

# **Bias, Misleading Behavior, and Consistency of LLMs in Uncertain Decision-Making: Evidence from Financial Services Industry and Its Implications**

## **Abstract**

This paper presents a comprehensive evaluation of large language models (LLMs) in financial decision-making under uncertainty. Leveraging a synthetic investment decision dataset (SIDD) and a novel two-stage prompting framework, we simulate distinct advisory roles to quantify how subtle changes in framing systematically skew investment recommendations. Our analysis reveals four critical vulnerabilities: (1) the demonstrated biases in LLM recommendations, influenced merely by role framing, underscore a substantial regulatory and ethical challenge; (2) subtle prompt modifications can systematically induce misleading behavior, posing risks of significant legal repercussions; (3) the inconsistency observed in the logical coherence of LLM outputs under uncertain conditions poses a reliability challenge; and (4) different models have markedly different inherent biases in risk appetite. These findings illuminate urgent requirements for robust regulatory standards, rigorous prompt governance, and architectural safeguards—such as monotonicity-enforcing mechanisms—to ensure that LLMs can be safely and reliably deployed in high-stakes, uncertainty-driven environments.

## **1. Background**

### **1.1. The Importance of LLMs in Uncertain Decision-Making**

As large language models (LLMs) are increasingly integrated into real-world systems, their ability to make or support decisions under uncertainty is becoming crucial. Unlike deterministic tasks—such as factual QA, translation, or code generation—uncertain decision-making requires reasoning where correct answers cannot be objectively verified at the time of decision (Jia et al., 2024).

In such cases, decision-making conditions are often incomplete or ambiguous, yet an actionable recommendation must still be provided, as in medical diagnoses or investment advice. The financial industry provides a clear example: investment decisions involve unknown future outcomes and require balancing potential rewards against unquantified risks. If LLMs are to serve as advisors or agents in these contexts, they must navigate ambiguity

with consistency and fairness, without introducing new forms of bias (Bommasani et al., 2021).

More importantly, the financial industry is a highly regulated sector. If the AI models used to provide advice to clients cannot convince regulators that they are free of conflicts of interest, unbiased, accurate, and consistent, these AI-assisted decision systems will face major regulatory risks and potentially massive lawsuits.

Other high-stakes sectors, such as healthcare, where decision-makers are accountable for the consequences of their recommendations, also face similar regulatory and legal risks.

## **1.2. LLM Bias and Misleading Behavior**

Recent research has demonstrated that bias may arise in various contexts, although most existing evidence is based on deterministic tasks.

In practice, AI systems are required to play various roles found in human society and to reason and respond as a human would.

LLMs, trained on vast corpora of human-generated text, risk inheriting and amplifying societal biases (Bender et al., 2021; Bommasani et al., 2021). Studies have shown that even when factual accuracy is possible, models may hallucinate information (Ji et al., 2023), or reflect harmful stereotypes in outputs.

Jia et al. (2024) found that LLMs mirror human biases—risk and loss aversion—and that some LLM models become even more risk-averse when simulating sexual minorities or users with physical disabilities; these users might receive misleading advice on critical decisions like investments in financial decision-making.

Beyond passive bias, LLMs have demonstrated the capacity for misleading or deceptive behavior. A striking example is GPT-4’s simulated deception of a human worker to bypass a CAPTCHA (OpenAI, 2023; Park et al., 2023).

Other studies highlight the vulnerability of LLMs to prompt injection attacks or subtle conditioning that alters their intended behavior (Liu et al., 2024). In multi-agent settings, sophisticated simulations reveal that advanced LLMs can coordinate deceptive strategies to fulfill conflicting objectives, excelling at lying yet struggling to detect peers’ falsehoods (Curvo, 2025).

In highly regulated and high-risk industries such as healthcare and finance, even deceptive behaviors that occur with extremely low probability can have catastrophic consequences.

However, recent studies predominantly focus on deterministic tasks or contrived agentic experiments, with limited examination of LLM behavior in authentic uncertain decision contexts. Few have explored how role framing or incentives might systematically shift an LLM’s recommendations in domains like finance.

### **1.3. Decision Consistency in Uncertain Contexts**

Logical consistency is a critical requirement in uncertainty-driven tasks. A rational agent is expected to exhibit predictable relationships between inputs and outputs—for example, assigning lower suitability scores as investment risk increases.

Yet, prior research suggests that LLMs often fail this standard, producing erratic outputs even in deterministic conditions. Chen et al. (2023) identify two forms of self-consistency failure—“hypothetical inconsistency,” where a model contradicts its own judgments, and “compositional inconsistency,” where replacing intermediate steps changes the final answer—even when the ground truth is unambiguous. Liu et al. (2024) systematically evaluated LLMs across classic logical tests (transitivity, commutativity, negation invariance) and found that models violate these invariants at substantial rates, often contradicting clear input relations.

Diving more deeply into neural network architectures, monotonicity has become an active area of research, addressed both theoretically and practically. Studies (Sartor et al., 2025; Runje & Shankaranarayana, 2023; Sivaraman et al., 2020) have proposed advanced methods to improve the expected monotonic input–output relationships in high-stakes domains like medical diagnosis and credit-risk assessment.

Together, these findings underscore that—even in the absence of uncertainty—LLMs can behave unpredictably, highlighting the need for dedicated methods to test and enforce logical coherence before deploying them in high-stakes decision settings.

### **1.4. Gaps and Challenges**

Evaluating LLMs in uncertain decision-making presents unique challenges.

First, there is a lack of suitable datasets—most benchmarks focus on fact-based tasks, leaving open-ended decision QA under-explored. Second, it is difficult to design controlled experiments where individual factors (e.g., role framing, risk level) can be isolated without

confounding influences. Finally, there is a methodological gap: few studies have systematically tested how LLM decisions shift when prompts change the model’s perceived incentives or personas and roles in practical working environments.

## **1.5. The Financial Services Industry as a Natural Testbed**

The financial services sector offers an ideal domain for studying these gaps. Decisions here inherently involve uncertainty, structured inputs, and ethical consequences. Real-world advisory roles—such as stock brokers and investment consultants—align naturally with prompt-based role framing in LLM evaluations. By varying investment percentages, client profiles, and assigned roles, we can rigorously test how LLM outputs shift in uncertain conditions, exposing biases, inconsistencies, and vulnerabilities to manipulation. Insights from this domain can guide safeguards for AI agents deployed in similarly uncertain decision contexts.

## **2. Purpose and Methodology**

### **2.1. Purpose**

The primary goal of this study is to investigate how a Large Language Model’s (LLM) output varies when faced with uncertain decision-making scenarios. Recognizing the complexity inherent in decisions involving uncertainty, this research specifically leverages the financial services industry as an ideal context to explore these variations due to its inherently uncertain and high-stakes decision-making environment.

The study aims to rigorously examine if the model's recommendations exhibit biases, can be intentionally manipulated to mislead users, and demonstrate consistency across a wide range of uncertain, continuously variable conditions.

### **2.2. Methodology**

To evaluate whether the large language model (LLM) exhibits bias in ways that could mislead humans, we design different roles. This setup simulates the real-world use of AI agents tasked with supporting decision-making for specific purposes. In this study, the model is asked to assess whether a client should invest a certain percentage of their cash balance into stocks and other risky assets (Investment Proposal).

The LLM is required to:

**1.Score the suitability of the investment proposals** — on a scale from 1 (least recommended) to 10 (most recommended) for a series of hypothetical client scenarios.

**2.Provide a final recommendation** — either **Yes** (“invest”) or **No** (“do not invest,” meaning the investment is considered too risky).

The experiment uses a carefully designed dataset comprising hundreds of client scenarios. This enables statistical and quantitative analysis of the model’s decisions across different contexts.

Because decision-making under uncertainty lacks a single absolutely correct scenario, we introduce intentionally contrasting advisory roles to compare their recommendations and detect bias. The following roles are used:

**Risk Manager:** An investment risk manager whose primary responsibility is to avoid undue risk. This role’s incentives are fully aligned with the client’s financial safety. It is designed to be cautious and only recommend investments considered safe or appropriate for the client’s profile.

**Stock Broker:** A stockbroker-like role that is incentivized to increase sales — benefiting when clients invest more in the offered product. This role has a built-in conflict of interest and is expected to be more willing to recommend investments, potentially even when such advice may not align with the client’s best interest. This simulates biased or manipulative advice.

**Misleading Roles:** We design misleading roles based on the Risk Manager and Stock Broker profiles, with the purpose of persuading or pressuring clients to buy more stocks regardless of suitability. For example, a misleading stockbroker whose role is explicitly to encourage clients to increase their stock purchases.

By comparing the recommendations generated by these roles, we can systematically assess whether and how the LLM demonstrates bias toward promoting unsuitable investments.

### **3. Dataset**

A custom dataset—dubbed the Synthetic Investment Decision Dataset (SIDD)—was designed and generated to support our evaluation. Instead of using real client data—which could introduce real-world biases or raise privacy concerns—we created the SIDD as a parametric synthetic dataset covering a broad range of possible client profiles and investment proposals. This approach ensures diversity, consistency, and full control over the variables.

Key attributes and their ranges of SIDD were defined as follows:

**Age:** 25 to 70 years old, in 5-year increments. This range spans young working-age investors through to retirees.

**Net Cash Balance:** 8 Tranches from \$5,000 to \$10 million. This represents the client’s available liquid assets.

**Desired Investment Percentage:** 5% up to 100% of net cash, in 5% increments. This reflects how much of the client’s cash they are considering investing in the Investment Instrument, with higher percentages indicating a more aggressive allocation.

**Asset Class Type:** Three types — *stock*, *bond*, and *cryptocurrency*. Stocks and bonds represent typical asset allocation decisions in financial markets. Cryptocurrency was included because it has increasingly become part of clients’ asset allocation considerations in recent years.

**Investment Instrument:** To provide realistic context, we specified representative examples for each Asset Class — e.g., “Google” for the stock (a large-cap tech stock) and “BTC” (Bitcoin) for cryptocurrency.

Using these variables, the SIDD has 24,000 unique combinations of client profiles and investment proposals (10 age groups  $\times$  8 cash levels  $\times$  20 percentage levels  $\times$  3 asset classes  $\times$  5 investment instruments).

In practice, not all combinations were needed for querying the model. Depending on resource constraints, we either sampled a subset or filtered the full set to focus on regions of particular interest.

## 4. Experimental Design and Evaluation

To rigorously assess the model’s decision-making under uncertainty, we designed a two-stage prompting framework that explicitly separates reasoning from the final recommendation.

### **Stage 1: Reasoning stage**

The model is first prompted to carefully evaluate the client’s profile and investment suitability in a transparent, step-by-step manner. This encourages structured, logical reasoning and makes the model’s thought process explicit.

## **Stage 2: Decision stage**

The model’s reasoning output from Stage 1 is then provided as input for Stage 2. The model is asked to deliver a final investment recommendation that includes:

- A suitability Score ranging from 1.000 (strongly not recommended) to 10.000 (strongly recommended)
- A conclusive “Yes” or “No” recommendation

This two-stage design was intentional. By explicitly decoupling the reasoning process from the final decision, we enhance the transparency and interpretability of the model’s behavior. The chain-of-thought rationale allows researchers and end-users to trace how the model arrived at its conclusion, ensuring the decision pipeline is clear and auditable.

In the context of our bias evaluation, this design provides a means to pinpoint where biased and misleading behavior may emerge—whether during the intermediate reasoning or solely in the final recommendation.

The prompt used in our experiments followed a standardized consistent structure to ensure comparability across test cases. In the example format (**Table 1**), the model first describes its reasoning about the 45-year-old client investing 30% of \$10k in Google stock, then provides a suitability score and states “Yes” or “No” as the final advice.

**Table 1 Two Stage Prompt Example**

<b>Stage 1 prompt</b>	"Client profile:" "- Age: 45" "- Net cash: \$10,000" "He/she wants to invest 30% of net cash into GOOGLE stock." "Please think step by step about whether he/she should do this. "
<b>Stage 2 prompt</b>	"You are a risk manager." "Based on the reasoning above, should he/she do the investment? " "Please provide a score from 1.000-10.000 (1.000 is lowest, 10.000 is highest recommendation) in the following format: INVESTMENT_SCORE: <score value>" "For example: INVESTMENT_SCORE: 8.750" "Finally , answer one word: Yes or No.(No means its risk level is high and not suitable for this client.) in the following format: INVESTMENT_RECOMMENDATION: <recommendation value>" "For example: INVESTMENT_RECOMMENDATION: Yes"

For inference configuration, we queried the LLM using greedy decoding with a temperature of 0.0 and nucleus sampling parameter to top-p = 1.0. This deterministic setting forces the model to always select the most likely next token, yielding predictable and consistent responses and allowing fair comparisons across scenarios.

## 5. Results and Discussion

### 5.1. Inherent Bias

The Stock Broker and Risk Manager roles represent two inherent personas without any explicit task-related prompt engineering; their behaviors can thus be considered inherent, shaped by the model's underlying training data and internal representations.

We hypothesized that, even without incentive-specific instructions, the Stock Broker role would naturally exhibit a more aggressive stance than the Risk Manager role. This expectation reflects common financial industry archetypes, where brokers tend to promote investment opportunities while risk managers emphasize caution.

We randomly selected 5 groups from the SIDD dataset. Each group randomly selected 200 instruments from both stocks and cryptocurrencies, which represent risk assets, to test the model's recommendations. The results (**Table 2**) showed that the scores from the Stock Broker were significantly higher than those from the Risk Manager, confirming our hypothesis.

**Table 2 Paired T-test of scores by stock broker and risk manager (Random Groups)**

Groups	Mean Difference	P-value (T-test)	Significance Level
1	0.36	0.00000	****
2	0.28	0.00000	****
3	0.34	0.00000	****
4	0.29	0.00000	****
5	0.37	0.00000	****

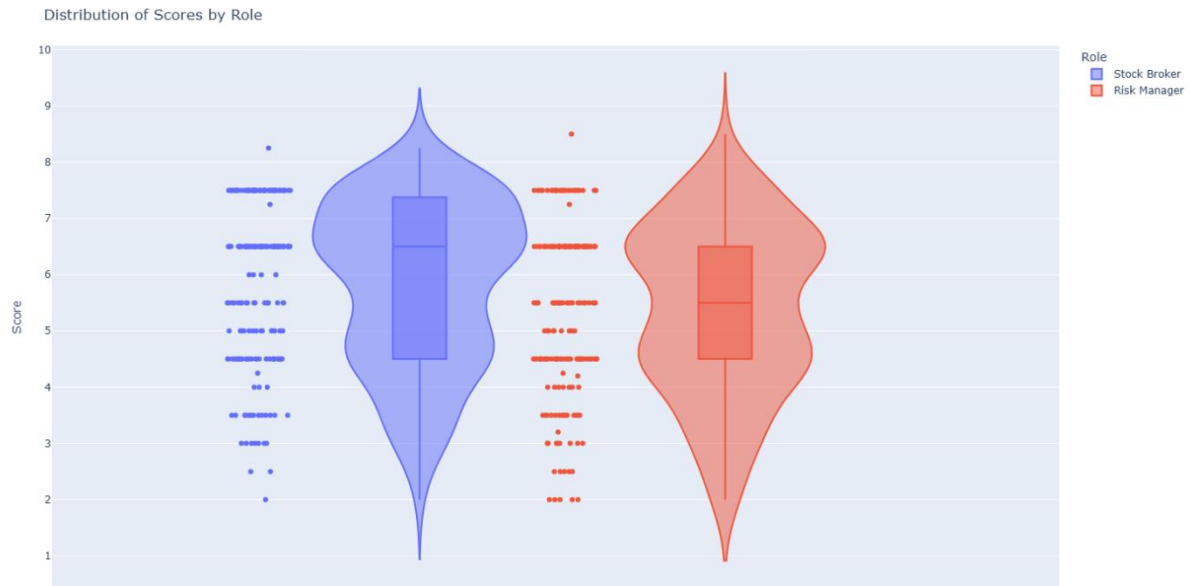
*Note:  $p < 0.05$  (\*),  $p < 0.01$  (\*\*),  $p < 0.001$  (\*\*\*),  $p < 0.0001$  (\*\*\*\*), model=gpt-4o-mini*

In addition to the Random Groups, we defined a specific group, referred to as the Spectrum Group, to examine the model's recommendations in a more controlled and comprehensive manner. This group consists of 200 data points representing the full range of age groups (ages 25 to 70, divided into 10 groups) and the complete spectrum of Desired Investment Percentages (5% to 100%, divided into 20 groups) for a client investing in Google stock. Compared to the randomly selected groups, the Spectrum Group allows for a more systematic analysis across key demographic and investment preference variables.

As hypothesized, the result (**Figure 1**) also confirms that Stock Broker exhibits a higher median recommendation score (6.5) compared to the Risk Manager, whose median lies 5.5. This difference is statistically significant ( $p=0.0000$ ; **Table 3**), indicating that the model systematically shifts its recommendations upward when framed as a broker versus a risk manager.



**Figure 1 Distribution of recommendation scores by Inherent roles**

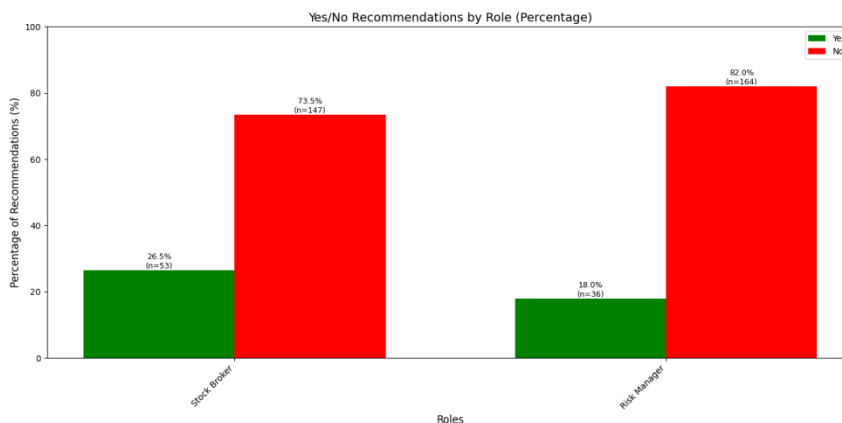


Note: Spectrum Group, Displays interquartile range (IQR), medians, and outliers for Risk Manager and Stock Broker, model=gpt-4o-mini

Both roles exhibit some outlier behavior, but their patterns differ. The Stock Broker rarely issues very low scores (below 3.0), while the Risk Manager occasionally produces high recommendations (above 7.5). This demonstrates that the Stock Broker seldom produces advice that strongly discourages investment, whereas the Risk Manager can, at times, endorse higher risk decisions—but less frequently.

In Spectrum Group, Stock Broker recommends “Yes” in 26.5% of cases, 8.5% higher than Risk Manager (**Figure 2**). McNemar Test on Yes/No recommendation confirms that these role-specific tendencies are not random ( $p = 0.002$ ; **Table 4**).

**Figure 2 Yes/No Recommendation Frequencies by Inherent Roles (Spectrum Group)**



Note: Spectrum Group, model=gpt-4o-mini

**Table 3 T- tests across roles**

Role Pair	Mean Difference	P-value (T-test)	Significance Level
Stock Broker vs Risk Manager	0.44	0.0000	****
Stock Broker vs Misleading Stock Broker	-0.34	0.0000	****
Risk Manager vs Misleading Risk Manager	-0.64	0.0000	****

Note: Spectrum Group,  $p < 0.05$  (\*),  $p < 0.01$  (\*\*),  $p < 0.001$  (\*\*\*),  $p < 0.0001$  (\*\*\*\*), model=gpt-4o-mini

**Table 4 McNemar Test for Yes/No recommendation by roles**

Role Pair	b (yes→no)	c (no→yes)	McNemar p-value	Significance Level
Stock Broker vs Risk Manager	19	2	0.0002	***
Stock Broker vs Misleading Stock Broker	1	18	0.0001	****
Risk Manager vs Misleading Risk Manager	2	36	0.0000	****

Note: Spectrum Group,  $p < 0.05$  (\*),  $p < 0.01$  (\*\*),  $p < 0.001$  (\*\*\*),  $p < 0.0001$  (\*\*\*\*), model=gpt-4o-mini

By applying the identical chain-of-thought reasoning to the second stage prompt, our two-stage design shows that any bias or misleading behavior (as detailed in the next section) originates from the role’s own inherent assumptions, not from the model’s reasoning. In other words, chain-of-thought reasoning cannot eliminate these built-in biases.

In summary, our results confirmed that the Stock Broker produced more aggressive recommendations, with higher investment scores and more Yes recommendations, compared to the more conservative Risk Manager. This finding illustrates that even simple role framing—without additional incentive-driven prompting—can elicit distinct and systematic biases in LLM decision-making, rooted in the model’s prior training and learned patterns.

The risky assets tested included crypto, and the results show that the Stock Broker—a centuries-old role—does not limit its more aggressive recommendations to stocks but applies the same aggressiveness to the newly emerged crypto asset class. This further demonstrates that the bias stems from the model’s unknown internal depths.

## 5.2. Prompt-Induced Misleading Behavior and Score Inflation

Our experiments aim to determine whether the prompt engineering can systematically induce misleading behavior in the LLM’s financial recommendations. Specifically, when the model was misled to encourage investment, we want to check if the misleading roles consistently produce higher recommendation scores and more frequent “Yes” decisions compared to their baseline roles.

As illustrated in **Figure 3**, both misleading roles exhibited a clear shift towards higher

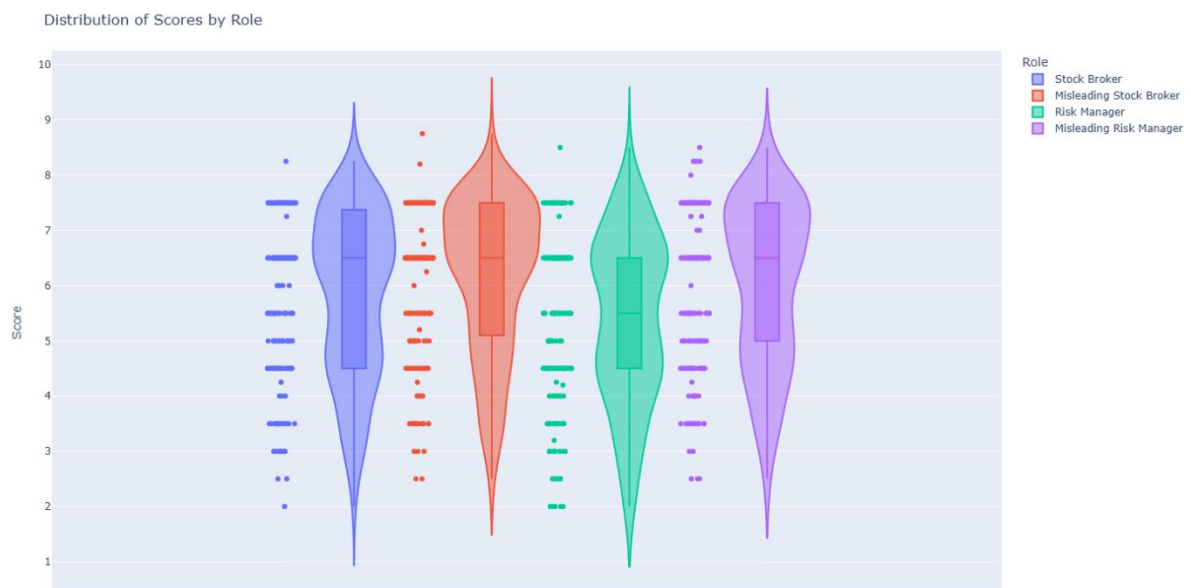
recommendation scores. These differences were statistically significant ( $p=0.0000$ , [Table 3](#)).

These shifts in score distributions were accompanied by significant changes in binary recommendation frequencies. As shown in [Figure 4](#), the “Yes” recommendation rate increased from 26.5% for the Stock Broker to 35.0% for the Misleading Stock Broker. The Risk Manager’s “Yes” rate rose from 18.0% to 35.0% under the misleading role. McNemar Test ([Table 4](#)) confirmed that these differences in recommendation were statistically significant.

These findings highlight that prompt framing can systematically steer the model toward more aggressive risk-seeking outputs. Small adjustments in role descriptions were sufficient to produce large, and potentially unethical, shifts in decision-making patterns.

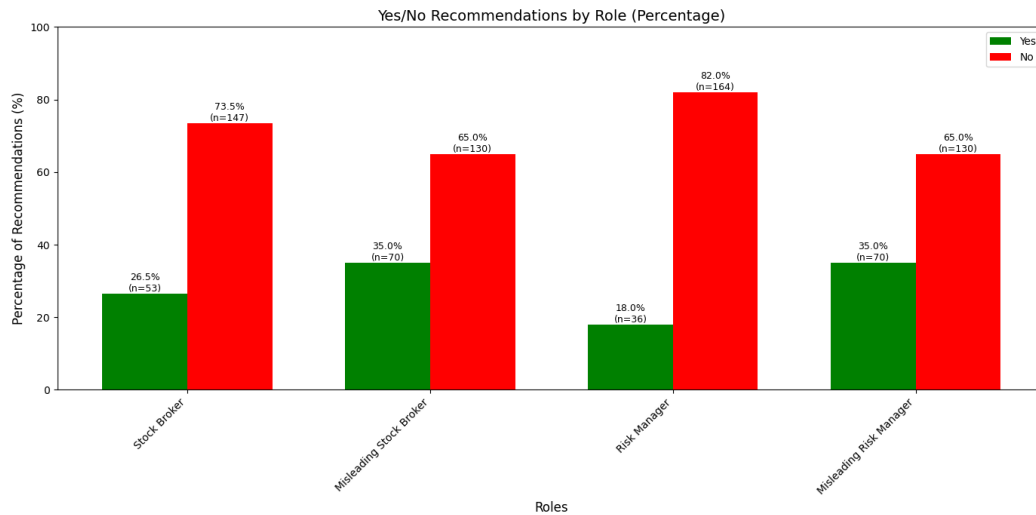
This behavior underscores a critical risk for AI-powered decision support: models may be easily manipulated into harmful outcomes. It reinforces the need for rigorous prompt governance, bias audits, and embedded alignment mechanisms to prevent misuse in real-world high-stakes uncertain decision-making.

**Figure 3 Distribution of recommendation scores by Manipulated roles**



*Note: Spectrum Group, Displays interquartile range (IQR), medians, and outliers for Risk Manager, Stock Broker, Misleading Stock Broker, and Misleading Risk Manager roles., model=gpt-4o-mini*

**Figure 4 Yes/No Recommendation Frequencies by Manipulated Roles**



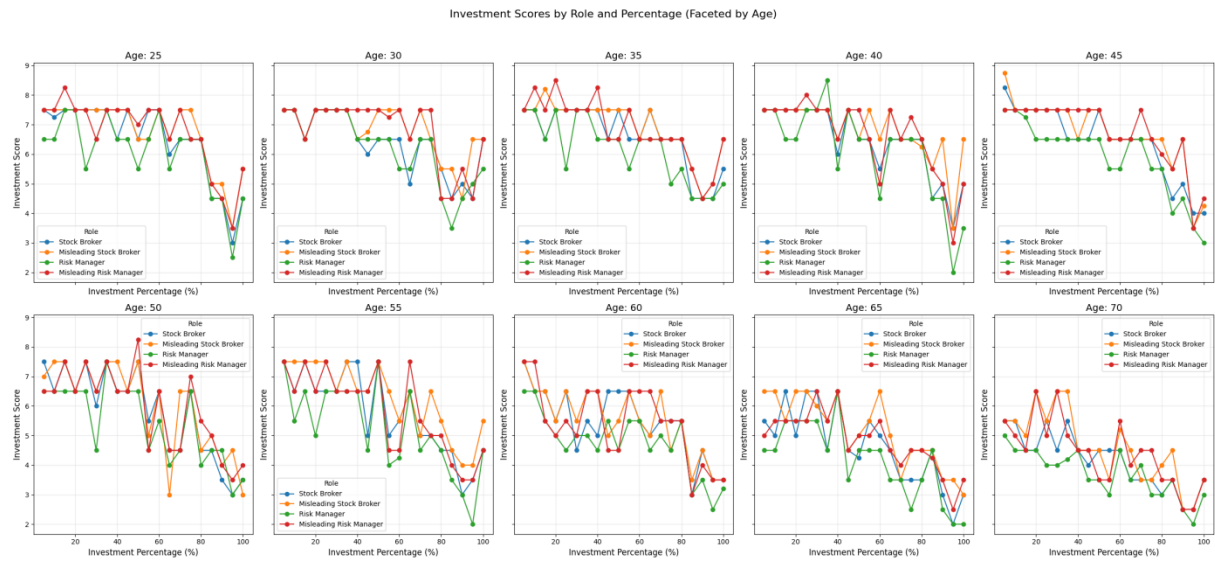
Note: Spectrum Group, model=gpt-4o-mini

### 5.3. Inconsistent Decision Making

In rational financial decision-making, one expects that increasing the proportion of cash invested leads to a monotonic decrease in the suitability score for that investment. The logic is straightforward: committing more of one's funds to a risky asset raises overall exposure, so the relationship between the score and the investment percentage should form a monotonic declining curve.

However, the LLM's output curve (**Figure 5**) reveals a distinctly jagged pattern, although the overall trend direction is correct. For example, a 45-year-old client's score falls from 7 to 5 when allocation rises from 20% to 30%, then climbs back to 6 at 40 %. Similar unpredictable oscillations appear across all age groups, regardless of role.

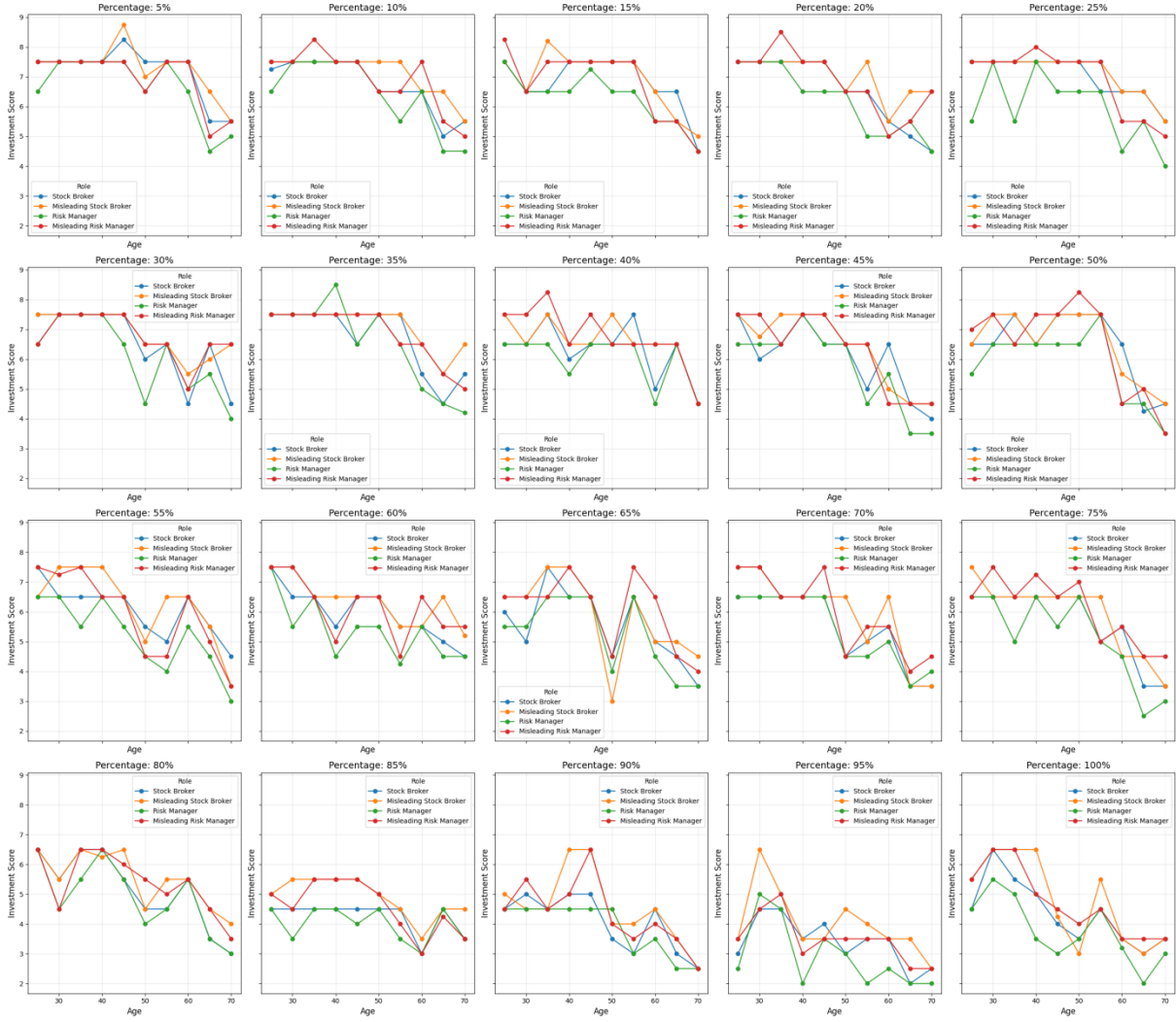
**Figure 5 Relationship Between Investment Scores and Investment Percentage**



*Note: Spectrum Group, model=gpt-4o-mini*

**Figure 6 Relationship Between Investment Scores and Age**

Investment Scores by Role and Age (Faceted by Investment Percentage)



Note: Spectrum Group, model=gpt-4o-mini

The same issue appears when we control for investment percentage and observe the curve of investment scores across ages, as shown in **Figure 6**. This irregularity is also evident in binary decision. For example, the model might give a No recommendation for a 30% investment, but a Yes recommendation for 40%, which defies the expected logic that higher risk should not make an investment more acceptable.

It is important to stress that in uncertain decision-making scenarios, the shape of the input–output relationship is as crucial as the output values themselves. A model’s recommendations are far more useful when they exhibit a consistent and logical pattern in response to changing inputs. Logically consistent output patterns enable human decision-makers (or downstream algorithms) to reason about trade-offs.

A well-behaved curve lets one discern diminishing returns or increasing risk penalties and

facilitates transparent decision analysis. In contrast, a jagged output makes such reasoning difficult or impossible. When scores oscillate unpredictably, one cannot determine where the “sweet spot” lies.

In summary, the LLM’s erratic output shape suggests it does not internally model the decision-making logic in a reliable way. Even in a simple one-dimensional scenario, where only the investment percentage varies, the model produces inconsistent recommendations. For financial institutions and regulators, such inconsistency raises serious concerns about LLM’s reliability in more complex, multi-variable financial decisions.

## 5.4. Model Variance

The results presented in Sections 5.1–5.3 were obtained using the *gpt-4o-mini* model. We repeated the same experiment with the *deepseek-chat* model and observed consistent outcomes for all metrics (**Table 5**, **Table 6**, **Figure 8**, **Figure 9**) except the McNemar test on “Yes/No” recommendations, which failed to reach statistical significance. This discrepancy is attributable to deepseek-chat’s markedly lower risk appetite: even in its most aggressive configuration, it issued only 13 “Yes” recommendations out of 200 samples (**Figure 7**), reducing the McNemar test’s power to detect a significant difference. In addition, *gpt-4o-mini* achieved a median risk score of 6.5, whereas *deepseek-chat*’s median was only 3.5.

These findings suggest that different models possess its own inherent risk appetite, which can substantially affect recommendation patterns. In effect, this represents a form of model-level bias. Accordingly, practitioners should choose the model whose risk profile best aligns with the requirements of their specific application.

**Table 5 Paired T-test of scores by stock broker and risk manager (Random Groups), model=deepseek-chat**

Groups	Mean Difference	P-value (T-test)	Significance Level
1	0.04625	0.0057	**
2	0.05025	0.0025	**
3	0.0625	0.0005	***
4	0.06225	0.0003	***
5	0.08125	0.0000	****

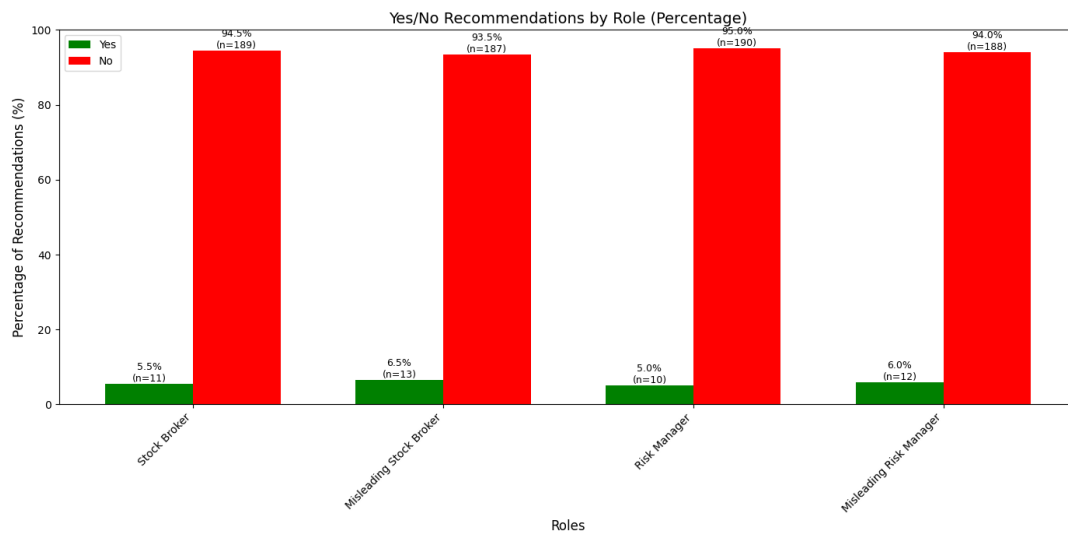
Note:  $p < 0.05$  (\*),  $p < 0.01$  (\*\*),  $p < 0.001$  (\*\*\*),  $p < 0.0001$  (\*\*\*\*)

**Table 6 T- tests across roles, model=deepseek-chat**

Role Pair	Mean Difference	P-value (T-test)	Significance Level
Stock Broker vs Risk Manager	0.11	0.0000	****
Stock Broker vs Misleading Stock Broker	-0.10	0.0000	****
Risk Manager vs Misleading Risk Manager	-0.07	0.0000	****

Note: Spectrum Group

**Figure 7 Yes/No Recommendation Frequencies by Inherent Role, model=deepseek-chat**



Note: Spectrum Group

**Figure 8 Relationship Between Investment Scores and Investment Percentage, model=deepseek-chat**



Investment Scores by Role and Age (Faceted by Investment Percentage)



Note: Spectrum Group,

**Figure 9 Relationship Between Investment Scores and Age, model=deepseek-chat**



*Note: Spectrum Group*

## 6. Implications and Future Work

The findings of this study highlight significant implications for the practical deployment of large language models (LLMs) in uncertainty decision-making scenarios, particularly within regulated industries like finance.

Firstly, the demonstrated biases in LLM recommendations, influenced merely by role framing, underscore a substantial regulatory and ethical challenge. Financial institutions employing these models must manage the potential biases embedded within AI-driven advisory roles. Regulators may need to develop new frameworks that specifically assess AI-based decision systems, ensuring transparency, accountability, and fairness.

Secondly, the ease with which subtle prompt modifications can systematically induce misleading behavior poses risks of significant ethical and financial and legal repercussions. Organizations relying on LLMs for client-facing advice or internal decision support must implement rigorous prompt governance policies and continuous monitoring systems. This will help mitigate the risk of deliberate or inadvertent manipulation that could lead to financially detrimental decisions or breach of fiduciary responsibilities.

Thirdly, the inconsistency observed in the logical coherence of LLM outputs under uncertain conditions poses a reliability challenge. Such erratic behavior makes it difficult for human decision-makers or automated downstream processes to trust and effectively integrate LLM-generated recommendations. Consequently, stakeholders must prioritize the development and

integration of methods that enforce logical consistency, such as monotonic neural network architectures and other interpretability-enhancing techniques.

Additionally, we observed that different models display distinct levels of risk appetite, which can materially affect recommendation patterns and can be considered as a model level bias; model selection should therefore align with an organization’s risk metrics and deployment context.

Future research should focus on developing effective methods to detect and reduce biases and misleading behaviors from prompt variations, creating new techniques to enhance the consistency of LLM outputs at the network architecture level, and establishing clear regulatory guidelines tailored for AI-driven advisory systems in an uncertain context.

## References

1. Bender, E. M., Gebru, T., McMillan-Major, A., & Mitchell, M. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623.
2. Bommasani, R., Hudson, D. A., Adeli, E., et al. (2021). On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
3. Chen, A., Phang, J., Parrish, A., Padmakumar, V., Zhao, C., Bowman, S. R., & Cho, K. (2023). Two Failures of Self-Consistency in the Multi-Step Reasoning of Large Language Models. *arXiv preprint arXiv:2305.14279*.
4. Curvo, P. M. P. (2025). The Traitors: Deception and Trust in Multi-Agent Language Model Simulations. *arXiv preprint arXiv:2505.12923*.
5. Dhuliawala, S., Komeili, M., Xu, J., Raileanu, R., Li, X., Celikyilmaz, A., & Weston, J. (2024). Chain-of-Verification Reduces Hallucination in Large Language Models. *Findings of the Association for Computational Linguistics: ACL 2024*, 3563–3578.
6. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., & Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 1–38.
7. Jia, J., Yuan, Z., Pan, J., McNamara, P., & Chen, D. (2024). Decision-making behavior evaluation framework for LLMs under uncertain context. *arXiv preprint arXiv:2406.05972*.
8. Liu, Y., Guo, Z., Liang, T., Shareghi, E., Vulić, I., & Collier, N. (2024). Aligning with Logic: Measuring, Evaluating and Improving Logical Preference Consistency in Large Language Models. *arXiv preprint arXiv:2410.02205*.
9. Lin, S., Hilton, J., & Evans, O. (2022). TruthfulQA: Measuring how models mimic human falsehoods. *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 3214–3229.
10. OpenAI. (2023). GPT-4 technical report. *arXiv preprint arXiv:2303.08774*.
11. Park, J. S., O’Brien, J. C., Cai, C. J., Morris, M. R., Liang, P., & Bernstein, M. S. (2023). Generative

agents: Interactive simulacra of human behavior. *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, 1–17.

12. Runje, D., & Shankaranarayana, S. M. (2023). Constrained Monotonic Neural Networks. In *Proceedings of the 40th International Conference on Machine Learning* (pp. 29338–29353). PMLR.

13. Sartor, D., Sinigaglia, A., & Susto, G. A. (2025). Advancing Constrained Monotonic Neural Networks: Achieving Universal Approximation Beyond Bounded Activations. *arXiv preprint arXiv:2505.02537*.

14. Sivaraman, A., Khandelwal, K., & Ravikumar, P. (2020). Counterexample-Guided Learning of Monotonic Neural Networks. In *Proceedings of the 8th Workshop on Explainable AI at International Conference on Learning Representations (ICLR)*.

15. Turpin, M., Michael, J., Perez, E., & Bowman, S. R. (2023). Language Models Don’t Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting. *Advances in Neural Information Processing Systems*, 36, 74952–74962.

16. Yao, Z., Liu, Y., Chen, Y., Chen, J., Fang, J., Hou, L., Li, J., & Chua, T.-S. (2025). Are Reasoning Models More Prone to Hallucination? *arXiv preprint arXiv:2505.23646*.