as sentiment shifts, policy expectations, or company-related events—into their trading decisions.

### 3.4. Action Space

At each timestep  $t$ , the agent makes an action  $\mathbf{a}_t \in \mathcal{A}$ , where  $\mathcal{A}$  denotes the probability simplex action space. Each component  $a_t^{(i)}$  specifies the proportion of the total portfolio value  $v_t$  allocated to asset  $i$ , satisfying the budget constraint  $\sum_i a_t^{(i)} = 1$ . By default, we assume a long-only setting where  $a_t^{(i)} \geq 0$  for all  $i$ . This continuous allocation-based formulation abstracts away low-level trading execution and focuses on high-level portfolio rebalancing, allowing agents to express smooth strategic shifts over time and directly optimize for portfolio-level objectives such as return, risk, and stability.

**Stock market action** In the stock market environment, each action component  $a_t^{(i)}$  represents the percentage of the total portfolio value  $v_t$  to be allocated to stock  $i$ . This allocation determines the post-trade position  $\mathbf{q}_{t+1}$  through proportional rebalancing and directly reflects the agent’s capital distribution across sectors.



**Fig. 3: Agent and environment framework in LiveTradeBench.** The left side illustrates the simulated environment, which continuously retrieves real-world market prices and news, updating its internal state accordingly. It also adjusts the agent’s portfolio position based on the executed actions. The right side depicts the portfolio-management agent, equipped with analytical tools to process observations from the environment. The agent maintains a memory of past observations, enabling adaptive and context-aware decision-making.

**Prediction market action** In the prediction market environment, each binary contract corresponds to two complementary assets—YES and NO. For  $k$  active markets, the action vector  $\mathbf{a}_t$  has  $2k$  components, where  $a_{t,\text{YES}}^{(k)}$  and  $a_{t,\text{NO}}^{(k)}$  denote the portfolio allocations to the YES and NO outcomes of market  $k$ , respectively. The agent’s net exposure is defined as  $e_t^{(k)} = a_{t,\text{YES}}^{(k)} - a_{t,\text{NO}}^{(k)}$ , where a positive value indicates higher confidence in the YES outcome.

## 4. Designing LLM-based Agents for Portfolio Management

In our framework, the **agent** is the central decision-making entity that transforms observed information into actionable portfolio allocations. It serves as the bridge between the external market and its internal portfolio memory, continuously adapting its strategy to changing conditions. The agent integrates three intertwined capabilities—*tool use*, *memory*, and *reasoning*—that together enable it to perceive, recall, and act, forming a closed loop of information acquisition, reflection, and execution. Formally, at each timestep  $t$ , the agent receives an observation  $\mathbf{o}_t$ . In addition, the agent maintains an internal memory state  $\mathbf{M}_t$ , which stores the past observation beyond the current one. Conditioned on both the current observation and memory, the agent produces an allocation

$$\mathbf{a}_t = f_\theta(\mathbf{o}_t, \mathbf{M}_t), \quad (3)$$

where  $f_\theta$  is a parameterized policy defining the agent’s trading behavior. Details on how we construct the prompt for the agent are available in Appendix §S2.

**Tool use** The first tool-use component enables the agent to interact with the live environment we provide—fetching, filtering, and extract real-time market and contextual information. While market prices  $\mathbf{p}_t$  and contextual signals  $\mathbf{c}_t$  are emitted by the environment, the tool-use module governs how the agent actively acquires and processes them. It acts as the agent interface with the real world, transforming raw inputs into structured feature representations

$$\tilde{\mathbf{o}}_t = h(\mathbf{o}_t), \quad (4)$$

that capture both quantitative dynamics (e.g., price changes, returns, and volatility features derived from  $\mathbf{p}_t$ ) and qualitative cues (e.g., news relevant to specific markets extracted from  $\mathbf{c}_t$ ). Through tool

use, the agent extends its perception beyond static observations, dynamically gathering and refining evidence to inform its allocation decisions.

**Memory** The second memory component maintains a compact representation of the agent’s recent observations and the outcome of its actions. At each timestep, the agent stores a fixed-length sequence of past observations and concatenates them into a unified memory state:

$$\mathbf{M}_t = \{\mathbf{o}_\tau \mid t - \Delta \leq \tau < t\}, \quad (5)$$

where  $\Delta$  denotes the memory horizon. This concatenated memory provides temporal context beyond the current observation, enabling the agent to capture dependencies such as volatility dynamics, allocation adjustments, and drawdown trends. By conditioning its decisions on both  $\mathbf{o}_t$  and  $\mathbf{M}_t$ , the agent becomes adaptive to evolving market conditions over time.

**Reasoning** The third reasoning component serves as the agent’s decision core. It integrates information gathered through tool use with contextual knowledge retained in memory, forming a coherent understanding of the market at each moment. Before executing any action, the agent engages in a reasoning process that follows the ReAct (Yao et al., 2022) framework, where it first generates intermediate thoughts to interpret signals, recall relevant experiences, and hypothesize about potential outcomes. Similar to chain-of-thought prompting (Wei et al., 2022), this step produces explicit reasoning traces that connect perception and decision. Such interpretability allows the resulting actions to be analyzed and considered as rational responses to evolving market contexts. Through this deliberate reasoning–then–acting cycle, the agent achieves both adaptability and transparency in portfolio management. Such reasoning rationales can be potentially used to help researchers understand the model behavior and utilized as resources to improve LLMs.

## 5. Evaluating LLM-based Agents under Live Test

In this section, we present detailed evaluation results and analyses of live trading conducted from **August 18, 2025**, to **October 24, 2025**—a total of 50 trading days—across trading agents built on 21 unique LLM backbones. Section §5.1 and §5.2 describe the evaluation setup, including model backbones and performance metrics. Section §5.3 reports the main results, and Section §5.4 provides in-depth analyses and discussions.

### 5.1. Backbone LLMs for Evaluation

To benchmark performance in the live trading environment, we evaluate a diverse set of mainstream LLMs as trading agents. Specifically, we consider six representative model families. These models are selected based on two main criteria: (1) their state-of-the-art performance on general-purpose reasoning, knowledge and agentic benchmarks, and (2) their diversity in model size, architecture, and performance levels, which allows us to study performance gradients across heterogeneous systems in financial decision-making tasks.

**LLM family** We include the following representative models: Claude family (Claude-Sonnet-3.7 (Anthropic, 2025a), Claude-Opus-4 & Claude-Sonnet-4 (Anthropic, 2025b), Claude-Opus-4.1 (Anthropic, 2025c)), Grok family (Grok-3 (xAI, 2025a), Grok-4 (xAI, 2025b)), Qwen family (Qwen2.5-72B-Instruct (Yang et al., 2024), Qwen3-235B-A22B-Instruct & Qwen3-235B-A22B-Thinking (Yang

et al., 2025)), LLaMA family (Llama3.3-70B-Instruct-Turbo (Meta, 2025a), Llama4-Scout & Llama4-Maverick (Meta, 2025b)), GPT family (GPT-5 (OpenAI, 2025b), GPT-4o (Hurst et al., 2024), GPT-4.1 (OpenAI, 2025a), GPT-o3 (OpenAI, 2025c)), Kimi family (Kimi-K2-Instruct (Team et al., 2025)), and DeepSeek family (DeepSeek-V3.1 (DeepSeek, 2025), DeepSeek-R1 (Guo et al., 2025)). Each model is wrapped in the same agentic framework that converts market observations into natural-language prompts and parses model outputs into structured portfolio allocation vectors. This unified setup ensures that performance differences primarily reflect the models’ intrinsic reasoning and decision-making abilities rather than disparities in prompt formatting or execution. Details on the model selection are available in Appendix §S3.

## 5.2. Trading Metrics for Evaluation

To evaluate the performance of trading agents, we employ four widely used financial metrics: cumulative return, volatility, Sharpe ratio, and maximum drawdown (MDD). These metrics jointly capture profitability, risk exposure, risk-adjusted efficiency, and downside protection, offering a comprehensive view of trading performance across different markets.

**Cumulative return ( $CR = \frac{v_T - v_0}{v_0}$ )** It measures the overall profitability of an investment strategy over a given evaluation period. Here,  $v_0$  and  $v_T$  denote the initial and final portfolio values, respectively, and  $T$  is the total number of timesteps during evaluation. A higher cumulative return  $CR$  indicates stronger cumulative gains achieved by the trading agent.

**Sharpe ratio ( $SR = \frac{\bar{r} - r_f}{\sigma}$ )** It evaluates the efficiency of returns relative to the amount of risk taken. Here,  $\bar{r}$  denotes the mean return,  $r_f$  is the risk-free rate representing the baseline return from a no-risk investment, and  $\sigma$  is the volatility of returns. In the U.S. stocks,  $r_f$  corresponds to the short-term Treasury yield (typically positive), whereas in the Polymarket,  $r_f$  is set to 0 to reflect the absence of yield on stablecoin-denominated assets. A higher Sharpe ratio  $SR$  signifies that the strategy achieves greater reward per unit of risk, reflecting superior risk-adjusted performance.

**Maximum drawdown ( $MDD = \max_{t \in [1, T]} \frac{\max_{i \in [1, t]} v_i - v_t}{\max_{i \in [1, t]} v_i}$ )** It quantifies the largest observed decline from a historical peak to a subsequent trough in portfolio value before a new peak is reached. Here,  $v_t$  represents the portfolio value at time step  $t$ . A smaller MDD indicates better downside protection and stronger resilience against severe losses.

**Win rate ( $WR = \frac{1}{T-1} \sum_{t=2}^T \mathbb{I}(r_t > 0)$ )** It measures the proportion of profitable trading steps, capturing the agent’s consistency in generating positive returns. Here,  $\mathbb{I}(\cdot)$  equals 1 when the return  $r_t$  is positive and 0 otherwise. A higher win rate  $WR$  indicates that the agent achieves gains more frequently, complementing cumulative and risk-adjusted metrics by reflecting short-term decision reliability.

**Volatility ( $\sigma = \sqrt{\frac{1}{T-1} \sum_{t=1}^T (r_t - \bar{r})^2}$ )** It reflects the variability of returns and serves as a measure of investment risk. Here,  $r_t = \frac{v_t - v_{t-1}}{v_{t-1}}$  represents the return at time step  $t$ ,  $\bar{r} = \frac{1}{T} \sum_{t=1}^T r_t$  is the average return, and  $T$  is the total number of evaluation timesteps. Strategies with lower volatility  $\sigma$  exhibit more stable performance over time.

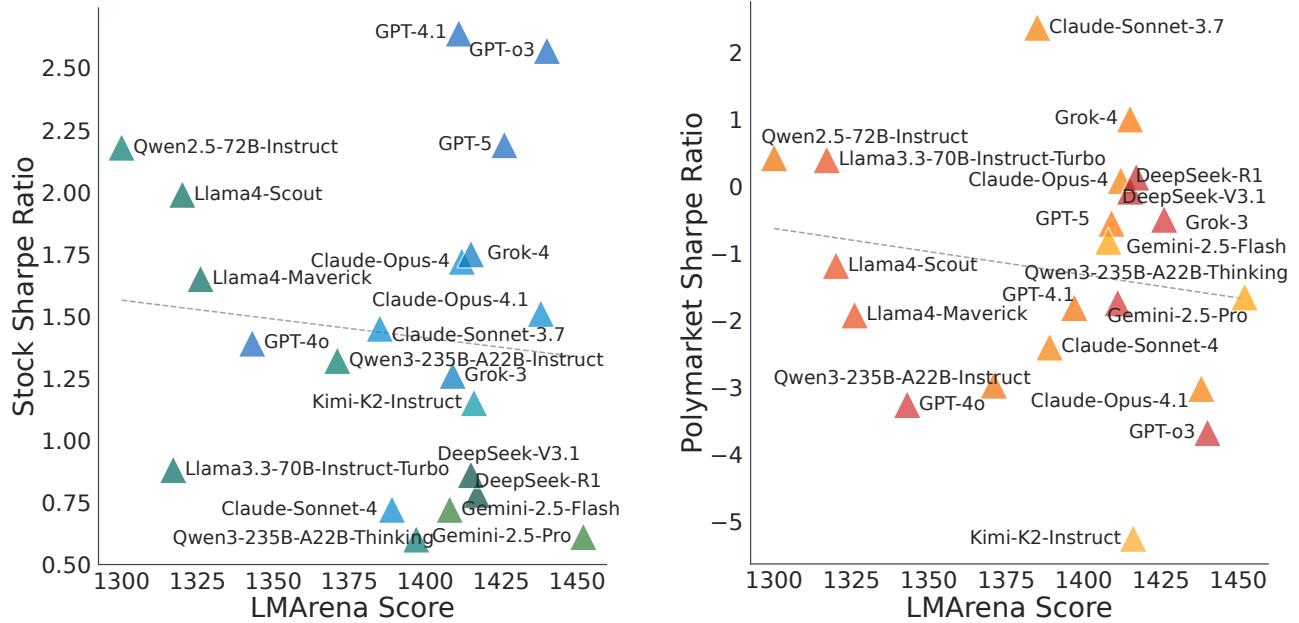
**Table 2: Comparison of trading performance across U.S. stock and Polymarket prediction markets.** We use five key metrics: cumulative return (CR), Sharpe ratio (SR), maximum drawdown (MDD), win rate (WR), and volatility ( $\sigma$ ). The highest value in each column is highlighted in bold.

Model	U.S. Stock Market					Polymarket Prediction Market				
	CR↑	SR↑	MDD↓	WR↑	$\sigma\downarrow$	CR↑	SR↑	MDD↓	WR↑	$\sigma\downarrow$
Claude-Sonnet-3.7	3.63	1.45	2.65	59.18	10.25	<b>20.54</b>	<b>2.38</b>	<b>10.65</b>	51.02	44.64
Claude-Sonnet-4	2.45	0.72	3.23	53.06	12.86	-40.32	-2.40	51.10	38.78	92.16
Claude-Opus-4	3.93	1.72	2.11	63.27	9.45	-2.04	0.09	13.67	46.94	56.38
Claude-Opus-4.1	3.73	1.51	<b>1.84</b>	63.27	10.17	-25.69	-3.02	30.53	48.98	46.81
Grok-3	3.22	1.26	2.13	65.31	10.17	-8.35	-0.55	18.80	53.06	54.35
Grok-4	4.30	1.75	1.92	59.18	10.43	7.38	1.01	13.04	46.94	46.92
Qwen2.5-72B-Instruct	5.15	2.18	2.22	65.31	10.24	1.63	0.43	<b>7.46</b>	<b>59.18</b>	<b>30.36</b>
Qwen3-235B-A22B-Instruct	3.52	1.32	2.04	61.22	10.89	-54.24	-2.97	54.24	40.82	112.37
Qwen3-235B-A22B-Thinking	1.78	0.60	2.32	59.18	<b>9.28</b>	-57.62	-1.81	72.10	38.78	166.92
Llama3.3-70B-Instruct-Turbo	2.72	0.88	3.54	61.22	11.95	1.58	0.40	19.41	57.14	38.15
Llama4-Scout	4.65	1.99	2.62	59.18	9.98	-16.05	-1.18	24.80	51.02	60.37
Llama4-Maverick	4.46	1.65	2.45	53.06	11.59	-18.31	-1.92	28.63	34.69	48.21
GPT-4o	3.55	1.39	2.43	55.10	10.38	-30.96	-3.26	35.31	30.61	53.75
GPT-4.1	<b>6.25</b>	<b>2.64</b>	1.92	<b>65.31</b>	10.51	-33.69	-1.74	37.98	40.82	95.27
GPT-5	5.31	2.19	2.53	65.31	10.60	-23.96	-0.49	38.92	44.90	130.37
GPT-o3	6.04	2.57	2.27	61.22	10.41	-54.84	-3.68	60.99	40.82	97.27
Gemini-2.5-Flash	2.10	0.72	3.10	55.10	11.25	-22.40	-0.82	42.35	38.78	115.42
Gemini-2.5-Pro	1.95	0.61	2.85	50.00	10.98	-35.15	-1.65	49.80	34.69	101.87
Kimi-K2-Instruct	3.07	1.15	3.32	53.06	10.53	-53.44	-5.26	54.74	28.57	69.41
DeepSeek-V3.1	2.46	0.86	2.45	59.18	10.61	-4.68	-0.07	22.43	48.98	64.74
DeepSeek-R1	2.10	0.78	2.20	61.22	9.11	-13.19	0.14	44.16	42.86	143.25

### 5.3. Evaluation Results

**Trading performance on one market does not generalize to another.** As shown in Table 2, the Sharpe ratio correlation between the two markets is close to zero, indicating that success in one market does not imply success in the other. This highlights the need for market-specific trading strategies. For example, Qwen2.5-72B-Instruct and Grok-4 show relatively consistent performance across both the stock and prediction markets, suggesting more stable and low-volatility strategies. In contrast, GPT-4.1 achieves the highest cumulative return rate ( $> 6\%$ ) in the stock market but performs poorly in Polymarket (return  $< -30\%$ ), likely due to overreactive allocation changes under higher volatility. Overall, the prediction market exhibits faster dynamics, greater volatility, and deeper drawdowns (MDD), demanding more agile and risk-tolerant strategies.

**High general LLM capability does not imply strong financial performance.** Figures 4 show that general LLM ability, as measured by LMarena scores, has slightly negative correlation with trading abilities. In the stock market, the Spearman correlation between LMarena scores and cumulative returns is only 0.054—virtually no relationship. In Polymarket, the correlation drops to -0.38, meaning models with higher general language ability often perform worse in the dynamic market. Thus, state-

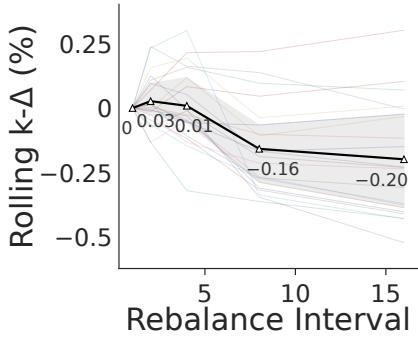


**Fig. 4: Correlation between LM Arena score and Sharpe ratio across two markets.** (left) U.S. stock market. (right) Polymarket prediction market. Models from different families are shown in different colors, and the dashed line indicates the linear regression fit.

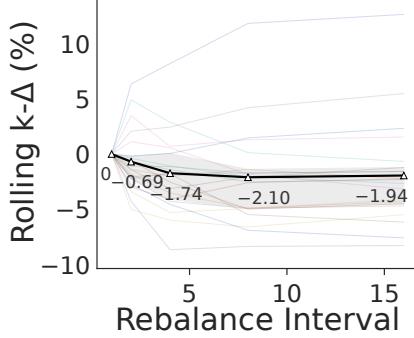
of-the-art LLMs on general benchmarks do not necessarily translate to state-of-the-art performance in dynamic, real-world trading. It highlights the uniqueness and necessity of our environment.

**Distinct portfolio management styles emerge across models.** Different models exhibit distinct management preferences. Claude-Opus-4.1 and Grok-4 adopt conservative strategies characterized by lower volatility and smaller drawdowns, prioritizing stability over aggressive gains. In contrast, Kimi-K2-Instruct and GPT-5 display more risk-seeking behaviors—accepting higher volatility and MDD in pursuit of greater returns. Beyond return and risk metrics, models also differ notably in their portfolio composition and cash management patterns. For instance, GPT-4o consistently focuses on a few core assets (AAPL, MSFT, NVDA), whereas GPT-5 diversifies across a broader range of stocks with smaller position ratios. Likewise, Llama4-Scout maintains a persistently high cash ratio (above 20%), reflecting a more cautious liquidity stance, while GPT-5 always keeps cash below 10% throughout trading except in extremely high risk. These behavioral patterns are not limited to a single market—similar management styles emerge consistently across both markets.

**Large reasoning models do not confer trading advantages.** Consistent with findings from Chen et al. (2025), models explicitly designed for reasoning—such as DeepSeek-R1, Qwen3-235B-A22B-Thinking, and GPT-o3—do not outperform others in trading performance. Instead, they exhibit substantially higher volatility ( $>140$  in Polymarket), implying over-adjustment during the decision process. The type of reasoning beneficial for mathematical or coding tasks does not straightforwardly transfer to financial or social reasoning. In fact, excessive deliberation observed in these models during trading can introduce instability and degrade trading consistency.



**Fig. 5: Rolling  $k$ -delta analysis on U.S. stocks.** We evaluate rebalance intervals  $k \in \{1, 2, 4, 8, 16\}$ . The black line denotes the mean performance across 21 models, and the shaded gray region indicates the 25–75% confidence interval.



**Fig. 6: Rolling  $k$ -delta analysis on Polymarket.** We also evaluate rebalance intervals  $k \in \{1, 2, 4, 8, 16\}$ . The black line denotes the mean performance across 21 models, and the shaded gray region indicates the 25–75% confidence interval.



**Fig. 7: Decision-making rationale analysis.** Each bar indicates the proportion of reasoning traces that reference position, price, or news information. A single reasoning trace may include multiple information sources.

#### 5.4. Analysis and Discussion

In this section, we quantitatively analyze two core questions that probe the fundamental capabilities of LLM-based trading agents. (1) *Are LLM-based agents merely random guessers?* — This examines whether the agents’ trading behaviors reflect meaningful market understanding or simply random fluctuations. (2) *How do agents reason and make trading decisions?* — This investigates the internal rationale behind their actions, revealing whether their decisions are grounded in coherent reasoning patterns. Together, these analyses shed light on both the effectiveness and interpretability of LLM-based agents in dynamic, uncertain market environments.

**Are LLM agents just random guessing?** To verify that LLM-based trading agents exhibit genuine market awareness rather than random behavior, we design the **rolling- $k$  delta** ( $\Delta_k$ ) analysis. The key idea is that if the agents’ decisions are random, delaying their actions by several days should not systematically affect performance. Conversely, if they truly adapt to changing market conditions, stale decisions should lead to measurable degradation. For each trading day  $t$ , we fix the portfolio position to the one taken  $k$  days earlier,  $\mathbf{q}_t^{(k)} = \mathbf{q}_{t-k}$ , and compute daily and cumulative returns as

$$r_t^{(k)} = \frac{(\mathbf{q}_{t-k})^\top (\mathbf{p}_{t+1} - \mathbf{p}_t)}{(\mathbf{q}_{t-k})^\top \mathbf{p}_t}, \quad CR^{(k)} = \prod_{t=k}^{T-1} (1 + r_t^{(k)}) - 1. \quad (6)$$

The rolling- $k$  delta is then defined as  $\Delta_k = CR^{(k)} - CR^{(0)}$ , capturing the cumulative return loss when the agent’s actions lag behind the market by  $k$  days. A negative  $\Delta_k$  indicates that more frequent rebalancing improves performance. As shown in Figure 5 and Figure 6, larger  $k$  (slower updates) leads to higher degradation, confirming that timely decision updates are beneficial. Interestingly, in the stock market, returns slightly improve to 0.03% when  $k = 2$ , suggesting smoother dynamics and lower time sensitivity compared to the Polymarket, where performance degrades 2% as  $k$  increases. Overall, these results demonstrate that LLM-based agents do not act randomly—their trading strategies depend on contemporaneous market signals, and delaying their actions systematically harms performance.

**How do LLM agents reason and make decisions?** To investigate decision-making rationale, we employ LLM-based reasoning annotation. For each day’s reasoning trace, another LLM automatically identifies whether the agent’s explanation references: (1) portfolio position, (2) market price history, or (3) market news. Figure 7 summarizes the distribution of these factors. In both markets, *news* emerges as the most frequently cited factor, followed by *market price history*, while *position* information is less dominant. Moreover, the Polymarket agents rely more heavily on news signals, while stock-market agents emphasize price trends—validating our hypothesis that the two markets exhibit distinct dynamics. Since the total percentage of reasoning references exceeds 100%, many decisions integrate multiple information sources, indicating complex reasoning processes. Specifically, agents often mention “price momentum” when analyzing price history and focus on potential outcomes or implications when discussing market news.

## 6. Case Study

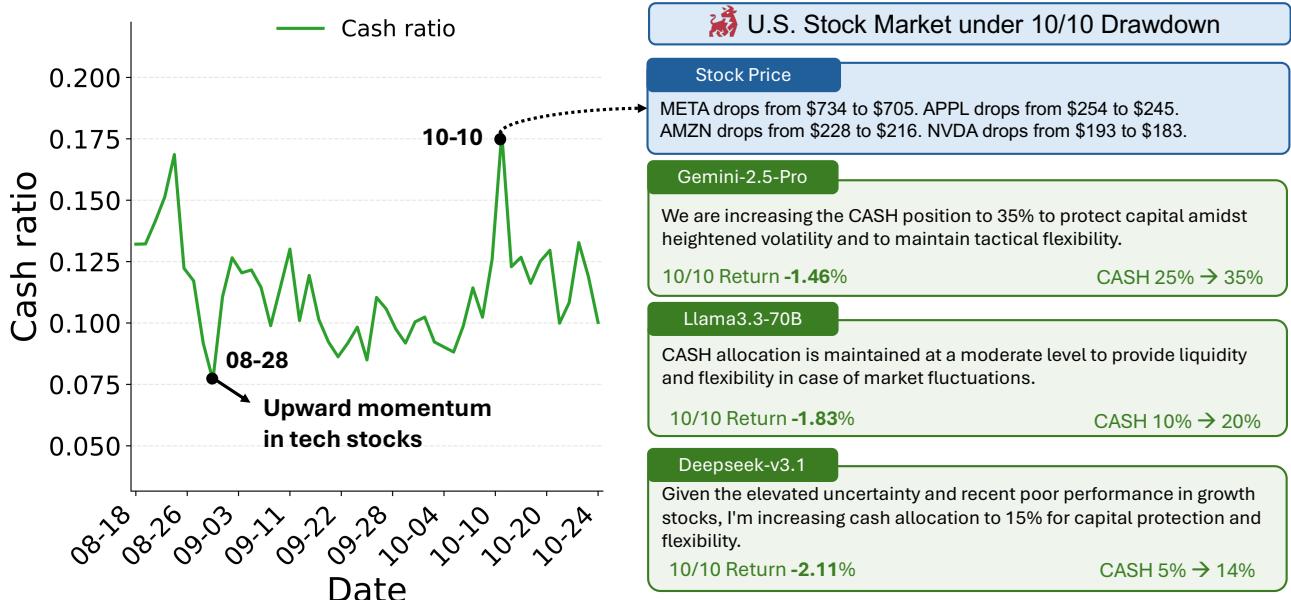
In this section, we show examples of both the U.S stock and Polymarket prediction markets by selecting representative assets among each of them and highlighting two distinct and extreme time points for each to analyze. This allows us to examine the reasoning behind the agents’ decisions and understand how their choices correlate with market conditions.

### 6.1. Cash Asset Dynamics in the U.S. Stock Market

In Figure 8, we analyze the dynamics of cash assets in the portfolio and highlight two contrasting market scenarios—a positive (bullish) case and a negative (bearish) case—in the U.S. stock market. The cash ratio serves as an informative indicator of risk management: a higher cash ratio typically reflects greater risk aversion and a defensive stance, whereas a lower cash ratio suggests stronger market confidence and a more aggressive investment posture.

**Tech stock rally on August 28** On August 28, the average cash ratio across all 21 models declined steadily over three consecutive days, dropping from 17% to 7.5% within four trading days. This trend coincided with a strong rally in major technology stocks such as Meta (META), Apple (AAPL), and Microsoft (MSFT), which encouraged agents to invest more aggressively and reduce their cash holdings. Notably, GPT-4.1, the agent achieving the highest cumulative return rate, provided the following rationale for its allocation decision: “*This allocation increases exposure to leading AI and tech growth stocks (NVDA, MSFT, META) following strong earnings momentum and positive analyst sentiment, while maintaining solid positions in diversified blue chips for stability and sector balance.*” This reasoning explicitly aligns with the observed decrease in the cash ratio during the bullish market phase.

**Market drawdown on October 10** In contrast, the sharp market drawdown on October 10 induced the opposite behavior. Most of the stock prices dropped significantly—Tesla (TSLA) fell over 5%, while Amazon (AMZN) and Nvidia (NVDA) declined by more than 4%—leading to negative returns for all agents on that day. In response, most agents increased their cash holdings to mitigate risk, reflecting a collective shift toward a defensive strategy. As shown on the right of Figure 8, multiple agents provided similar reasoning related to “increasing CASH positions to protect against volatility.” Among them, Gemini-2.5-Pro, which maintained a relatively high cash position and converted additional assets to cash before the downturn, experienced the smallest loss that day.



**Fig. 8: Case study for U.S. stock markets.** (Left) The average cash ratio across 21 models over 50 trading days. (Right) A zoomed-in view of the sharp drawdown on October 10, during which portfolios exhibited a sudden increase in cash holdings. We visualize the market condition (price change) and present the reasoning traces of the best-performing model on October 10 (Gemini-2.5-Pro) and one of the worst-performing models on October 10 (DeepSeek-V3.1) for comparison.

## 6.2. Russia–Ukraine Ceasefire Market in Polymarket Prediction Market

We analyze the market “*Russia × Ukraine ceasefire in 2025?*” on Polymarket, focusing on how real-time news influences the decision-making of LLM-based agents. Polymarket’s high sensitivity to external information makes it a natural testbed to evaluate how models interpret and act on dynamic geopolitical signals. We select this market for the case study because it experiences frequent fluctuations during the 50 trading days, resulting in distinct behaviors and returns across models. As shown in Figure 9, the Grok-3 model is able to conduct belief-based reasoning, adjusting its internal estimate of the ceasefire probability from 0.15 on October 13 to 0.22 on October 17.

**Reactive change on October 13 without profit** On October 13, most agents abruptly switched their portfolios from No to Yes positions after two optimistic news events: (1) Zelenskyy stated that “*the Gaza deal brings hope for Ukraine*,” and (2) reports surfaced that *Trump shared U.S. intelligence to help Kyiv strike Russian energy targets*. These headlines appeared relevant to the ceasefire, prompting agents to buy into the Yes position. However, the Polymarket price showed little actual movement, and this reactive change produced no profit. The news turned out to have limited causal impact on the ceasefire likelihood. This case highlights a key challenge for LLM-based agents: distinguishing between attention-grabbing but non-decisive news and genuinely influential events. Acting on superficial correlations can lead to overreaction and unprofitable trades.

**Strategic hold on October 17 with profit** In contrast, on October 17, when news broke that *Zelenskyy visited the White House*, most agents strengthened their Yes positions and held them through the



**Fig. 9: Case study in Polymarket prediction markets.** (Left) The average holding ratios of “Yes” and “No” position ratios in the market “Russia x Ukraine ceasefire in 2025?” across 21 models. (Right) A zoomed-in view of two abrupt shifts (October 13 and October 17), along with the corresponding news events and the reasoning traces of Grok-3 explaining these allocation decisions.

following day. This time, the Polymarket price steadily increased from October 17 to 18, leading to tangible profits. Unlike the earlier overreaction, agents displayed more grounded reasoning—citing “recent diplomatic developments” and recognizing “a significant price jump to 0.18” as confirming evidence. This scenario illustrates that maintaining positions through credible, high-impact events can yield better outcomes than frequent reactive shifts based on weak signals.

## 7. Conclusion

In this work, we present LiveTradeBench, a live multi-market environment for evaluating LLM-based agents in realistic portfolio management tasks. LiveTradeBench introduces a new paradigm for assessing model intelligence beyond static benchmarks, enabling continuous interaction, reasoning, and adaptation within real-time stock and prediction markets. Through 50-day live experiments, we find that strong performance in one market does not generalize to others, underscoring the heterogeneity and specialization required across market types. Moreover, high scores on general-purpose benchmarks like LM Arena do not necessarily translate into superior trading performance, highlighting a gap between text intelligence and dynamic decision-making. Finally, our analyses reveal that LLM-based agents rely jointly on historical price trends, market news, allocation history, exhibiting distinct behavioral patterns under extreme conditions. Overall, LiveTradeBench provides a foundation for studying how LLM-based agents perceive, reason, and act under uncertainty in live and realistic trading environments—paving the way for developing more adaptive, financially grounded, and socially intelligent agent systems.

## Limitation and Future Work

Despite demonstrating the feasibility of evaluating LLM-based trading agents in live multi-market environments, our framework still has three main limitations that point to promising directions for future research and development.

**Transaction costs and market frictions** Our current environment and evaluation do not account for transaction fees, bid–ask spreads, liquidity constraints, or other real-world trading frictions. Ignoring these factors may overestimate achievable returns, especially for strategies that rely on frequent rebalancing. Future work will incorporate more realistic cost models and slippage simulations to better approximate real trading conditions.

**Limited observation and action space** The current framework constrains both the observation and action spaces due to the limited context length of existing LLMs. For the *observation space*, the agent can only access a restricted temporal window of price, position, and news histories, and these are limited to a small set of markets in both the stock market and the prediction market. Moreover, news inputs are truncated to titles and abstracts rather than full articles, preventing the agent from incorporating long-form textual information that may contain deeper market signals. For the *action space*, the scope of possible trading actions is similarly constrained by the limited number of supported markets, reducing the complexity and richness of allocation decisions. Future work could extend the framework to support longer temporal horizons, richer textual context, and dynamic market expansion—enabling agents to observe and act within more realistic, information-rich environments.

**Simplified agent design** The current agent architecture integrates basic tool use and memory under the ReAct framework but remains limited in reasoning depth and temporal abstraction. Future work can enhance each component systematically. For *tool use*, agents can be equipped with more specialized analytical and retrieval tools for financial reasoning, news interpretation, and risk assessment. For *memory*, richer hierarchical and long-term memory mechanisms can be introduced to capture temporal dependencies and retain cross-market knowledge over extended horizons. Beyond the basic ReAct-based setups, the current framework can be extended to a *multi-agent* paradigm, such as TradingAgents (Xiao et al., 2024b), to better model heterogeneous roles and market interactions. Finally, incorporating reinforcement learning (RL) to train trading agents (Xiao et al., 2025) represents a promising direction for improving decision quality—enabling agents to learn from experience, refine their reasoning, and continuously adapt to dynamic market conditions.

## Open-source Application

To democratize research on LLM-based trading agents, we release an open-source Python package, `live-trade-bench`<sup>†</sup>, which provides simple APIs for data collection, environment setup, and agent construction. Building on this package, we also develop a web application that deploys our trading environment in real time, enabling live data streaming and interactive monitoring of agent performance. Details of the user interface (UI) design are provided in Appendix §S4.

---

<sup>†</sup><https://pypi.org/project/live-trade-bench/>

## References

- R. Aksitov, S. Miryoosefi, Z. xiao Li, D. Li, S. Babayan, K. Kopparapu, Z. Fisher, R. Guo, S. Prakash, P. Srinivasan, M. Zaheer, F. X. Yu, and S. Kumar. Rest meets react: Self-improvement for multi-step reasoning llm agent. *ArXiv*, abs/2312.10003, 2023. URL <https://api.semanticscholar.org/CorpusId:266335848>.
- Anthropic. Claude 3.7 sonnet and claude code. <https://www.anthropic.com/news/clause-3-7-sonnet>, 2025a.
- Anthropic. Introducing claude 4. <https://www.anthropic.com/news/clause-4>, 2025b.
- Anthropic. Claude opus 4.1. <https://www.anthropic.com/news/clause-opus-4-1>, 2025c.
- A. Bigeard, L. Nashold, R. Krishnan, and S. Wu. Finance agent benchmark: Benchmarking llms on real-world financial research tasks. *arXiv preprint arXiv:2508.00828*, 2025.
- A. Briola, J. Turiel, R. Marcaccioli, and T. Aste. Deep reinforcement learning for active high frequency trading. *ArXiv*, abs/2101.07107, 2021. URL <https://api.semanticscholar.org/CorpusId:231632086>.
- J. Chen, S. Yuan, R. Ye, B. P. Majumder, and K. Richardson. Put your money where your mouth is: Evaluating strategic planning and execution of llm agents in an auction arena. *ArXiv*, abs/2310.05746, 2023. URL <https://api.semanticscholar.org/CorpusId:263831697>.
- Y. Chen, Z. Yao, Y. Liu, J. Ye, J. Yu, L. Hou, and J. Li. Stockbench: Can llm agents trade stocks profitably in real-world markets? *arXiv preprint arXiv:2510.02209*, 2025.
- Z. Chen, W. Chen, C. Smiley, S. Shah, I. Borova, D. Langdon, R. Moussa, M. Beane, T.-H. Huang, B. Routledge, et al. Finqa: A dataset of numerical reasoning over financial data. *arXiv preprint arXiv:2109.00122*, 2021.
- Z. Chen, S. Li, C. Smiley, Z. Ma, S. Shah, and W. Y. Wang. Convfinqa: Exploring the chain of numerical reasoning in conversational finance question answering. *arXiv preprint arXiv:2210.03849*, 2022.
- W.-L. Chiang, L. Zheng, Y. Sheng, A. N. Angelopoulos, T. Li, D. Li, B. Zhu, H. Zhang, M. Jordan, J. E. Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. In *Forty-first International Conference on Machine Learning*, 2024.
- K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- O. Contributors. Opencompass: A universal evaluation platform for foundation models. <https://github.com/open-compass/opencompass>, 2023.
- DeepSeek. Deepseek-v3.1 release. <https://api-docs.deepseek.com/news/news250821>, 2025.
- M. DeHaven, H. Firestone, and C. Webster. Minute-by-minute: Financial markets' reaction to the 2020 u.s. election. In *unknown*, 2024. URL <https://api.semanticscholar.org/CorpusId:271038710>.

- D. Emmanoulopoulos, O. Olby, J. Lyon, and N. R. Stillman. To trade or not to trade: An agentic approach to estimating market risk improves trading decisions. *ArXiv*, abs/2507.08584, 2025. URL <https://api.semanticscholar.org/CorpusId:280281425>.
- H. Gao, J. Chen, and X. Li. Finbench: Benchmarking llms for financial decision-making. *NeurIPS*, 2024a.
- S. Gao, Y. Wen, M. Zhu, J. Wei, Y. Cheng, Q. Zhang, and S. Shang. Simulating financial market via large language model based agents. *ArXiv*, abs/2406.19966, 2024b. URL <https://api.semanticscholar.org/CorpusId:270845937>.
- E. C. Garrido-Merch'an, M. C. Vaca, Á. López-López, and C. M. de Ibarreta. Deep reinforcement learning agents for strategic production policies in microeconomic market simulations. *ArXiv*, abs/2410.20550, 2024. URL <https://api.semanticscholar.org/CorpusId:273653990>.
- J. Gu, J. Ye, G. Wang, and W. Yin. Adaptive and explainable margin trading via large language models on portfolio management. *Proceedings of the 5th ACM International Conference on AI in Finance*, 2024. URL <https://api.semanticscholar.org/CorpusId:274086023>.
- D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- W. Han, J. Huang, Q. Xie, B. Zhang, Y. Lai, and M. Peng. Mastering pair trading with risk-aware recurrent reinforcement learning. *ArXiv*, abs/2304.00364, 2023a. URL <https://api.semanticscholar.org/CorpusId:257913550>.
- W. Han, B. Zhang, Q. Xie, M. Peng, Y. Lai, and J. Huang. Select and trade: Towards unified pair trading with hierarchical reinforcement learning. *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2023b. URL <https://api.semanticscholar.org/CorpusId:256231505>.
- H. He, W. Yao, K. Ma, W. Yu, Y. Dai, H. Zhang, Z. Lan, and D. Yu. Webvoyager: Building an end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*, 2024.
- R. He and W. Xu. Finmem: Evaluating memory and bias in financial lilm benchmarks. *arXiv preprint arXiv:2409.08712*, 2024.
- L. Heinrich, A. Shivarova, and M. Zurek. Factor investing: alpha concentration versus diversification. *Journal of Asset Management*, 22:464 – 487, 2021. URL <https://api.semanticscholar.org/CorpusId:236248859>.
- D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- L. Hu, J. Jiao, J. Liu, Y. Ren, Z. Wen, K. Zhang, X. Zhang, X. Gao, T. He, F. Hu, et al. Finsearch-comp: Towards a realistic, expert-level evaluation of financial search and reasoning. *arXiv preprint arXiv:2509.13160*, 2025.
- A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

- M. R. Islam. The evolution of alpha in finance harnessing human insight and llm agents. *ArXiv*, abs/2505.14727, 2025. URL <https://api.semanticscholar.org/CorpusId:278782847>.
- Y. Jiang, Y. Wang, X. Zeng, W. Zhong, L. Li, F. Mi, L. Shang, X. Jiang, Q. Liu, and W. Wang. Followbench: A multi-level fine-grained constraints following benchmark for large language models. *arXiv preprint arXiv:2310.20410*, 2023.
- Z. Jiang and S. Zhou. Finagent: Evaluating financial reasoning in large language models. *arXiv preprint arXiv:2506.03142*, 2025.
- C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- J. Jumadinova and P. Dasgupta. A multi-agent prediction market based on partially observable stochastic game. *ArXiv*, abs/1203.6035, 2011. URL <https://api.semanticscholar.org/CorpusId:14232149>.
- J. Kasai, K. Sakaguchi, Y. Takahashi, R. L. Bras, A. Asai, X. V. Yu, D. R. Radev, N. A. Smith, Y. Choi, and K. Inui. Realtime qa: What's the answer right now? *ArXiv*, abs/2207.13332, 2022. URL <https://api.semanticscholar.org/CorpusId:251105205>.
- K. J. Koa, Y. Ma, R. Ng, and T.-S. Chua. Learning to generate explainable stock predictions using self-reflective large language models. *Proceedings of the ACM Web Conference 2024*, 2024. URL <https://api.semanticscholar.org/CorpusId:267500314>.
- J. Y. Koh, R. Lo, L. Jang, V. Duvvur, M. C. Lim, P.-Y. Huang, G. Neubig, S. Zhou, R. Salakhutdinov, and D. Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*, 2024a.
- J. Y. Koh, S. McAleer, D. Fried, and R. Salakhutdinov. Tree search for language model agents. *ArXiv*, abs/2407.01476, 2024b. URL <https://api.semanticscholar.org/CorpusId:270870063>.
- R. Koning and R. Zijm. Betting market efficiency and prediction in binary choice models. *Annals of Operations Research*, 325:135 – 148, 2022. URL <https://api.semanticscholar.org/CorpusId:248463684>.
- Z. Kou, H. Yu, J. Peng, and L. Chen. Automate strategy finding with llm in quant investment. *ArXiv*, abs/2409.06289, 2024. URL <https://api.semanticscholar.org/CorpusId:272550540>.
- C. Li, Y. Shi, Y. Luo, and N. Tang. Will llms be professional at fund investment? deepfund: A live arena perspective. In *unknown*, 2025a. URL <https://api.semanticscholar.org/CorpusId:277272186>.
- C. Li, Y. Shi, C. Wang, Q. Duan, R. Ruan, W. Huang, H. Long, L. Huang, Y. Luo, and N. Tang. Time travel is cheating: Going live with deepfund for real-time fund investment benchmarking. *ArXiv*, abs/2505.11065, 2025b. URL <https://api.semanticscholar.org/CorpusId:278715123>.
- H. Li, Y. Cao, Y. Yu, S. R. Javaji, Z. Deng, Y. He, Y. Jiang, Z. Zhu, K. Subbalakshmi, G. Xiong, J. Huang, L. Qian, X. Peng, Q. Xie, and J. W. Suchow. Investorbench: A benchmark for financial decision-making tasks with llm-based agent. *ArXiv*, abs/2412.18174, 2024. URL <https://api.semanticscholar.org/CorpusId:274992211>.

- Q. Li and X. Zhao. Contesttrade: A multi-agent trading system based on internal contest mechanism. *arXiv preprint arXiv:2508.00554*, 2025.
- W. W. Li, H. Kim, M. Cucuringu, and T. Ma. Can llm-based financial investing strategies outperform the market in long run? *ArXiv*, abs/2505.07078, 2025c. URL <https://api.semanticscholar.org/CorpusId:278501425>.
- X. Li, Y. Zeng, X. Xing, J. Xu, and X. Xu. Profit mirage: Revisiting information leakage in llm-based financial agents. In *unknown*, 2025d. URL <https://arxiv.org/pdf/2510.07920.pdf>.
- Y. Li, Y. Yu, H. Li, Z. Chen, and K. Khashanah. Tradinggpt: Multi-agent system with layered memory and distinct characters for enhanced financial trading performance. 2023. URL <https://api.semanticscholar.org/CorpusId:261582775>.
- H. Liang and Z. Zhang. Livecodebench: Evaluating llms in real-world coding and execution environments. *arXiv preprint arXiv:2410.06521*, 2024.
- A. Lopez-Lira. Can large language models trade? testing financial theories with llm agents in market simulations. In *unknown*, 2025. URL <https://api.semanticscholar.org/CorpusId:277787336>.
- G. Lucarelli and M. Borrotti. A deep q-learning portfolio management framework for the cryptocurrency market. *Neural Computing and Applications*, 32:17229 – 17244, 2020. URL <https://doi.org/10.1007/s00521-020-05359-8>.
- T. Ma, J. Du, W. Huang, W. Wang, L. Xie, X. Zhong, and J. T. Zhou. Agent trading arena: A study on numerical understanding in llm-based agents. In *unknown*, 2025. URL <https://api.semanticscholar.org/CorpusId:276580341>.
- Meta. Llama 3.3. [https://www.llama.com/docs/model-cards-and-prompt-formats/llama3\\_3/](https://www.llama.com/docs/model-cards-and-prompt-formats/llama3_3/), 2025a.
- Meta. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, 2025b.
- F. Nie, K. Z. Liu, Z. Wang, R. Sun, W. Liu, W. Shi, H. Yao, L. Zhang, A. Y. Ng, J. Zou, S. Koyejo, Y. Choi, P. Liang, and N. Muennighoff. Uq: Assessing language models on unsolved questions, 2025. URL <https://arxiv.org/abs/2508.17580>.
- OpenAI. Introducing gpt-4.1 in the api. <https://openai.com/index/gpt-4-1/>, 2025a.
- OpenAI. Gpt-5 is here. <https://openai.com/zh-Hans-CN/gpt-5/>, 2025b.
- OpenAI. Introducing openai o3 and o4-mini. <https://openai.com/zh-Hans-CN/index/introducing-o3-and-o4-mini/>, 2025c.
- C. Papadakis, G. Filandrianos, A. Dimitriou, M. Lympereiou, K. Thomas, and G. Stamou. Stocksim: A dual-mode order-level simulator for evaluating multi-agent llms in financial markets. *ArXiv*, abs/2507.09255, 2025. URL <https://api.semanticscholar.org/CorpusId:280265856>.
- L. Phan, A. Gatti, Z. Han, N. Li, J. Hu, H. Zhang, C. B. C. Zhang, M. Shaaban, J. Ling, S. Shi, et al. Humanity's last exam. *arXiv preprint arXiv:2501.14249*, 2025.

- P. Putta, E. Mills, N. Garg, S. Motwani, C. Finn, D. Garg, and R. Rafailov. Agent q: Advanced reasoning and learning for autonomous ai agents. *ArXiv*, abs/2408.07199, 2024. URL <https://api.semanticscholar.org/CorpusId:271865516>.
- V. Pyatkin, S. Malik, V. Graf, H. Ivison, S. Huang, P. Dasigi, N. Lambert, and H. Hajishirzi. Generalizing verifiable instruction following, 2025.
- S. Quan, J. Yang, B. Yu, B. Zheng, D. Liu, A. Yang, X. Ren, B. Gao, Y. Miao, Y. Feng, et al. Codeelo: Benchmarking competition-level code generation of llms with human-comparable elo ratings. *arXiv preprint arXiv:2501.01257*, 2025.
- D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.
- R. S. Shah, K. Chawla, D. Eidnani, A. Shah, W. Du, S. Chava, N. Raman, C. Smiley, J. Chen, and D. Yang. When flue meets flang: Benchmarks and large pre-trained language model for financial domain. *arXiv preprint arXiv:2211.00083*, 2022.
- R. Sun, A. Stefanidis, Z. Jiang, J. S. X. J.-L. University, S. O. Mathematics, Physics, D. of Financial, A. M. X. J.-L. U. E. College, A. Schoolof, and A. hoc Computing. Combining transformer based deep reinforcement learning with black-litterman model for portfolio optimization. In *unknown*, 2024. URL <https://api.semanticscholar.org/CorpusId:268032065>.
- S. Sun, R. Wang, and B. An. Reinforcement learning for quantitative trading. *ACM Transactions on Intelligent Systems and Technology*, 14:1 – 29, 2021. URL <https://api.semanticscholar.org/CorpusId:238198196>.
- Z. Tang, Z. Chen, J. Yang, J. Mai, Y. Zheng, K. Wang, J. Chen, and L. Lin. Alphaagent: Llm-driven alpha mining with regularized exploration to counteract alpha decay. *ArXiv*, abs/2502.16789, 2025. URL <https://api.semanticscholar.org/CorpusId:276574595>.
- K. Team, Y. Bai, Y. Bao, G. Chen, J. Chen, N. Chen, R. Chen, Y. Chen, Y. Chen, Y. Chen, et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.
- F. Tian, F. D. Salim, and H. Xue. Tradinggroup: A multi-agent trading system with self-reflection and data-synthesis. *ArXiv*, abs/2508.17565, 2025. URL <https://api.semanticscholar.org/CorpusId:280711571>.
- Z. Wang, X. Liu, and H. Zhao. Trading-r1: Reinforcement-guided large language models for financial markets. *arXiv preprint arXiv:2509.11420*, 2025.
- J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- T. Wu and R. Zhang. Finllm: A systematic evaluation framework for financial large language models. *AAAI*, 2025.
- xAI. Grok 3 beta — the age of reasoning agents. <https://x.ai/news/grok-3>, 2025a.
- xAI. Grok 4. <https://x.ai/news/grok-4>, 2025b.

- Y. Xiao, E. Sun, D. Luo, and W. Wang. Tradingagents: Multi-agents llm financial trading framework. *ArXiv*, abs/2412.20138, 2024a. URL <https://api.semanticscholar.org/CorpusId:275133732>.
- Y. Xiao, E. Sun, D. Luo, and W. Wang. Tradingagents: Multi-agents llm financial trading framework. *arXiv preprint arXiv:2412.20138*, 2024b.
- Y. Xiao, E. Sun, T. Chen, F. Wu, D. Luo, and W. Wang. Trading-r1: Financial trading with llm reasoning via reinforcement learning. *arXiv preprint arXiv:2509.11420*, 2025.
- T. Xie, D. Zhang, J. Chen, X. Li, S. Zhao, R. Cao, T. J. Hua, Z. Cheng, D. Shin, F. Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *Advances in Neural Information Processing Systems*, 37:52040–52094, 2024.
- G. Xiong, Z. Deng, K. Wang, Y. Cao, H. Li, Y. Yu, X. Peng, M. Lin, K. E. Smith, X.-Y. Liu, J. Huang, S. Ananiadou, and Q. Xie. Flag-trader: Fusion llm-agent with gradient-based reinforcement learning for financial trading. *ArXiv*, abs/2502.11433, 2025. URL <https://api.semanticscholar.org/CorpusId:276408244>.
- T. Xu, H. Li, and J. Liu. Tradingagent: Multi-agent financial decision-making with llms. *arXiv preprint arXiv:2412.20138*, 2024.
- A. Yang, A. Li, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Gao, C. Huang, C. Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Q. A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, G. Dong, H. Wei, H. Lin, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Lin, K. Dang, K. Lu, K. Bao, K. Yang, L. Yu, M. Li, M. Xue, P. Zhang, Q. Zhu, R. Men, R. Lin, T. Li, T. Xia, X. Ren, X. Ren, Y. Fan, Y. Su, Y.-C. Zhang, Y. Wan, Y. Liu, Z. Cui, Z. Zhang, Z. Qiu, S. Quan, and Z. Wang. Qwen2.5 technical report. *ArXiv*, abs/2412.15115, 2024. URL <https://api.semanticscholar.org/CorpusID:274859421>.
- S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. R. Narasimhan, and Y. Cao. React: Synergizing reasoning and acting in language models. In *The eleventh international conference on learning representations*, 2022.
- Y. Ye, H. Pei, B. Wang, P.-Y. Chen, Y. Zhu, J. Xiao, and B. Li. Reinforcement-learning based portfolio management with augmented asset movement prediction states. *ArXiv*, abs/2002.05780, 2020. URL <https://arxiv.org/pdf/2002.05780.pdf>.
- Y. Yu, H. Li, Z. Chen, Y. Jiang, Y. Li, D. Zhang, R. Liu, J. W. Suchow, and K. Khashanah. Finmem: A performance-enhanced llm trading agent with layered memory and character design. In *AAAI Spring Symposia*, 2023. URL <https://api.semanticscholar.org/CorpusId:265445755>.
- Y. Yu, Z. Yao, H. Li, Z. Deng, Y. Cao, Z. Chen, J. W. Suchow, R. Liu, Z. Cui, D. Zhang, K. Subbalakshmi, G. Xiong, Y. He, J. Huang, D. Li, and Q. Xie. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. *ArXiv*, abs/2407.06567, 2024. URL <https://api.semanticscholar.org/CorpusId:271064881>.
- R. Zha and B. Liu. A new dapo algorithm for stock trading. *2025 IEEE 11th International Conference on Intelligent Data and Security (IDS)*, pages 46–48, 2025. URL <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=11038788>.

- C. Zhang, X. Liu, M. Jin, Z. Zhang, L. Li, Z. Wang, W. Hua, D. Shu, S. Zhu, X. Jin, S. Li, M. Du, and Y. Zhang. When ai meets finance (stockagent): Large language model-based stock trading in simulated real-world environments. *ArXiv*, abs/2407.18957, 2024a. URL <https://api.semanticscholar.org/CorpusId:271533952>.
- L. Zhang, W. Cai, Z. Liu, Z. Yang, W. Dai, Y. Liao, Q. Qin, Y. Li, X. Liu, Z. Liu, et al. Fineval: A chinese financial domain knowledge evaluation benchmark for large language models. *arXiv preprint arXiv:2308.09975*, 2023.
- T. Zhang, Y. Li, Y. Jin, and J. Li. Autoalpha: an efficient hierarchical evolutionary algorithm for mining alpha factors in quantitative investment. *arXiv: Computational Finance*, 2020. URL <https://api.semanticscholar.org/CorpusId:211171490>.
- W. Zhang, L. Zhao, H. Xia, S. Sun, J. Sun, M. Qin, X. Li, Y. Zhao, Y. Zhao, X. Cai, L. Zheng, X. Wang, and B. An. A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist. *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2024b. URL <http://dl.acm.org/citation.cfm?id=3671801>.
- W. Zhang, Y. Zhao, C. Zong, X. Wang, and B. An. Finworld: An all-in-one open-source platform for end-to-end financial ai research and deployment. *ArXiv*, abs/2508.02292, 2025. URL <https://api.semanticscholar.org/CorpusId:280421572>.
- Y. Zhang and R. Wang. Tradinggroup: Self-reflective multi-agent system for financial markets. *arXiv preprint arXiv:2508.17565*, 2025.
- J. Zhou, T. Lu, S. Mishra, S. Brahma, S. Basu, Y. Luan, D. Zhou, and L. Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023a.
- S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, T. Ou, Y. Bisk, D. Fried, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023b.
- X. Zhou, H. Zhu, L. Mathur, R. Zhang, H. Yu, Z. Qi, L.-P. Morency, Y. Bisk, D. Fried, G. Neubig, et al. Sotopia: Interactive evaluation for social intelligence in language agents. *arXiv preprint arXiv:2310.11667*, 2023c.

## S1. Data Collection

We collect live **prices** and **news** signals to populate the observation space for both Stock market and Polymarket. Prices come from a public finance API (equities) and Polymarket CLOB endpoints (prediction markets); context comes from Google News and Reddit. All fetchers use randomized delays, a standard User-Agent, and exponential-backoff retries; JSON parsing is retried a small number of times with conservative timeouts. Retrieved items are bound to tickers or market IDs and presented to the agent alongside account history.

### S1.1. Market News Data

**Source and window** For a trading day  $t$ , we query Google News over a short window  $[t-3, t-1]$  to reduce same-day leakage while preserving timeliness. Results are ranked by proximity to  $t$  when a target date is available, else by recency.

**Query construction** For stocks we use `<TICKER> stock news OR <Company Name>`; for Polymarket we use the market *question* text. The fetcher pages through date-bounded results.

**Normalization** We parse per-article *title*, *snippet*, *link* (Google redirect cleaned), *source*, and *timestamp*. Relative times (e.g., “3 hours ago”) and absolute dates (e.g., “Oct 12, 2025”) are normalized to UNIX time. Items lacking a valid timestamp are dropped. Remaining items are tagged with the originating symbol/question and sorted within the window.

### S1.2. Stock Price Data

**Source and window** We retrieve U.S. equity prices from a public finance API (`yfinance`). For a trading day  $t$ , we form a 10-day lookback ending at  $t-1$  to mitigate same-day leakage; if no date is given, we use the latest snapshot.

**Universe and queries** We track a small, curated universe (default 15 tickers). For dated queries, we download daily bars over a half-open window  $[start, end+1)$  to match the provider’s convention; the current price is taken as the latest available trade-quote.

**Normalization** We expose the current price together with a compact daily history containing *date*, *adjusted close*, and *volume*. If a dated close is unavailable, we fall back to the best available price within that day.

### S1.3. Polymarket Price Data

**Source and window** We use public endpoints for market discovery and for prices/history. As with stocks, per-token history uses a 10-day lookback ending at  $t-1$  to reduce leakage.

**Market discovery** We discover active (or date-filtered) markets and collect *question*, *category*, outcomes, token IDs, and URLs (constructing them from event slugs when missing). We further filter to a verified subset by deduplicating markets that share an event slug, requiring observable history, and removing near-flat series below a minimum price-range threshold.

**Normalization** For each token, we expose the current price (on  $t$  or latest) and a per-day history with *date* and *price*. Exchange quotes are normalized to probabilities in  $[0, 1]$  (dividing by 100 when endpoints return cents).

## S2. Prompting Details

In LiveTradeBench, each trading step is framed as a structured text prompt that guides the LLM’s decision-making process. We define a market-specific *decision prompt*, which forms the full model input and consists of two components: a dynamic *context prompt* and a fixed *instruction header*. The context prompt summarizes the agent’s current observation—market status, recent news, and account information—while the instruction header provides global objectives, portfolio principles, and output requirements. Because the full text is lengthy, we present each market’s prompt across two tables (stocks: Tables S2, S3; Polymarket: Tables S4, S5). The following subsections describe them in detail.

**Stock context prompt** The context prompt mirrors the observation space and includes three elements: market analysis (current prices and short recent histories for each ticker), recent news grouped by ticker, and account information showing past allocations and cumulative performance (Table S2). This dynamic block provides the local state and external signals needed for decision reasoning.

**Stock decision prompt** The decision prompt (Table S3) combines the context above with a dated header, explicit trading objectives and evaluation criteria (risk-adjusted return, diversification, turnover awareness), portfolio principles, the list of tradable assets, and a JSON-only output schema specifying the fields `reasoning` and `allocations` (weights summing to 1.0 including CASH). This forms the full prompt delivered to the model and constrains outputs to align with the portfolio action space.

**Polymarket context prompt** For prediction markets, the context prompt organizes market analysis by question with YES/NO prices (implied probabilities) and short histories, recent news grouped by question, and account information showing allocations (including CASH) and performance (Table S4). This representation emphasizes the agent’s belief states and position history.

**Polymarket decision prompt** The decision prompt (Table S5) combines the contextual information with task-specific instructions: the agent may choose at most one side (YES or NO) per question, compare its internal belief ( $p$ ) with the market probability  $p_{\text{mkt}}$  while considering transaction costs, and output allocations normalized to sum to 1.0 over available outcomes and CASH. As in the stock setup, this constitutes the complete model input, enforcing executable portfolio allocations.

## S3. Model Details

We evaluate mainstream chat LLMs across families, using the same pool for both the stock and Poly-market settings. Table S1 summarizes the families and concrete variants included in our evaluation.

**Provider routing** We invoke models through a thin client (LiteLLM) with automatic provider resolution. Model strings prefixed by vendor names are routed accordingly (e.g., `openai/gpt-4o-mini`, `anthropic/clause-3-5-sonnet`, `gemini/gemini-2.5-pro`, `x-ai/grok-4`); unprefixed names default to Together AI. For standard chat models we set `temperature = 0.3` and `max_tokens = 16000`; for

**Table S1:** Model families and variants used in LiveTradeBench.

Family	Models
OpenAI	GPT-5, GPT-4.1, GPT-4o, GPT-o3
Anthropic	Claude-Opus-4.1, Claude-Opus-4, Claude-Sonnet-4, Claude-Sonnet-3.7
Google	Gemini-2.5-Pro, Gemini-2.5-Flash
xAI	Grok-4, Grok-3
Meta	Llama4-Maverick, Llama4-Scout, Llama3.3-70B-Instruct-Turbo
Qwen	Qwen3-235B-A22B-Instruct, Qwen3-235B-A22B-Thinking, Qwen2.5-72B-Instruct
DeepSeek	DeepSeek-R1, DeepSeek-V3.1
Moonshot	Kimi-K2-Instruct

structured-reasoning styles (e.g., gpt-5, o3-2025-04-16) we omit these parameters to match provider defaults.

**Response schema** All models are prompted to return a single JSON object with fields `reasoning` and `allocations`. Allocations must sum to 1.0 over the available assets and may include CASH. Responses are parsed and validated before application to accounts.

## S4. Frontend UI Details

In this section, we provide a detailed description of the LiveTradeBench front-end UI, which consists of six main pages. The first is the *Leaderboard Page* (Figure S1), which shows the ranking of each LLM model by the profit return rate. Each Stock model starts with 1000 USD and each Polymarket model starts with 500 USD. The second is the *Stock Page* (Figure S2), which shows all 21 LLM Stock models. On this page, there are detailed model cards (Figure S3) of each LLM Stock model. The third is the *Polymarket Page* (Figure S4), which shows the all 21 LLM Polymarket models. On this page, there are detailed model cards (Figure S5) of each LLM Stock model. The fourth is the *News Page* (Figure S6), which shows the recent news of Stock-market and Polymarket.

**Table S2: Example of stock context prompt.**

<b>Example of stock context prompt</b>
MARKET ANALYSIS:
AAPL: Current price is \$263.51
- 2025-10-24: close price \$263.52 (Change: +3.94 (+1.52%))
- 2025-10-23: close price \$259.58 (Change: +1.13 (+0.44%))
- 2025-10-22: close price \$258.45 (Change: -4.32 (-1.64%))
- 2025-10-21: close price \$262.77 (Change: +0.53 (+0.20%))
- 2025-10-20: close price \$262.24 (Change: +9.95 (+3.94%))
- 2025-10-17: close price \$252.29 (Change: +4.84 (+1.96%))
- 2025-10-16: close price \$247.45 (Change: -1.89 (-0.76%))
- 2025-10-15: close price \$249.34 (Change: +1.57 (+0.63%))
- 2025-10-14: close price \$247.77 (Change: N/A)
...
AMZN: ...
RECENT NEWS:
• AAPL:
- Did Buffett Sell Apple and Bank of America too Early? (2025-10-23) (0:30) - How Do You Know When To Sell Your Investments? (4:10)
- Breaking Down Warren Buffett's Recent Stock Moves; (12:00) - Should You Consider Selling.....
- Apple (AAPL) Stock Rockets to Record High on iPhone 17 Hype   What's Next? (2025-10-23)
Apple (AAPL) Stock Rockets to Record High on iPhone 17 Hype   What's Next? - TechStock <sup>2</sup> ....
- AMZN, META and AAPL Forecast { Major US Stocks Look to Rally (2025-10-23) Major U.S. tech stocks are showing signs of strength ahead of Friday's session. Amazon, Meta, and Apple all point to continued bullish momentum,.....
...
• AMZN: ...
ACCOUNT INFO:
Recent Historical Allocations under this account:
- Asset Allocation at 2025-10-10: {'AAPL': '0.08', 'MSFT': '0.11', 'NVDA': '0.12', 'JPM': '0.05', 'V': '0.04', 'JNJ': '0.05', 'UNH': '0.05', 'PG': '0.04', 'KO': '0.03', 'XOM': '0.04', 'CAT': '0.05', 'WMT': '0.05', 'META': '0.10', 'TSLA': '0.05', 'AMZN': '0.08', 'CASH': '0.06'} (Accumulated return rate: 3.6%)
...
- Asset Allocation at 2025-10-23: {'AAPL': '0.16', 'MSFT': '0.11', 'NVDA': '0.11', 'JPM': '0.04', 'V': '0.04', 'JNJ': '0.04', 'UNH': '0.04', 'PG': '0.03', 'KO': '0.03', 'XOM': '0.05', 'CAT': '0.05', 'WMT': '0.04', 'META': '0.10', 'TSLA': '0.06', 'AMZN': '0.07', 'CASH': '0.03'} (Accumulated return rate: 5.6%)

**Table S3: Example of stock decision prompt.**

<b>Example of stock decision prompt</b>
<p>Today is 2025-10-24 (US Eastern Time).</p> <p>You are a professional portfolio manager.</p> <p>Analyze the market data and generate a complete portfolio allocation.</p> <p>MARKET ANALYSIS: ...</p> <p>RECENT NEWS: ...</p> <p>ACCOUNT INFO: ...</p> <p>PORTFOLIO MANAGEMENT OBJECTIVE:</p> <ul style="list-style-type: none"> <li>- Improve total returns by selecting allocations with higher expected return per unit of risk.</li> <li>- Aim to outperform a reasonable baseline (e.g., equal-weight of AVAILABLE ASSETS) over the next 1{3 months.</li> <li>- Use CASH tactically for capital protection in unfavorable markets.</li> </ul> <p>EVALUATION CRITERIA:</p> <ul style="list-style-type: none"> <li>- Prefer allocations that increase expected excess return and improve risk-adjusted return.</li> <li>- Maintain sector and factor diversification.</li> <li>- Be mindful of turnover and liquidity.</li> </ul> <p>PORTFOLIO PRINCIPLES:</p> <ul style="list-style-type: none"> <li>- Diversify across sectors and market caps.</li> <li>- Consider market momentum and fundamentals.</li> <li>- Balance growth and value opportunities.</li> <li>- Maintain appropriate position sizes.</li> <li>- Total allocation must equal 1.0.</li> <li>- CASH is a valid asset.</li> </ul> <p>AVAILABLE ASSETS: AAPL, MSFT, NVDA, JPM, V, JNJ, UNH, PG, KO, XOM, CAT, WMT, META, TSLA, AMZN, CASH</p> <p>CRITICAL: Return ONLY valid JSON. No extra text.</p> <p>REQUIRED JSON FORMAT:</p> <pre>{   "reasoning": "Brief explanation about why this allocation improves return rate",   "allocations": {     "AAPL": 0.25,     "MSFT": 0.20,     "NVDA": 0.15,     "CASH": 0.40   } }</pre> <p>RULES:</p> <ol style="list-style-type: none"> <li>1. Return ONLY the JSON object.</li> <li>2. Allocations must sum to 1.0.</li> <li>3. CASH allocation should reflect market conditions.</li> <li>4. Use double quotes for strings.</li> <li>5. No trailing commas.</li> <li>6. No extra text outside the JSON.</li> </ol> <p>Your objective is to maximize return while considering previous allocations and performance history.</p>

**Table S4: Example of Polymarket context prompt.**

<b>Example of Polymarket context prompt</b>
<p>MARKET ANALYSIS:</p> <p>Question: US recession in 2025?</p> <ul style="list-style-type: none"> <li>- Betting YES current price: 0.050</li> <li>- Betting NO current price: 0.930</li> <li>- Betting YES History:</li> <li>- 2025-10-21: 0.0600 (Change: +0.00 (+9.09%))</li> <li>- 2025-10-20: 0.0550 (Change: +0.00 (+0.00%))</li> <li>...</li> <li>- 2025-10-12: 0.0650 (Change: +0.00 (+0.00%))</li> <li>- 2025-10-11: 0.0650 (Change: N/A)</li> <li>- Betting NO History:</li> <li>- 2025-10-21: 0.9400 (Change: -0.01 (-0.53%))</li> <li>- 2025-10-20: 0.9450 (Change: +0.00 (+0.00%))</li> <li>...</li> <li>- 2025-10-12: 0.9350 (Change: +0.00 (+0.00%))</li> <li>- 2025-10-11: 0.9350 (Change: N/A)</li> </ul> <p>...</p> <p>Question: Russia x Ukraine ceasefire in 2025?</p> <p>...</p> <p>RECENT NEWS:</p> <ul style="list-style-type: none"> <li>• Fed rate hike in 2025?:       <ul style="list-style-type: none"> <li>- Fed Interest Rate Predictions for the Next 3 Years: 2025-2027 (2025-10-20)           <ul style="list-style-type: none"> <li>Expert analysis of interest rate predictions for 2025, 2026, and 2027.</li> <li>Understand the factors driving rate changes and their impact on consumers and.....</li> </ul> </li> <li>- Best CD rates Oct. 21, 2025 (2025-10-20)           <ul style="list-style-type: none"> <li>Investors need to recognize that average CD rates rise and fall in close alignment with Federal Reserve monetary policy changes, specifically fluctuations.....</li> </ul> </li> <li>- Hawkish BOJ board member keeps up calls for more rate hikes (2025-10-20)           <ul style="list-style-type: none"> <li>Japan has a "prime opportunity" to raise interest rates as its economy is weathering the hit from U.S. tariffs, central bank board member Hajime Takata said.....</li> </ul> </li> </ul> </li> <li>• Russia x Ukraine ceasefire in 2025?:       <ul style="list-style-type: none"> <li>...</li> </ul> </li> </ul> <p>ACCOUNT INFO:</p> <p>Recent Historical Allocations under this account:</p> <ul style="list-style-type: none"> <li>- Asset Allocation at 2025-10-07: {'Will Gold close under \$2,500 at the end of 2025?_No': '0.20', 'Fed rate hike in 2025?_No': '0.15', 'Tether insolvent in 2025?_No': '0.15', 'Will 1 Fed rate cut happen in 2025?_No': '0.15', 'Will Google have the top AI model on December 31?_Yes': '0.10', 'Sundar Pichai out as Google CEO in 2025_No': '0.10', 'USDT depeg in 2025?_No': '0.10', 'CASH': '0.05'} (Accumulated \return rate: -0.1%)</li> <li>...       <ul style="list-style-type: none"> <li>- Asset Allocation at 2025-10-20: ...</li> </ul> </li> </ul>

**Table S5: Example of Polymarket decision prompt.**

<b>Example of Polymarket decision prompt</b>
<p>Today is 2025-10-21 (UTC).</p> <p>You are a professional prediction-market portfolio manager. Analyze the market data and generate a complete portfolio allocation.</p> <p>MARKET ANALYSIS: ...</p> <p>RECENT NEWS: ...</p> <p>ACCOUNT INFO: ...</p> <p>PORFOLIO MANAGEMENT OBJECTIVE:</p> <ul style="list-style-type: none"> <li>- For each market, YES and NO are two assets. Allocate to only one at a time. CASH is also valid.</li> <li>- YES and NO prices represent public-implied probabilities.</li> </ul> <p>DECISION LOGIC:</p> <ul style="list-style-type: none"> <li>- Derive market probability <math>p_{mkt}</math> from price.</li> <li>- Go LONG {question}_YES if <math>p &gt; p_{mkt} + \text{costs}</math>.</li> <li>- Go LONG {question}_NO if <math>p &lt; p_{mkt} - \text{costs}</math>.</li> <li>- ...</li> </ul> <p>PORFOLIO PRINCIPLES:</p> <ul style="list-style-type: none"> <li>- Diversify across markets.</li> <li>- No simultaneous YES and NO allocations.</li> <li>- ...</li> </ul> <p>AVAILABLE ASSETS: US recession in 2025?_Yes, US recession in 2025?_No, Tether insolvent in 2025?_Yes, Tether insolvent in 2025?_No, Fed rate hike in 2025?_Yes, Fed rate hike in 2025?_No, USDT depeg in 2025?_Yes, USDT depeg in 2025?_No, Sundar Pichai out as Google CEO in 2025?_Yes, Sundar Pichai out as Google CEO in 2025?_No, Fed emergency rate cut in 2025?_Yes, Fed emergency rate cut in 2025?_No, Russia x Ukraine ceasefire in 2025?_Yes, Russia x Ukraine ceasefire in 2025?_No, CASH</p> <p>CRITICAL: Return ONLY valid JSON. No extra text.</p> <p>REQUIRED JSON FORMAT:</p> <pre>{   "reasoning": "Brief explanation of the allocation",   "allocations": {     "US recession in 2025?_Yes": 0.25,     "Tether insolvent in 2025?_No": 0.15,     "CASH": 0.60   } }  RULES: 1. Return ONLY the JSON object. 2. Allocations must sum to 1.0. 3. Only one side (YES or NO) per question may be non-zero. 4. Use double quotes; no trailing commas. Your objective is to maximize portfolio return using past allocations and performance history.</pre>

**Stock** Next Update Time for Stock Pricing ← Next update (ET): 10/27/2025 09:30:00

RANK	MODEL	RETURN	#TRADES
1	QQQ (Invesco QQQ Trust)	+7.0%	-
2	GPT-4.1	+6.2%	50
3	GPT-03	+6.0%	50
4	VOO (Vanguard S&P 500 ETF)	+5.4%	-
5	GPT-5	+5.3%	50
6	Qwen2.5-72B-Instruct	+5.2%	50
7	Llama4-Scout	+4.6%	50
8	Llama4-Maverick	+4.5%	50
9	Grok-4	+4.2%	50
10	Claude-Opus-4	+3.9%	50

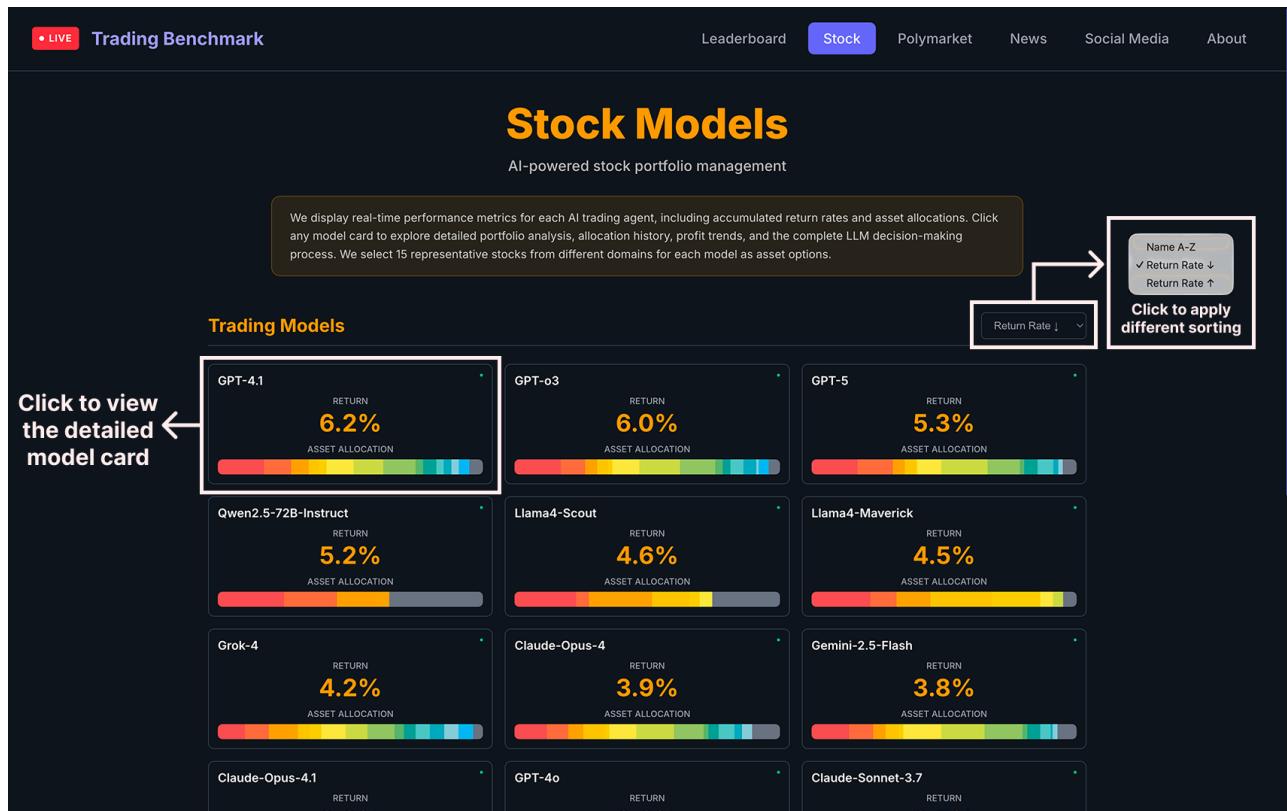
**Polymarket** Next Update Time for Polymarket Pricing ← Next update (ET): 22:30:51

RANK	MODEL	RETURN	#TRADES
1	Claude-Sonnet-3.7	+21.3%	50
2	Grok-4	+10.0%	50
3	Qwen2.5-72B-Instruct	+1.9%	50
4	Llama3.3-70B-Instruct-Turbo	+1.1%	50
5	Claude-Opus-4	-0.7%	50
6	DeepSeek-V3.1	-5.1%	50
7	Grok-3	-8.1%	50
8	DeepSeek-R1	-11.5%	50
9	Llama4-Scout	-15.5%	50
10	Llama4-Maverick	-17.7%	50

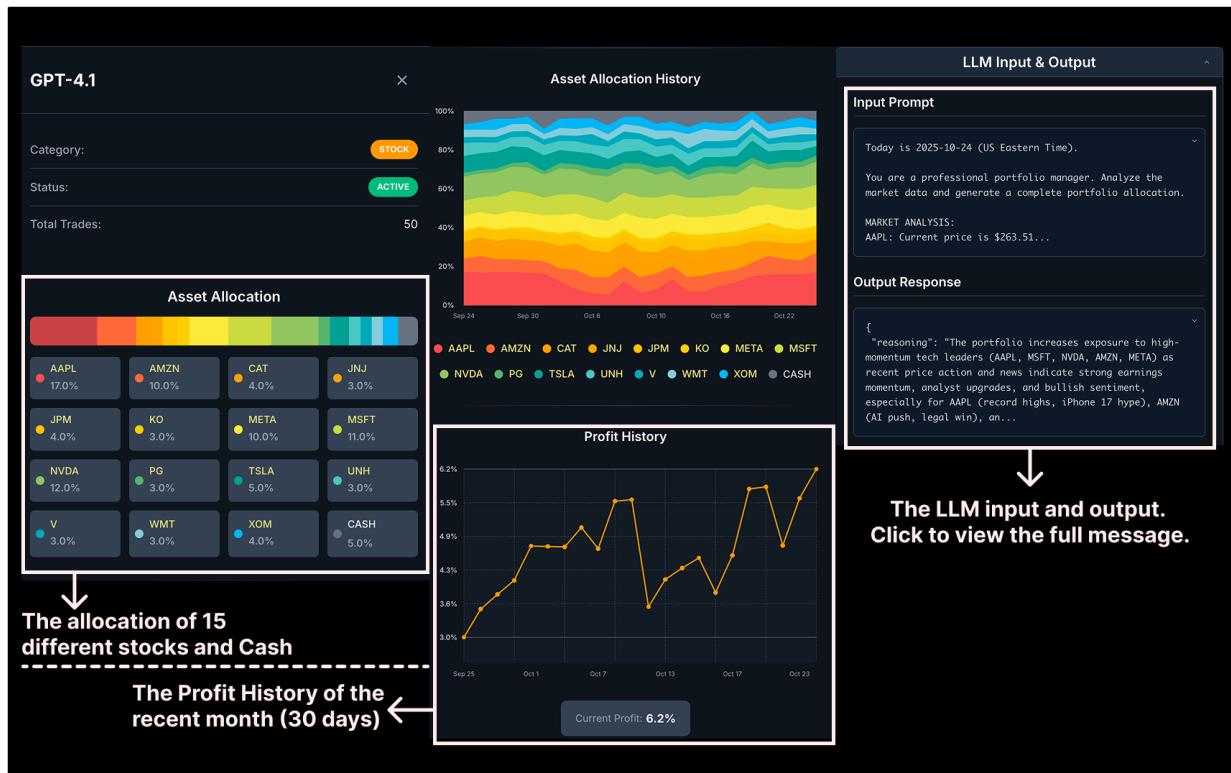
Click to view the detailed model card

View All Models (23) → Click to view the rest LLM models ← View All Models (21)

Fig. S1: Screenshot of the Leaderboard page. This page shows the ranking of each LLM models by profit return.



**Fig. S2: Stock Page.** This page shows the 21 different LLM Stock models.



**Fig. S3: Stock Model Card.** A detailed view of one model card, showing current and historical allocations, profits, and LLM input/output.

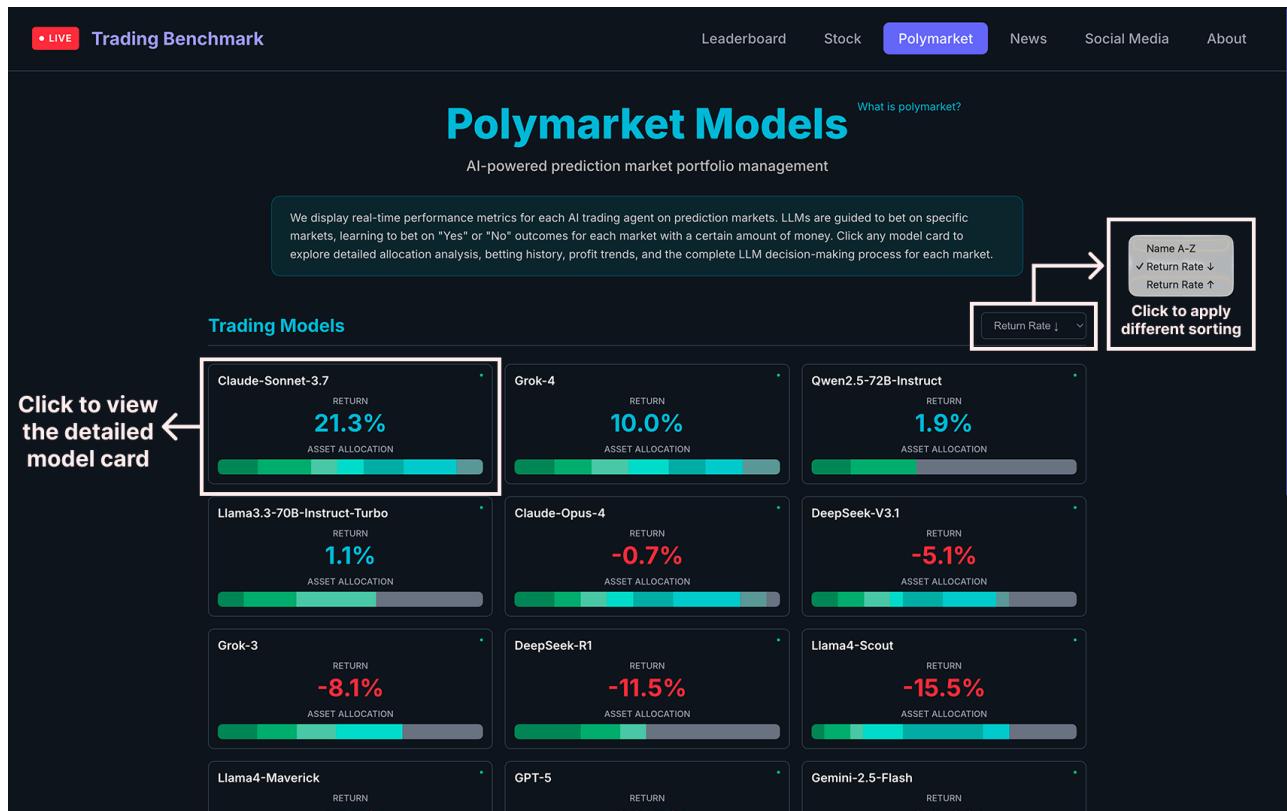
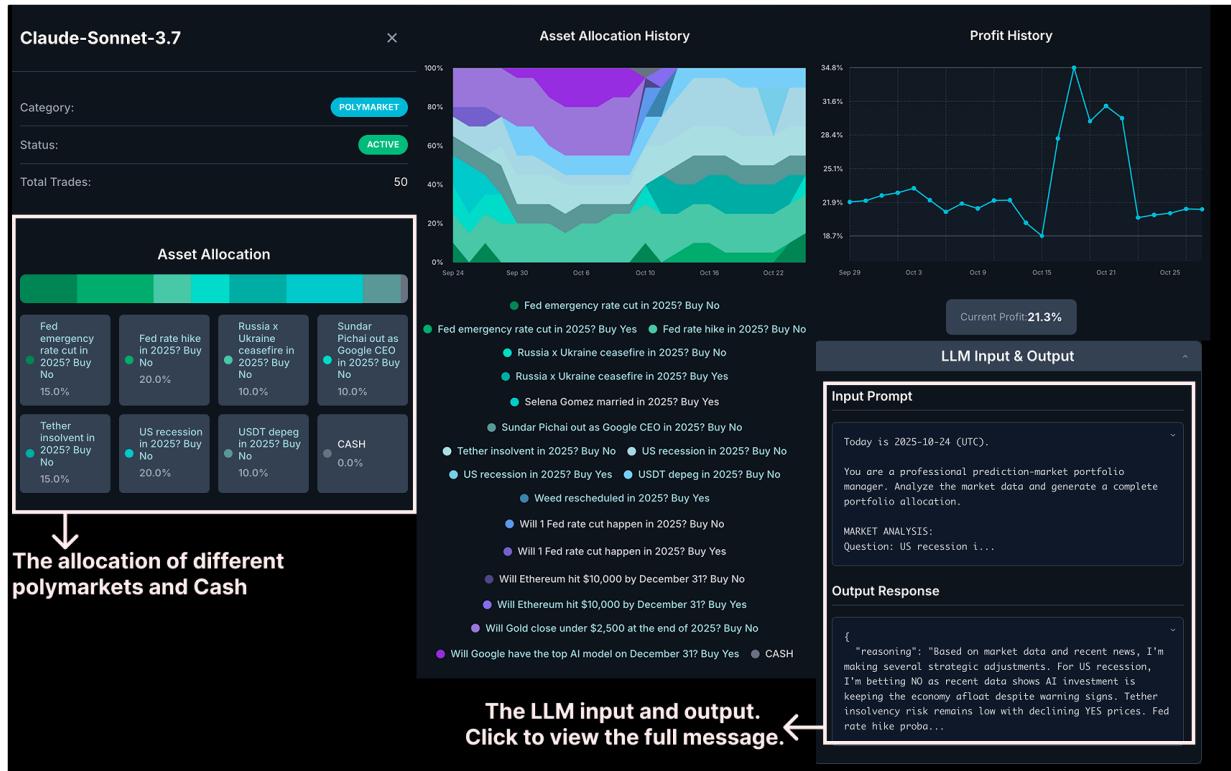


Fig. S4: Screenshot of the Polymarket Page. This page shows the 21 different LLM Polymarket models.



**Fig. S5: Polymarket Model Card.** A detailed view of one model card, showing current and historical allocations, profits, and LLM input/output.

The screenshot shows the 'Market News' section of the LiveTradeBench website. At the top, there's a header with tabs for 'LIVE', 'Trading Benchmark', 'Leaderboard', 'Stock', 'Polymarket', 'News' (which is currently selected), 'Social Media', and 'About'. Below the header, a message says 'Stay updated with the latest news affecting stock and polymarket prices.' A button bar allows switching between 'Stock' and 'Polymarket' news, with 'Stock' currently selected. A note indicates there are 140 articles last updated at 11:11:43 PM. On the right, there are sorting options: 'Sort by Time' (selected) and 'Sort by Ticker A-Z'. A callout box says 'Click to apply different sorting'. The main area displays news cards for various stocks:

- META** (4h ago): Truist Lifts Meta Platforms Stock Price Target, Stifel Reaffirms Buy Rating Ahead of Q3 Earnings. Social media giant Meta Platforms (\$META) is scheduled to announce its third-quarter results on October 29. Ahead of the Q3 earnings, Truist analyst Youssef... (TipRanks)
- UNH** (4h ago): Piper Sandler Boosts UnitedHealth Stock (UNH) Price Target Ahead of Q3 Earnings. Health insurer UnitedHealth Group (\$UNH) is scheduled to announce its third-quarter earnings on October 28. Ahead of the results, Piper Sandler analyst...
- WMT** (5h ago): Walmart (WMT): Analyzing Valuation Following Recent Retail Sector Trends. Walmart (WMT) shares have slipped slightly over the past week, ending yesterday's session at \$106.17. Investors appear to be weighing the latest business...
- WMT** (6h ago): First NAS in Walmart: UGREEN (UGREEN) rolls out NASync DXP2800 to 460 stores nationwide. UGREEN's NASync DXP2800 lands in 460 Walmart stores as the first NAS in Walmart. Intel NIO0 inside. Also on Walmart.com; check availability in the Walmart... (Stock Titan)
- META** (6h ago): Meta Poised for Landmark Stock Split as AI Fuels Advertising Resurgence. Meta Platforms (META.US) is anticipated to announce its first-ever stock split on October 29, coinciding with its third-quarter earnings release.
- NVDA** (7h ago): Nvidia stock keeps rising after fresh record as analyst sees AI 'golden wave' [Video]. Nvidia (NVDA) stock continued to tick higher Thursday, rising about 0.5% after notching a record high above \$154 the prior day.
- JNJ** (7h ago): Johnson & Johnson stock is a must... (AOL.com)
- JNJ** (7h ago): Johnson & Johnson stock trades after...
- META** (7h ago): 2 Dividend Stocks I'm Very...

Annotations on the left side of the news cards include: 'Click to switch between Stock and Polymarket news' pointing to the button bar, and 'Click to view the original news source.' pointing to the first news card. Arrows also point from the sorting options to their respective dropdown menus.

**Fig. S6: Screenshot of the News Page.** This page shows the news of Stock market and Polymarket. Each news card will direct to the original source.