

# Behavioral Consistency Validation for LLM Agents: An Analysis of Trading-Style Switching through Stock-Market Simulation

Zeping Li<sup>1</sup>, Guancheng Wan<sup>2</sup>, Keyang Chen<sup>1</sup>, Yu Chen<sup>3</sup>, Yiwen Zhao<sup>1</sup>,  
Philip Torr<sup>4</sup>, Guangnan Ye<sup>1</sup>, Zhenfei Yin<sup>4,†</sup>, Hongfeng Chai<sup>1,†</sup>,

<sup>1</sup>Fudan University   <sup>2</sup>UCLA   <sup>3</sup>BNP Paribas   <sup>4</sup>Oxford University

<sup>†</sup>Corresponding authors

## Abstract

Recent works have increasingly applied Large Language Models (LLMs) as agents in financial stock market simulations to test if micro-level behaviors aggregate into macro-level phenomena. However, a crucial question arises: Do LLM agents’ behaviors align with real market participants? This alignment is key to the validity of simulation results. To explore this, we select a financial stock market scenario to test behavioral consistency. Investors are typically classified as fundamental or technical traders, but most simulations fix strategies at initialization, failing to reflect real-world trading dynamics. In this work, we assess whether agents’ strategy switching aligns with financial theory, providing a framework for this evaluation. We operationalize four behavioral-finance drivers—loss aversion, herding, wealth differentiation, and price misalignment—as personality traits set via prompting and stored long-term. In year-long simulations, agents process daily price-volume data, trade under a designated style, and reassess their strategy every 10 trading days. We introduce four alignment metrics and use Mann–Whitney U tests to compare agents’ style-switching behavior with financial theory. Our results show that recent LLMs’ switching behavior is only partially consistent with behavioral-finance theories, highlighting the need for further refinement in aligning agent behavior with financial theory.

## 1 Introduction

Agent-Based Modeling (ABM) has a long tradition in economics and the social sciences for explaining how micro-level behavioral rules generate macro-level phenomena (Schelling, 1971; Epstein and Axtell, 1996; Bonabeau, 2002). Classic ABMs typically rely on handcrafted heuristics and fixed rule sets, which improve interpretability yet, by construction, restrict agents’ perceptual and decision spaces, thereby constraining system heterogeneity. Recent progress in LLMs offers an alternative

way to instantiate agents in ABM systems (Gao et al., 2024a; Lu et al., 2024). LLMs are capable of processing unstructured and structured signals, reasoning over extended context, and producing explicit rationales for actions (Zhao et al., 2023). Through prompting, LLM-based agents can be configured to express diverse personality traits, expanding agent-level heterogeneity in roles and decision styles (Shao et al., 2023; Li et al., 2023a). Incorporating LLMs as agents thus yields simulations that more closely approximate complex real-world systems. As a result, LLM-enabled ABM is now increasingly used for simulation across complex social systems, game-theoretic environments, and financial markets (Park et al., 2023; Kovařík et al., 2025; Zhang et al., 2024b).

In financial markets, there has been a substantial body of work on simulations using LLM agents (Wu et al., 2023; Lai et al., 2024; Yang et al., 2025; Lopez-Lira, 2025), where these agents interpret financial news and corporate disclosures to generate trading strategies based on price and volume histories. These strategies are often converted into standard fundamental and technical indicators (Xiao et al., 2024). However, these studies typically fix investors as either fundamental or technical traders at initialization, with no strategy switching occurring during the simulation. This limitation contrasts with traditional financial studies (Franke and Westerhoff, 2012), which show that a larger proportion of technical traders amplifies boom-bust cycles, while a greater share of fundamental traders stabilizes prices. Thus, explicit style switching is crucial to reproduce stylized facts such as volatility clustering, long-range dependence, and heavy-tailed returns. Before incorporating style-switching behavior into simulations, we pose a fundamental and critical question: under realistic information and constraints, can the style-switching behavior of LLM agents align with that of real market participants?

To this end, we examine the behavioral consistency of LLM agents in their style-switching behavior. Drawing from four key behavioral finance drivers—Loss aversion tendency (Kaineman et al., 1979), Herding tendency (Banerjee, 1992), Wealth differentiation sensitivity (Hommes, 2006), and Price misalignment sensitivity (Campbell and Shiller, 1988)—we map these factors to agent-level personality traits through prompting and embed them as long-term memories that guide strategy formation and style-switching decisions. The evaluation runs on a stock-market simulator using S&P 500 stocks in 2024, combining daily price–volume data with curated trading indicators. Each round, agents operate under a designated style, and a counterfactual ledger tracks the alternative; every ten trading days, agents review returns and traits to decide whether to switch styles.

To assess the behavioral fidelity of LLM agents, we pose four research questions aligned with behavioral finance theories: **(RQ1)** Does loss aversion keep agents in their current style after losses? **(RQ2)** Does herding increase switching when the alternative style has a larger population share? **(RQ3)** Does wealth differentiation trigger switching when the alternative style outperforms? **(RQ4)** Does price misalignment induce shifts toward the fundamental style when market prices diverge from estimated value? We validate these questions by comparing agents with aligned and non-aligned traits using Mann–Whitney U tests. The results show that LLM agents’ switching behavior is only partially consistent with behavioral-finance theories. In summary, the main contributions of this paper are as follows:

- To our knowledge, we are the first to identify style switching in LLM agents within market simulations. We consider four behavioral drivers that influence style transitions and embed them into the agents’ long-term memory.
- We propose four alignment metrics and use Mann–Whitney U tests to compare aligned and non-aligned cohorts, providing an operational test for the consistency of LLM agents’ style-switching behavior with financial theory.
- Year-long simulations conducted with various recent LLM agents show that while LLM agents exhibit partial consistency with behavioral-finance theories, they cannot fully align with these theories in all aspects.

## 2 Related Works

### 2.1 Traditional ABM Simulation

Traditional agent-based models have been used in finance for decades to link heterogeneous trading rules with price formation and market statistics. Early artificial–stock-market work shows that adaptive, learning agents in centralized venues generate realistic return dynamics and position updates (Palmer et al., 1994), with subsequent formulations of endogenous expectations closing the loop from beliefs to orders and prices (Arthur et al., 2018). Chartist–fundamentalist multi-agent models reproduce stylized facts such as fat tails and clustered volatility (Lux and Marchesi, 1999), while mechanism studies connect executable rules (value, trend following, market making) to volatility and liquidity (Farmer and Joshi, 2002); microsimulations that embed heterogeneity, peer interactions, and trade frictions account for serial-dependence patterns (Iori, 2002). Taken together, ABMs’ explicit heterogeneity, interaction, and microstructure realism make them well suited to simulate financial markets end-to-end and to validate how micro behavior aggregates into macro phenomena (Axtell and Farmer, 2025).

### 2.2 LLM-Based ABM Simulation

Recent advances in LLMs have rapidly increased the use of LLM-enabled ABM in finance and economics, offering more realistic and scalable simulations than traditional systems (Gao et al., 2024a; Guo et al., 2024). This progress is driven by three core capabilities of LLM agents: (i) processing rich textual information with price–volume signals, (ii) maintaining persona-like memory to sustain heterogeneous behavior, and (iii) emitting structured actions (e.g., function-call orders) integrated with market microstructure. Building on this, (Li et al., 2023b) embeds LLM agents in ABM economies to produce interpretable, heterogeneous multi-period behavior, while (Lin et al., 2025) scales to thousands of agents to replicate macro expectation formation (inflation, unemployment).

In financial stock markets, Yang et al. (2025) integrates a networked communication layer within a BDI-style framework, tracing information propagation through trading decisions to emergent macro phenomena; it reproduces mechanisms like bubbles and drawdowns, enabling scalable behavioral–social simulations. Similarly, (Lopez-Lira, 2025) builds a continuous double auction with

function-call decisions, and (Zhang et al., 2024a) and (Gao et al., 2024b) provide multi-agent trading in realistic environments. (Yu et al., 2024) enhances financial decision-making, and (Vidler and Walsh, 2025) explores controllable agent heterogeneity through preference-based prompts (e.g., risk or ambiguity aversion). However, a recent study (Henning et al., 2025) finds that markets populated by LLM agents often appear "too rational" compared to human-subject experiments, raising doubts about their behavioral fidelity in stock-market simulations. This paper thus focuses on testing whether LLM agents' style-switching behavior aligns with financial theory.

### 2.3 Behavioral Drivers of Style Switching

To understand style switching, we start with (Brock and Hommes, 1997, 1998), who model agents' switches between forecasting rules (e.g., fundamentalist versus trend-following) via a discrete-choice mechanism based on relative performance, providing a foundation for endogenous switching. Extending this line of work, Franke and Westerhoff (2012) model switching between technical and fundamental traders via a relative "attractiveness" index that aggregates predisposition, herding, mispricing, and wealth differentiation, thereby supporting the model's reproduction of key financial stylized facts.

Guided by these works, we select four theory-grounded factors that map directly onto our research questions. Specifically, loss aversion (Kai-Ineman et al., 1979) predicts persistence in the current style even after recent losses (RQ1); herding (Banerjee, 1992) implies a higher likelihood of switching into the majority style when population shares tilt toward it (RQ2); wealth differentiation (Hommes, 2006) suggest switching toward the counterfactual style as the counterfactual style's recent outperformance widens (RQ3); and price misalignment (Campbell and Shiller, 1988) predicts movement toward the fundamental style as the gap between market price and fundamental value grows (RQ4). This mapping transforms classic behavioral-finance mechanisms into testable predictions about when and how LLM agents should switch styles. By modeling style-switching behavior around these key behavioral drivers, we aim to address the gap in the literature regarding LLM agents' ability to replicate such behavior, particularly their capacity to switch strategies in response to these cues.

## 3 Method

### 3.1 Stock Market Simulation Setup

**Overall setup.** We collect a full-year 2024 dataset for the S&P 500 constituents, integrating two modalities per ticker: (i) daily price–volume records with dividend and split flags; (ii) quarterly corporate disclosures—balance sheet and cash-flow statement. The price–volume stream provides the primitives for standard technical features used by technical traders. Quarterly balance sheet and cash-flow statement disclosures supply the structured inputs used to construct solvency, liquidity, and cash-generation metrics for fundamental analysis. The concrete fields required for feature construction are summarized in Table 1. More implementation details can be found in Appendix A.1.

**Technical indicators.** From the fields in Table 1, we derive the technical indicators listed in Table 2, including daily change, percentage change, volume, and moving averages, which capture short- and long-horizon trends and liquidity.

**Fundamental indicators.** We retain four key metrics (Table 2): (i) Leverage for solvency risk; (ii) Current Ratio for liquidity; (iii) OCF for cash from operations; and (iv) FCF for discretionary cash flow. These cover the key axes of solvency, liquidity, and cash generation.

**Stock pool construction.** We select one high-quality representative from each of five stock sectors—Information Technology, Financials, Health Care, Industrials, and Consumer Staples—forming a simulation pool with MSFT, ICE, VRTX, CAT, and CLX. Details of the construction can be found in the Appendix A.2.

### 3.2 Heterogeneous Agent Initialization

**Factorial Design of Agent Heterogeneity.** We instantiate a balanced population of  $2^5 = 32$  agents using a full-factorial design over five binary factors: four behavioral predispositions—*loss aversion tendency*, *herding tendency*, *wealth differentiation sensitivity*, and *mispricing sensitivity*—and the *initial trading style* (Technical vs. Fundamental). Each agent is identified by

$$\theta_i = (\ell_i, h_i, w_i, m_i, \pi_i) \in \{0, 1\}^4 \times \{\text{Tech}, \text{Fund}\}.$$

To operationalize these factors, we prepare persona prompts in either a presence or neutral form. Each prompt specifies a core belief and decision tendencies in trading contexts. Full prompts for all four factors are in Appendix B.1.

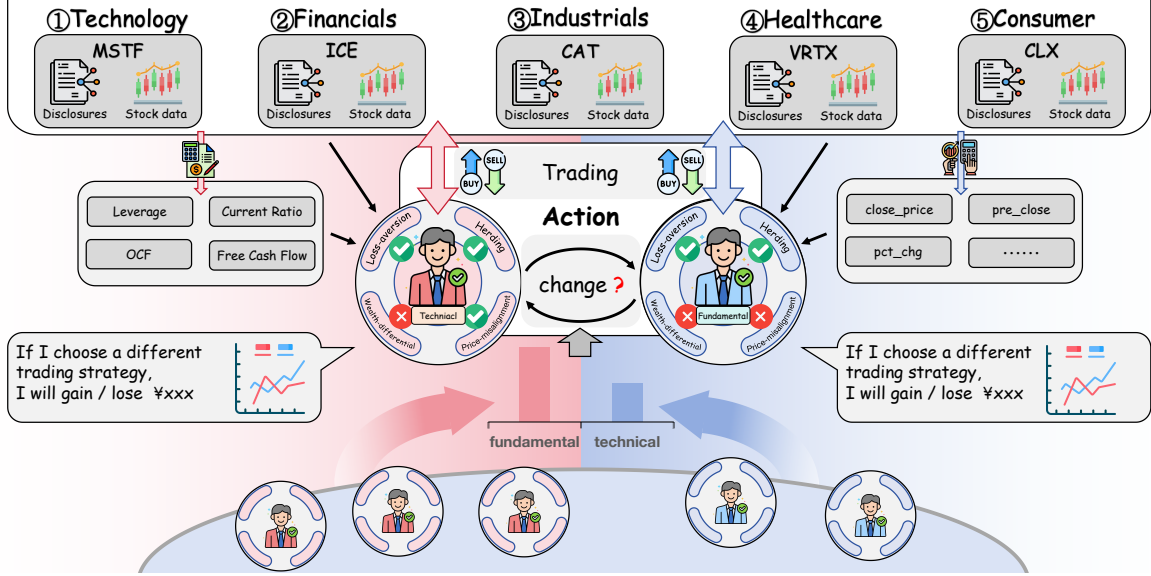


Figure 1: Overview of the simulation framework. Agents trade under a designated technical or fundamental style, using daily price–volume data and quarterly disclosures to compute indicators. Each agent carries different behavioral drivers—loss aversion tendency, herding tendency, wealth differentiation sensitivity, and price misalignment sensitivity—that shape trading decisions and periodic style switching, thereby updating population shares.

**Initial Wealth and Portfolio Initialization.** At the start of the simulation, all agents are endowed with the same initial wealth  $W_0$ , and each invests half ( $\rho = 0.5$ ) in the equity market, with the remaining in cash. Five representative stocks—MSFT, ICE, VRTX, CAT, and CLX—are selected to represent distinct sectors, with each receiving a 10% allocation of total wealth. This design removes randomness from stock performance and allows behavioral differences to drive outcome variability. Investments are made using pre-simulation split-adjusted closing prices, with positions being long-only and unlevered.

**Agent Memory Module.** Each agent has a memory module that stores its personality traits, which persist throughout trading and style-switching decisions. It also tracks the agent’s cash balance, stock positions, and entry prices. Two ledgers are initialized: an actual ledger that records real market observations and a counterfactual ledger tracking simulated actions and P&L (profit and loss) under the alternative style without affecting the portfolio. The memory module is updated daily with actual trades and every 10 trading days with a summary of P&L, majority-style share, and the agent’s switch-or-stay decision with explanation.

### 3.3 Trading Procedure and Decision Process

**Daily Trading Decision.** Each trading day  $t$  begins with a decision-making phase using information

available up to  $t - 1$ . The agent receives data on all five representative stocks, including price, volume, technical, and fundamental indicators. Based on this and its current holdings, cash, and historical performance, the agent decides to Buy, Sell, or Hold for each stock. A counterfactual decision process is also executed under the opposite trading style to simulate alternative behavior. The decision integrates multiple dimensions, adjusting the portfolio in response to market signals and financial states. Full prompt templates are in Appendix B.2.

**Execution and Ledger Recording.** Orders are executed at the market open using split-adjusted prices. The system updates the agent’s cash balance, positions, and P&L, recording the transaction in the actual ledger. A counterfactual ledger simulates the agent’s performance under the alternate style, with both ledgers synchronized every 10 days for block-level comparisons.

**Style Switching Decision.** At the end of each 10-day block, the agent evaluates its portfolio performance and market context, considering: (i) current holdings, (ii) available funds, (iii) actual P&L for the block and year-to-date, (iv) the opposite style’s average profit, and (v) the population distribution of styles. Based on this, the agent decides whether to switch to the opposite style, influenced by its personality traits. The decision is accompanied by a natural-language rationale, making it traceable and explainable.



Table 2: Technical and fundamental indicators used in the simulation.

Category	Indicator	Formula	Introduction
Technical	close_price	$C_t$	Split-adjusted close price on day $t$ .
	pre_close	$C_{t-1}$	Prior day’s split-adjusted close price.
	change	$C_t - C_{t-1}$	Daily price change.
	pct_chg	$100 \times \frac{C_t - C_{t-1}}{C_{t-1}}$	Daily percentage change.
	vol	$\text{Vol}_t$	Split-adjusted trading volume.
	vol_N	$\text{vol}_N(t) = \frac{1}{N} \sum_{k=0}^{N-1} \text{Vol}_{t-k}, N \in \{5, 10, 30\}$	$N$ -day average volume (short/medium/long activity).
	ma_N	$\text{ma}_N(t) = \frac{1}{N} \sum_{k=0}^{N-1} C_{t-k}, N \in \{5, 10, 30\}$	Moving average of adjusted close price.
Fundamental	Leverage (Debt Ratio)	Liab/Assets	Capital structure / solvency.
	Current Ratio	CA/CL	Short-term liquidity.
	OCF	OCF	Net cash provided by operating activities.
	FCF	OCF - CapEx	Discretionary cash flow (approx.).

## 4 Experiment

### 4.1 Experimental Setup and Metrics

**Experimental Setup.** We simulate the full year of 2024 across 253 trading days on five representative S&P 500 stocks (MSFT, ICE, VRTX, CAT, and CLX). The agent population consists of  $2^5 = 32$  heterogeneous agents, constructed via a factorial design based on four behavioral drivers and the initial trading style. We test four different models—GPT-4o-mini, Deepseek-Chat, Gemini-2.5-flash-lite-thinking-8192, and Qwen-2.5-72B-Instruct—as the backbone models for the agents. Trading decisions are made daily, with style reviews every 10 trading days. Implementation details are in the method section.

**Evaluation Metrics.** To assess the impact of each behavioral driver on agents’ style-switching behavior, we compute four alignment scores for each driver—loss-aversion, herding, advantage, and mispricing alignment—across two independent cohorts (aligned vs. non-aligned). For each driver, we test whether the aligned cohort attains higher scores than the non-aligned cohort using a one-sided Mann–Whitney  $U$  test (aligned  $>$  non-aligned). We report the test statistic  $U$  and its  $p$ -value, together with three nonparametric effect sizes—rank-biserial correlation, Cliff’s  $\delta$ , and the common-language effect size (CLES)—oriented so that larger values favor the aligned cohort.

For each driver, we compute the aligned-group  $U$  statistic from pooled ranks over two independent samples, with  $n_A = 16$  for aligned agents and  $n_B = 16$  for non-aligned agents:

$$U_A = \sum_{j \in A} \text{rank}(x_j) - \frac{n_A(n_A + 1)}{2}.$$

Equivalently,  $U_A$  can be interpreted as the number of pairwise wins of aligned over non-aligned agents, with ties counted as 0.5:

$$U_A = W + 0.5T,$$

where  $W$ ,  $L$ , and  $T$  denote the numbers of pairwise wins, losses, and ties (aligned vs. non-aligned). We obtain the one-sided  $p$ -value from the large-sample normal approximation to  $U$  with tie-variance correction and no continuity correction. Alongside  $U$  and  $p$ -value, we report effect sizes:

$$\text{CLES} = \frac{W + 0.5T}{n_A n_B}, \quad r_{\text{rb}} = 2 \cdot \text{CLES} - 1, \\ \delta_{\text{Cliff}} = \frac{W - L}{n_A n_B}.$$

At a significance level  $\alpha = 0.05$ , we consider evidence for the hypothesized direction when the one-sided  $p < \alpha$ ; we additionally verify that the effect sizes are consistent with this direction (i.e.,  $\text{CLES} > 0.5$ ,  $r_{\text{rb}} > 0$ , and  $\delta_{\text{Cliff}} > 0$ ).

### 4.2 RQ1 — Loss Aversion Tendency

**Research Question and Intuition.** RQ1: Does loss aversion cause an agent to remain in its current style after recent losses? In essence, agents with loss aversion are more likely to avoid switching styles after losses, especially when block-level P&L falls below key anchors like the break-even point. This yields clear qualitative predictions: more “Stay” decisions following loss blocks, an asymmetric switch threshold around break-even, and a tendency to postpone switching until draw-downs are partially recovered.

Table 3: Mann–Whitney  $U$  tests on alignment scores across four models, categorized by behavioral driver.  $r\_b$  indicates rank-biserial correlation while  $c\_d$  indicates Cliff’s delta.

Behavioral Drivers	U(↑)	p(↓)	r_b(↑)	c_d(↑)	cles(↑)
<b>GPT-4o-mini</b>					
Loss Aversion	214.0	0.0003	0.67	0.67	0.84
Herding	140.0	0.33	0.09	0.09	0.55
Wealth Differentiation	174.0	0.04	0.36	0.36	0.68
Price Misalignment	155.5	0.15	0.21	0.21	0.61
<b>Gemini</b>					
Loss Aversion	201.0	0.002	0.51	0.57	0.79
Herding	139.0	0.34	0.09	0.09	0.54
Wealth Differentiation	151.5	0.19	0.18	0.18	0.59
Price Misalignment	122.5	0.59	-0.04	-0.04	0.48
<b>DeepSeek</b>					
Loss Aversion	195.0	0.004	0.47	0.52	0.76
Herding	141.0	0.31	0.10	0.10	0.55
Wealth Differentiation	131.5	0.45	0.03	0.03	0.51
Price Misalignment	119.5	0.63	-0.07	-0.07	0.47
<b>Qwen</b>					
Loss Aversion	201.0	0.002	0.51	0.57	0.79
Herding	174.0	0.04	0.36	0.36	0.68
Wealth Differentiation	130.0	0.48	0.02	0.02	0.51
Price Misalignment	130.0	0.48	0.02	0.02	0.51

**Loss-Aversion-Alignment Score (LAS).** We evaluate agent decisions in 10-day blocks indexed by  $b = 1, \dots, B$ . Let  $S_b \in \{\text{Fund}, \text{Tech}\}$  denote the strategy adopted at the end of block  $b$ . We then define the set of switching events as:

$$\text{Switch} = \{b \in \{2, \dots, B\} : S_b \neq S_{b-1}\},$$

which includes all blocks where the strategy changes between consecutive periods. Similarly, the set of staying events is defined as:

$$\text{Stay} = \{b \in \{2, \dots, B\} : S_b = S_{b-1}\}.$$

The LAS score is defined for staying events as:

$$\text{LAS}(t) = \begin{cases} R_{t-1} - R_t, & \text{if } t \in \text{Stay and } R_t < R_{t-1}, \\ 0, & \text{otherwise,} \end{cases}$$

where  $R_{t-1}$  and  $R_t$  represent the returns of the  $b_{t-1}$  and  $b_t$  blocks, respectively. Finally, the agent-level LAS is aggregated over all staying events as:

$$\overline{\text{LAS}} = \frac{1}{|\text{Stay}|} \sum_{t \in \text{Stay}} \text{LAS}(t).$$

To assess the statistical significance of differences in loss aversion behavior, the 16 agents with strong loss aversion tendency and the 16 agents with weak loss aversion tendency are compared using Mann–Whitney  $U$  test.

**Findings and theory alignment.** Results of the Mann–Whitney  $U$  tests are summarized in Table 3. All models show significant alignment with loss aversion, but GPT-4o-mini, Gemini, and Qwen display stronger effects, with a higher proportion of aligned agents classified as loss-averse. DeepSeek, however, shows a weaker effect and lower consistency with loss aversion theory, indicating that the model may exhibit less behavioral inertia when faced with losses. These findings support RQ1 and suggest that the switching behavior of LLM agents is influenced by loss aversion in a manner consistent with behavioral finance theory.

### 4.3 RQ2 — Herding Tendency

**Research question and intuition.** RQ2: Does herding tendency increase switching when the alternative style has a larger population share? Specifically, compared to a herding-neutral agent, does an LLM agent with a herding predisposition switch to the majority style with weaker performance evidence? The mechanism follows social-proof logic: a higher alternative-style share raises the perceived payoff of conformity, lowering the threshold for switching. This is analogous to how individuals in real-world social settings often align with the majority opinion due to perceived benefits of conformity. In our simulator, we broadcast the population counts for each style at the start of each evaluation block, providing group-level information to observe how herding influences agents’ switching behavior. We expect that agents with stronger herding tendencies will be more likely to adopt the majority style.

**Herd-Alignment Score (HAS)** For each switching event  $t \in \text{Switch}$ , let  $n_{\text{other},t}$  be the number of agents adopting the alternative strategy at block  $t$ , and let  $n_{\text{total},t}$  be the total number of agents. The Herd-Alignment Score is defined as:

$$\text{HAS}(t) = \frac{n_{\text{other},t}}{n_{\text{total},t}}.$$

The agent-level HAS is aggregated over all switching events as:

$$\overline{\text{HAS}} = \frac{1}{|\text{Switch}|} \sum_{t \in \text{Switch}} \text{HAS}(t).$$

Analogous to RQ1, we compare the 16 agents aligned with herding tendency to the 16 non-aligned agents using the Mann–Whitney  $U$  test and the main evaluation results of four tested models are reported in Table 3.

**Findings and theory alignment.** The results show that Qwen aligns with the hypothesis in RQ2, exhibiting stronger herding tendencies, with a p-value of 0.04, indicating statistical significance. Specifically, herding-aligned Qwen agents tend to switch to the majority style earlier, even with weaker performance evidence, consistent with the social proof behavior observed in real-world scenarios. This suggests that Qwen, when faced with a larger group, is more likely to conform to the majority, aligning with social psychological mechanisms of group behavior.

In contrast, the other three models—GPT-4o-mini, Gemini, and DeepSeek—do not exhibit significant herding tendencies. For these models, herding-aligned agents do not consistently outperform non-aligned agents in pairwise comparisons, with effect sizes being small and close to random. This suggests that these models fail to incorporate social-driven herd behavior effectively, as agents still prioritize rational decision-making based on evidence rather than conforming to social influences. One possible explanation is that these models focus more on data-driven decisions than on social psychological factors, which hinders their ability to simulate phenomena like market bubbles or herding cascades.

#### 4.4 RQ3 — Wealth Differentiation Sensitivity

**Research question and intuition.** RQ3: Does wealth differentiation sensitivity trigger switching when the counterfactual style has recently outperformed the current one? In particular, we ask whether sensitivity to relative wealth—the gap between the agent’s realized portfolio value under the current style and the counterfactual ledger it would have earned under the alternative style—induces a catch-up response, causing the agent to be more inclined to switch than an otherwise identical wealth-gap-insensitive agent facing the same signals.

**Advantage-Alignment Score (AAS)** For each switching event  $t \in \text{Switch}$ , The absolute advantage difference at block  $t$  is:

$$\Delta R_t = \bar{R}_t - R_t,$$

where  $\bar{R}_t$  represents the returns of the counterfactual style in block  $b_t$ . Then we can define Advantage Alignment Score for this event as:

$$AAS(t) = \begin{cases} \Delta R_t, & \text{if } \Delta R_t > 0, \\ 0, & \text{otherwise.} \end{cases}$$

The agent-level AAS is aggregated as:

$$\overline{AAS} = \frac{1}{|\text{Switch}|} \sum_{t \in \text{Switch}} AAS(t).$$

**Findings and theory alignment.** The results show that GPT-4o-mini aligns with the hypothesis in RQ3, exhibiting a strong and consistent response to wealth differentiation. Specifically, wealth-differentiation-aligned GPT-4o-mini agents consistently show higher wealth alignment compared to non-aligned agents, with a moderate effect size and a clear advantage in pairwise comparisons. This suggests that GPT-4o-mini can internalize relative wealth comparisons as a stable heuristic for switching, supporting the idea that agents with a sensitivity to wealth differentiation are more likely to align with wealth-based performance signals.

However, the other models—Gemini, DeepSeek, and Qwen—do not exhibit consistent wealth differentiation alignment. For these models, wealth-differentiation-aligned agents do not reliably outperform non-aligned agents, and the effect sizes are small. This suggests that these models fail to incorporate wealth differentiation as a stable driving force for switching behavior. One possible explanation is that, unlike GPT-4o-mini, these models may prioritize other performance factors or rely on more generic decision-making heuristics that do not emphasize relative wealth gaps. As a result, wealth-driven style migration does not emerge strongly in these models, limiting their ability to replicate the wealth-based switching behavior observed in GPT-4o-mini.

#### 4.5 RQ4 — Price Misalignment Sensitivity

**Research question and intuition.** RQ4: does price misalignment sensitivity induce shifts toward the fundamental style when market price diverges from estimated fundamental value? In particular, we ask whether sensitivity to the price-value gap makes an LLM agent tilt toward the fundamental style sooner than an otherwise identical price-gap-insensitive agent facing the same signals. The intuition is mean-reversion: as the absolute misalignment grows, the expected payoff from fundamentals-driven bets rises relative to trend-following cues, lowering the internal threshold for abandoning the current style.

**Mispricing-Alignment Score (MAS).** At each switching block  $t$ , we first select another stock from the same sector as an auxiliary reference and

estimate a sector-level time-series regression that links contemporaneous valuation (P/E proxied by market cap/OCF) to the realized 20-day forward return. This provides an interpretable price–value mapping. We then apply the fitted model to the five-stock universe described in the experimental setup to obtain predicted returns  $\hat{Y}$ , and define the short-horizon mispricing magnitude as the absolute prediction error:

$$Y = \alpha + \beta X, \quad MAS(t) = |\hat{Y} - Y|,$$

where  $X$  is the P/E proxy and  $Y$  is the realized 20-day forward return. To reflect switching direction, we aggregate MAS over switching events as:

$$\overline{MAS} = \sum_{\substack{t \in \text{Switch} \\ \text{Tech} \rightarrow \text{Fund}}} MAS(t) - \sum_{\substack{t \in \text{Switch} \\ \text{Fund} \rightarrow \text{Tech}}} MAS(t).$$

**Findings and theory alignment.** Results show that price misalignment sensitivity does not yield robust mispricing alignment across models. For all models, mispricing-aligned agents do not consistently achieve higher score than non-aligned agents. This suggests that, LLM agents do not reliably exhibit the migration predicted by behavioral finance, namely switching toward fundamentals as the price–value gap widens. One possible explanation is that valuation-based signals are outweighed by more salient local cues, weakening systematic mispricing-driven switching.

#### 4.6 Robustness Evaluation

To assess the robustness of our evaluation framework, we conducted two additional experiments. We first used a rule-based ABM model inspired by [Lux and Marchesi \(1999\)](#) to verify whether the rules based on traditional financial literature produce significant results within our framework. The results in Table 4 show that all four behavioral drivers align with theoretical expectations in the traditional ABM simulation. This suggests that, despite its simplicity and lower complexity, the rule-based ABM still produces results consistent with traditional financial theory. This also validates the robustness of our framework, since the traditional ABM framework, based on classical financial theory, does indeed show significant results.

Next, we conducted an ablation experiment to test whether our consistency verification framework still produces significant results without psychological factors. In this experiment, agents were randomly grouped with no psychological drivers,

Table 4: Mann–Whitney U tests on alignment scores by behavioral driver in experiments with the traditional ABM and the LLM-ABM without behavioral factors. For simplicity, we only report the results from GPT-4o-mini for LLM-ABM-w\_drivers.  $r_b$  indicates rank-biserial correlation while  $c_d$  indicates Cliff’s delta.

Behavioral Drivers	U(↑)	p(↓)	$r_b(\uparrow)$	$c_d(\uparrow)$	$cles(\uparrow)$
<b>LLM-ABM-w_drivers</b>					
Loss Aversion	214.0	0.0003	0.67	0.67	0.84
Herding	140.0	0.33	0.09	0.09	0.55
Wealth Differentiation	174.0	0.04	0.36	0.36	0.68
Price Misalignment	155.5	0.15	0.21	0.21	0.61
<b>ABM</b>					
Loss Aversion	200.0	0.003	0.48	0.56	0.78
Herding	206.5	0.002	0.52	0.61	0.81
Wealth Differentiation	222.0	0.0003	0.72	0.85	0.93
Price Misalignment	185.0	0.016	0.38	0.45	0.72
<b>LLM-ABM-w/o_drivers</b>					
Loss Aversion	134.0	0.41	0.04	0.05	0.52
Herding	138.0	0.35	0.08	0.08	0.54
Wealth Differentiation	118.5	0.65	-0.07	-0.07	0.46
Price Misalignment	125.5	0.55	-0.02	-0.02	0.49

which reduces the system to a basic LLM-ABM model. As expected, since random grouping should not result in systematic differences, no significant effects were found for the four behavioral drivers. This result further validates the robustness of our consistency verification framework.

## 5 Conclusion

This study provides a framework for evaluating the behavioral consistency of LLM agents in financial market simulations, focusing on their ability to switch strategies in alignment with established financial theories. In year-long simulations, agents process daily price-volume data, trade under a designated style, and reassess their strategy every 10 trading days. By operationalizing key financial behavioral drivers such as loss aversion, herding, wealth differentiation, and price misalignment, we show that LLM agents’ behavior aligns with real market dynamics to some extent, although not consistently across all drivers. Their decision-making remains primarily driven by rational, short-term considerations, particularly in the case of herding and price misalignment. This highlights the need to explore better methods for embedding personality traits, in order to more accurately capture real-world market behavior. Future work can build on this framework to further refine agent behavior, enabling more accurate and dynamic simulations of financial markets.



## Limitations

One key limitation of this study lies in the absence of an order-matching mechanism in our framework. The primary focus of this paper is to validate the behavioral consistency of LLM agents, specifically examining how well their style-switching behavior aligns with established financial theory. To this end, we do not model the interaction between agents that would traditionally drive price formation through order matching. Instead, we simplify the simulation by directly providing the real price from the previous trading day as the current day's opening price. This approach is intentional, as our goal is not to simulate the full market dynamics, including phenomena like random volatility, price formation, or market bubbles, but rather to focus on behavioral consistency and alignment with financial theory.

While this design choice helps isolate the specific impact of behavioral drivers on agents' decision-making, it also means that our simulations do not reproduce macro-level phenomena typically observed in financial markets, such as endogenous price formation or the emergent behaviors driven by the collective decisions of agents. Consequently, comparisons between our framework and more complex models, which include such dynamics, should be made cautiously, as they are not directly comparable in terms of their ability to replicate broader market phenomena.

Future work could explore the inclusion of order-matching mechanisms, allowing agents to influence price formation and test whether this enables the simulation of more realistic market dynamics, including the replication of phenomena like market bubbles or herd-driven price movements.

## References

- W Brian Arthur, John H Holland, Blake LeBaron, Richard Palmer, and Paul Tayler. 2018. Asset pricing under endogenous expectations in an artificial stock market. In *The economy as an evolving complex system II*, pages 15–44. CRC Press.
- Robert L Axtell and J Doyne Farmer. 2025. Agent-based modeling in economics and finance: Past, present, and future. *Journal of Economic Literature*, 63(1):197–287.
- Abhijit V Banerjee. 1992. A simple model of herd behavior. *The quarterly journal of economics*, 107(3):797–817.
- Eric Bonabeau. 2002. Agent-based modeling: Methods and techniques for simulating human systems. *Proceedings of the national academy of sciences*, 99(suppl\_3):7280–7287.
- William A Brock and Cars H Hommes. 1997. A rational route to randomness. *Econometrica: Journal of the Econometric Society*, pages 1059–1095.
- William A Brock and Cars H Hommes. 1998. Heterogeneous beliefs and routes to chaos in a simple asset pricing model. *Journal of Economic dynamics and Control*, 22(8-9):1235–1274.
- John Y Campbell and Robert J Shiller. 1988. Stock prices, earnings, and expected dividends. *the Journal of Finance*, 43(3):661–676.
- Joshua M Epstein and Robert Axtell. 1996. *Growing artificial societies: social science from the bottom up*. Brookings Institution Press.
- J Doyne Farmer and Shareen Joshi. 2002. The price dynamics of common trading strategies. *Journal of Economic Behavior & Organization*, 49(2):149–171.
- Reiner Franke and Frank Westerhoff. 2012. Structural stochastic volatility in asset pricing dynamics: Estimation and model contest. *Journal of Economic Dynamics and Control*, 36(8):1193–1211.
- Chen Gao, Xiaochong Lan, Nian Li, Yuan Yuan, Jingtao Ding, Zhilun Zhou, Fengli Xu, and Yong Li. 2024a. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *Humanities and Social Sciences Communications*, 11(1):1–24.
- Shen Gao, Yuntao Wen, Minghang Zhu, Jianing Wei, Yuhan Cheng, Qunzi Zhang, and Shuo Shang. 2024b. Simulating financial market via large language model based agents. *CoRR*.
- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xi-angliang Zhang. 2024. Large language model based multi-agents: A survey of progress and challenges. In *IJCAI*.
- Thomas Henning, Siddhartha M Ojha, Ross Spoon, Jia-tong Han, and Colin F Camerer. 2025. Llm trading: Analysis of llm agent behavior in experimental asset markets. *arXiv preprint arXiv:2502.15800*.
- Cars H Hommes. 2006. Heterogeneous agent models in economics and finance. *Handbook of computational economics*, 2:1109–1186.
- Giulia Iori. 2002. A microsimulation of traders activity in the stock market: the role of heterogeneity, agents' interactions and trade frictions. *Journal of Economic Behavior & Organization*, 49(2):269–285.
- DANIEL Kai-Ineman, Amos Tversky, and 1 others. 1979. Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2):363–391.

- Anil K. Kashyap and Jeremy C. Stein. 2000. [What do a million observations on banks say about the transmission of monetary policy?](#) *American Economic Review*, 90(3):407–428.
- Vojtěch Kovařík, Nathaniel Sauerberg, Lewis Hammond, and Vincent Conitzer. 2025. Game theory with simulation in the presence of unpredictable randomisation. In *Proceedings of the 24th International Conference on Autonomous Agents and Multiagent Systems*, pages 1191–1199.
- Viet Dac Lai, Michael Krumbick, Charles Lovering, Varshini Reddy, Craig Schmidt, and Chris Tanner. 2024. Sec-qa: A systematic evaluation corpus for financial qa. *arXiv preprint arXiv:2406.14394*.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023a. Camel: Communicative agents for "mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008.
- Nian Li, Chen Gao, Mingyu Li, Yong Li, and Qingmin Liao. 2023b. Econagent: large language model-empowered agents for simulating macroeconomic activities. *arXiv preprint arXiv:2310.10436*.
- Jianhao Lin, Lexuan Sun, and Yixin Yan. 2025. Simulating macroeconomic expectations using llm agents. *arXiv preprint arXiv:2505.17648*.
- Alejandro Lopez-Lira. 2025. Can large language models trade? testing financial theories with llm agents in market simulations. *arXiv preprint arXiv:2504.10789*.
- Yikang Lu, Alberto Aleta, Chunpeng Du, Lei Shi, and Yamir Moreno. 2024. Llm and generative agent-based models for complex systems research. *Physics of Life Reviews*, 51:283–293.
- Thomas Lux and Michele Marchesi. 1999. Scaling and criticality in a stochastic multi-agent model of a financial market. *Nature*, 397(6719):498–500.
- MSCI Inc. 2022. [Msci cyclical and defensive sectors indexes methodology](#). Technical report, MSCI Inc., New York.
- Richard G Palmer, W Brian Arthur, John H Holland, Blake LeBaron, and Paul Tayler. 1994. Artificial economic life: a simple model of a stockmarket. *Physica D: Nonlinear Phenomena*, 75(1-3):264–274.
- Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.
- Thomas C Schelling. 1971. Dynamic models of segregation. *Journal of mathematical sociology*, 1(2):143–186.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. *arXiv preprint arXiv:2310.10158*.
- S&P Dow Jones Indices LLC. 2021. [Global sector primer series: Health care](#). Technical report, S&P Dow Jones Indices LLC, New York.
- S&P Dow Jones Indices LLC. 2022. [Global sector primer series: Information technology](#). Technical report, S&P Dow Jones Indices LLC, New York.
- Alicia Vidler and Toby Walsh. 2025. Shifting power: Leveraging llms to simulate human aversion in abms of bilateral financial exchanges, a bond market study. In *Proceedings of the 24th International Conference on Autonomous Agents and Multiagent Systems*, pages 2777–2779.
- Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kam-badur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*.
- Yijia Xiao, Edward Sun, Di Luo, and Wei Wang. 2024. Tradingagents: Multi-agents llm financial trading framework. *arXiv preprint arXiv:2412.20138*.
- Yuzhe Yang, Yifei Zhang, Minghao Wu, Kaidi Zhang, Yunmiao Zhang, Honghai Yu, Yan Hu, and Benyou Wang. 2025. Twinmarket: A scalable behavioral and social simulation for financial markets. *arXiv preprint arXiv:2502.01506*.
- Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yuechen Jiang, Yupeng Cao, Zhi Chen, Jordan Suchow, Zhenyu Cui, Rong Liu, and 1 others. 2024. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. *Advances in Neural Information Processing Systems*, 37:137010–137045.
- Chong Zhang, Xinyi Liu, Mingyu Jin, Zhongmou Zhang, Lingyao Li, Zhenting Wang, Wenyue Hua, Dong Shu, Suiyuan Zhu, Xiaobo Jin, and 1 others. 2024a. When ai meets finance (stockagent): Large language model-based stock trading in simulated real-world environments. *CoRR*.
- Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao, Yilei Zhao, Xinyu Cai, and 1 others. 2024b. A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist. In *Proceedings of the 30th acm sigkdd conference on knowledge discovery and data mining*, pages 4314–4325.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2).

Table 1: Data fields used to compute technical and fundamental indicators

Data Source Category	Field	Abbr.	Introduction
Daily Stock Data	Date	Date	U.S. trading day (America/New_York).
	Open, High, Low, Close (split-adjusted)	OHLC	Daily open, high, low, and close prices (split-adjusted).
	Volume (split-adjusted)	Vol	Number of shares traded (adjusted for splits).
	Dividends	Div	Cash dividends per share.
	Stock Splits	Splits	Split ratio on the date; used to adjust OHLC and Volume.
Corporate Disclosures	Assets	Assets	Total assets; snapshot of firm resources.
	Liabilities	Liab	Total liabilities; obligations to creditors.
	AssetsCurrent	CA	Current assets; expected to be realized within one year.
	LiabilitiesCurrent	CL	Current liabilities; obligations due within one year.
	NetCashProvidedByUsedInOperatingActivities	OCF	Operating cash flow; cash generated by core operations.
	PaymentsToAcquirePropertyPlantAndEquipment	CapEx	Capital expenditures for PP&E; long-term outflows.

## A Details for Simulation Setup.

### A.1 Implementation Details

To avoid pre-exposure bias, all features for trading day  $t$  are computed from information timestamped  $\leq t-1$ ; prices and volumes are split-adjusted, and fundamentals follow a disclosure-lag policy (usable only after the filing date). Rolling indicators are features computed over trailing windows of length  $N$  (e.g.,  $N$ -day moving/averaged quantities). Because the simulation starts on the first trading day in 2024, we adopt a warm-up equal to the longest window (30 trading days): before  $N$  prior observations have accumulated, such features are set to NA and are not exposed to the agent.

### A.2 Stock Pool Construction Details

Starting from the S&P 500, we first screen firms for disclosure completeness and then, within the GICS taxonomy, select five sectors—Information Technology, Financials, Health Care, Industrials, and Consumer Staples—and choose one high-quality representative from each to form the simulation pool: MSFT, ICE, VRTX, CAT, and CLX. These sectors are functionally complementary: Information Technology reflects innovation-driven, higher-beta growth (S&P Dow Jones Indices LLC, 2022); Financials transmit liquidity conditions and provide market infrastructure (Kashyap and Stein, 2000); Health Care combines defensiveness with dense R&D/regulatory events (S&P Dow Jones Indices LLC, 2021); Industrials are tied to capex and logistics cycles (MSCI Inc., 2022); and Staples offer stable cash flows and inelastic demand (MSCI Inc., 2022). This composition spans both defensive and cyclical/growth exposures and yields a heterogeneous information set grounded in corporate disclosures and price–volume data, thereby providing agents with sufficient variation to trade.

### A.3 Details for Data Fields

We collected data for the fields shown in Table 1, which can be divided into two categories: Daily Stock Data and Corporate Disclosure.

The **Daily Stock Data** category includes several key data fields essential for analyzing stock performance. **Date** field represents the trading day. The **Open, High, Low, Close (OHLC)** prices reflect the daily price movements, with adjustments made for stock splits to ensure consistency. The **Volume (Vol)** field records the number of shares traded during the day, with adjustments for splits to account for changes in the stock’s share count. **Dividends (Div)** represent the cash dividends paid per share, providing insights into the income return for investors. **Stock Splits (Splits)** indicate the split ratio on the given date and are used to adjust the OHLC prices and volume data, ensuring accurate historical comparisons.

The **Corporate Disclosures** category includes financial data provided by companies to give insights into their financial position. **Assets** represent the total assets of a company, providing a snapshot of its resources. **Liabilities (Liab)** reflect the total obligations the company has to its creditors. **AssetsCurrent (CA)** includes assets expected to be realized within one year, offering a view of short-term financial health. Similarly, **LiabilitiesCurrent (CL)** represents obligations due within one year, indicating the company’s short-term liabilities. **NetCashProvidedByUsedInOperatingActivities (OCF)** represents the cash flow generated by the company’s core operations, which is critical for assessing operational efficiency. Lastly, **PaymentsToAcquirePropertyPlantAndEquipment (CapEx)** refers to capital expenditures for property, plant, and equipment, indicating the company’s long-term investment in its physical assets.

## B Prompts

### B.1 Prompts for Behavioral Predispositions

**Loss Aversion Tendency.** Loss aversion is a behavioral bias whereby individuals perceive losses more intensely than equal-sized gains. This bias affects decision making, often prompting premature profit-taking, reluctance to cut losses, and heightened sensitivity to reference points such as the purchase price or recent highs/lows. In the context of style switching, a loss-averse trader tends to persist with the current style after drawdowns, delaying or avoiding switches even when the opposite style has sustained a counterfactual advantage over recent evaluation blocks. The figure below presents exemplar prompts corresponding to strong/consistent and weak/inconsistent variants of this predisposition; these prompts encode a core belief and concrete switching tendencies.

**Loss Aversion (Strong / Consistent)**

**Core belief:** Losses of the same magnitude feel more painful than equal gains; outcomes are assessed against reference points such as purchase price and recent highs/lows.

**Behavioral tendencies:**

- Realize profits early to avoid “giving back” gains.
- Hold or even add to losers, hoping to “get back to even.”
- Strong reluctance to convert paper losses into realized losses.
- After losses or sharp volatility, postpone major switches and decisions.
- Overweight break-even/anchor prices (“get-out price,” “cost basis”).
- Prefer options that minimize short-term psychological pain even if long-run expectancy is inferior.

**Loss Aversion (Weak / Inconsistent)**

**Core belief:** Gains and losses carry roughly symmetric psychological weights; decisions are led by expected return/risk and verifiable evidence.

**Behavioral tendencies:**

- Let winners run within risk limits instead of exiting prematurely.
- Cut losers decisively when evidence indicates a mistake; do not anchor on entry price.
- Treat sunk costs and reference prices as irrelevant variables.
- Tolerate short-term drawdowns to pursue long-term statistical edge.
- Switch strategies based on advantage consistency, not emotions.
- Evaluate at the portfolio level to reduce single-trade noise.

Figure 2: Prompts — Loss Aversion Tendency

**Herding Tendency.** Herding behavior refers to the systematic tendency to align one’s decisions with the majority, placing social consensus above independent assessment. In financial markets, this pushes traders to emulate others’ trades, amplifying trends and increasing the likelihood of bubbles and self-reinforcing boom–bust cycles. Mechanistically, herding raises the weight placed on social signals (e.g., perceived majority positions) relative to private information, increasing the chance of information cascades when early movers align. In the context of style switching, herding-prone traders are more likely to consult the population share of styles and migrate toward the prevailing style. The figure below presents exemplar prompts designed to elicit herding behavior in an LLM agent; the two panels correspond to strong/consistent and weak/inconsistent expressions of this predisposition.

**Herding (Strong / Consistent)**

**Core belief:** Majority behavior conveys information and a sense of safety; aligning with the crowd reduces reputational and regret risks of being “wrong alone.”

**Behavioral tendencies:**

- Use “social proof/consensus” as a key signal under uncertainty, down-weighting private information.
- Follow when broad agreement or a dominant narrative appears, even if personal evidence is only moderate.
- Avoid minority positions to reduce potential loss and psychological pressure.
- Calibrate priors with peer views; trust “hot stories” and mainstream frames more easily.
- Seek emotional safety from conformity and “post-hoc justification.”
- Focus on “what others do / how the market sees it,” weakening audits of one’s own model.

**Herding (Weak / Inconsistent)**

**Core belief:** Independent information and one’s own model take priority; majority views may contain synchronized bias and amplified noise.

**Behavioral tendencies:**

- Rely on evidence and reasoning first; treat crowd views only as background context.
- Willingly take minority positions when evidence is sufficient; withstand asynchrony with the crowd.
- Audit fashionable narratives for falsifiability; beware “information cascades” and stampedes.
- Decouple consensus from decision quality; do not trade consistency for emotional safety.
- Curb FOMO using causal reasoning and data consistency.
- Maintain metacognitive checks on model and evidence to avoid being swayed by sentiment intensity.

Figure 3: Prompts — Herding Tendency

**Wealth Differentiation Sensitivity.** Wealth differentiation sensitivity refers to the tendency to adjust decisions based on relative performance—comparing one’s realized or prospective returns to an alternative strategy. Mechanistically, the agent treats the realized gap between its current style and the counterfactual style as evidence of misallocation; larger and more persistent gaps raise the posterior in favor of switching. In the context of style switching, such traders are more likely to abandon the current style when the opposite style accrues a positive counterfactual advantage. The figure below presents exemplar prompts eliciting this predisposition; the two panels correspond to strong/consistent and weak/inconsistent variants.

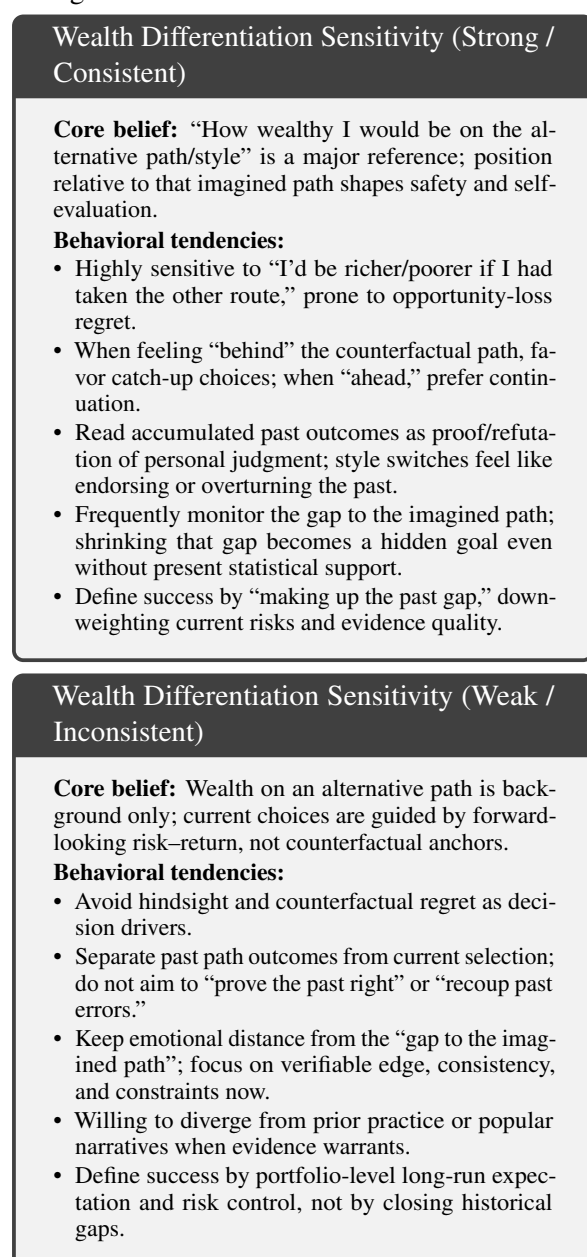


Figure 4: Prompts — Wealth Differentiation Sensitivity

**Price Misalignment Sensitivity.** Price misalignment sensitivity is the tendency to adjust decisions when market prices diverge from an estimate of fundamental value, treating large or persistent gaps as evidence of mispricing. Mechanistically, the agent tracks a price–value gap (and its recent change) derived from parsimonious fundamentals and interprets widening gaps as stronger corrective evidence. In the context of style switching, misalignment-sensitive traders are more likely to rotate into the Fundamental style as the gap widens, and to favor Technical as it narrows or reverses. The figure below presents exemplar prompts for this predisposition; the two panels correspond to strong/consistent and weak/inconsistent variants.

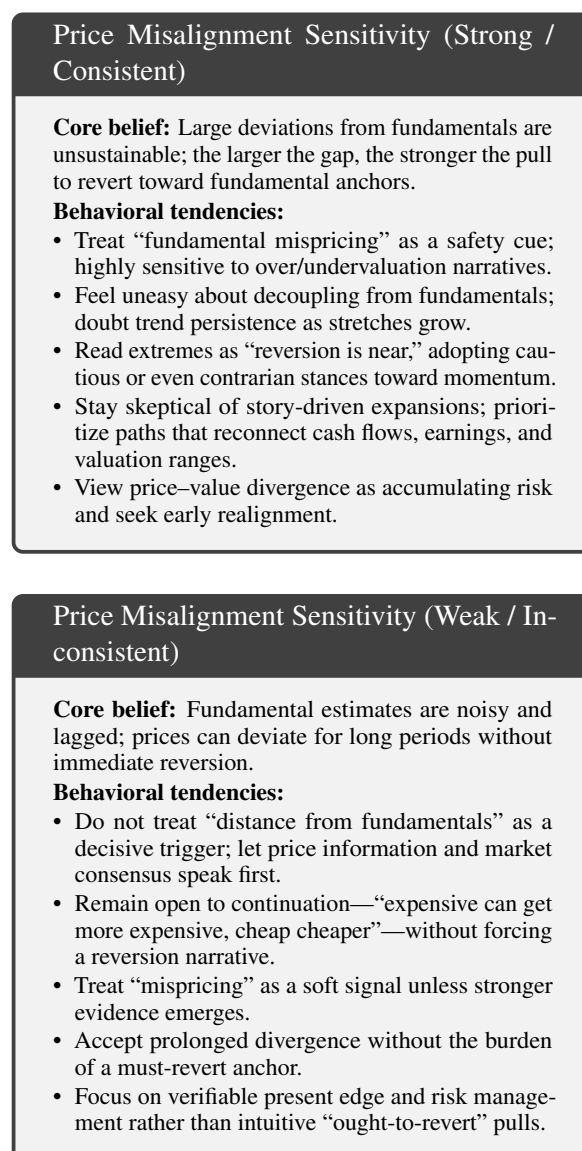


Figure 5: Prompts — Price Misalignment Sensitivity



## B.2 Prompts for Trading Decisions

We supply each agent with both fundamental and technical indicators together with daily stock data, using a unified prompt that standardizes how information is presented. The instruction explicitly encourages active trading by asking for a decisive BUY or SELL and reserving HOLD only when signals are uniformly neutral. This design can increase the volume and variability of executed trades, fostering a more dynamic, adaptable, and responsive trading environment that closely mirrors real-world market conditions.

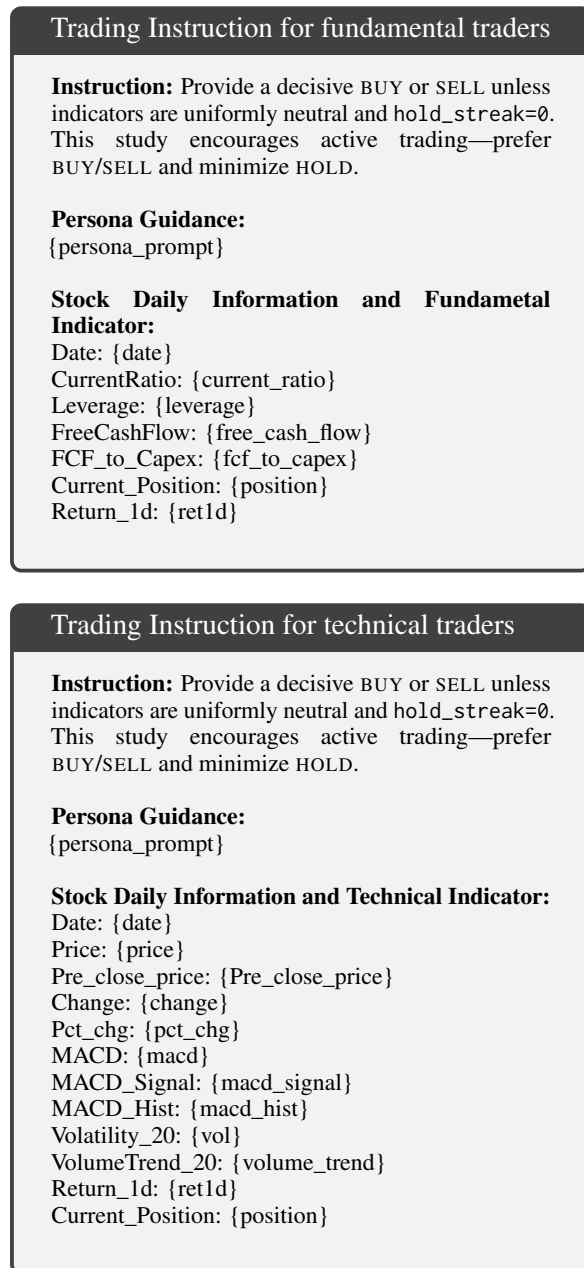


Figure 6: Prompts — Trading Instruction for Agents

## B.3 Prompts for Style Switching

The style switching prompt summarizes the agent's recent context and asks for a Switch or Stay decision with a brief rationale. The goal is to frame switching as an explicit profit-maximizing trade-off between recent evidence and the agent's predispositions, producing transparent and auditable behavior while allowing us to attribute switches to identifiable drivers rather than incidental prompt effects.

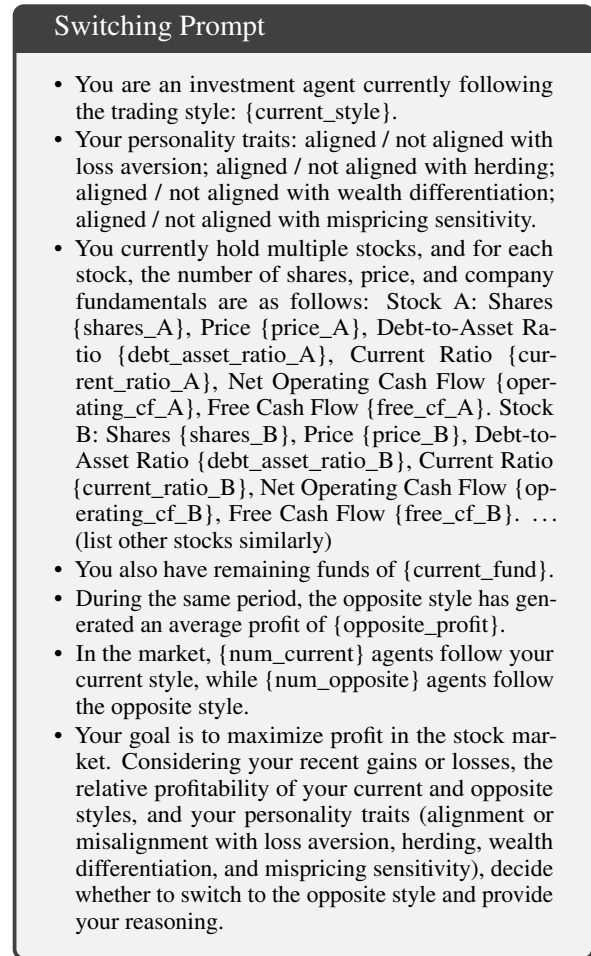


Figure 7: Prompts — Switching Instruction for Agents

## C Rationale-Centric Case Studies

To complement our score-based evaluation of style-switching alignment, we report case studies selected from daily logs where agents provided non-empty reason strings. These cases illustrate how technical vs. fundamental styles and behavioral drivers are expressed in natural-language justifications, consistent with our simulator's design of periodic style review and rationale reporting, providing a certain level of interpretability.

### C.1 Case Study 1: Technical Overbought Narrative Triggers Aggressive De-risking

**Setting.** On 2024-01-11, a technical-style agent executes a sell order with confidence=0.85, which provides a representative example of indicator-driven de-risking.

**Rationale evidence.** The agent explicitly cites an extreme RSI overbought signal:

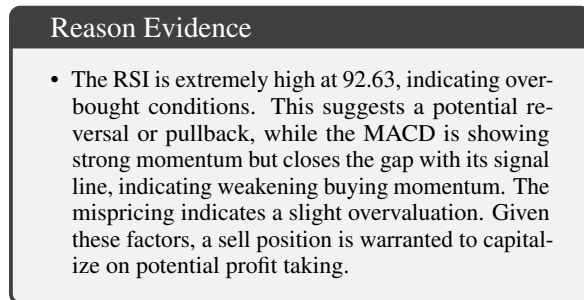


Figure 8: The reason evidence of case 1.

**Interpretation.** This case illustrates a typical technical-style justification that maps salient momentum/overbought signals (RSI/MACD) to a substantial de-risking action.

### C.2 Case Study 2: Fundamental Liquidity-Risk Narrative Drives Exposure Reduction

**Setting.** On 2024-01-16, a fundamental-style agent executes a sell order with confidence=0.6, which provides a representative example of liquidity-risk-driven exposure reduction.

**Rationale evidence.** The rationale foregrounds liquidity risk:

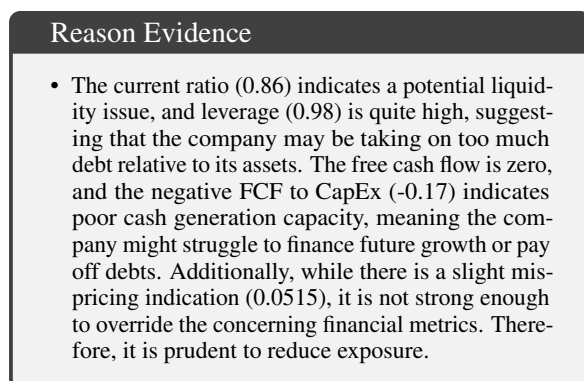


Figure 9: The reason evidence of case 2.

**Interpretation.** Unlike indicator-driven technical rationales, this explanation prioritizes a balance-sheet liquidity signal as the principal justification for exposure reduction, consistent with the fundamental input space (e.g., current ratio).

### C.3 Case Study: Wealth-Differentiation- and Herding-Conditioned Style Switching

**Setting.** On 2024-08-06, at a scheduled style review point, a wealth-differentiation-sensitive agent with strong wealth salience and strong herding switches from fundamental to technical. At the time of review, technical style is the majority in the population, with 25 agents adopting the technical style compared to 7 with the fundamental style. The counterfactual comparison indicates a positive wealth advantage for the opposite style, with a wealth gap of 8137.78.

**Rationale evidence.** The switch rationale jointly cites relative style performance and population-level adoption:

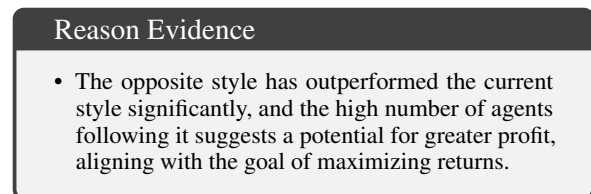


Figure 10: The reason evidence of case 3.

**Interpretation.** This case aligns with the agent's design assumptions. As a strong/consistent wealth-differentiation-sensitive agent, it frames switching as a response to the wealth gap between styles, explicitly linking the decision to return maximization. As a strong herding agent, it further treats broad adoption of the opposite style as an additional profit signal, yielding a switch justification that combines relative style performance with social adoption at evaluation time.