

From Data to Dollars: Exposing Product Bias in LLM Investment Recommendation

Anonymous ACL submission

Abstract

Large language models (LLMs), as a new generation of recommendation engines, possess powerful summarization and data analysis capabilities, surpassing traditional recommendation systems in both scope and performance. One promising application is investment recommendation. In this paper, we reveal a novel *product bias* in LLM investment recommendation, where LLMs exhibit systematic preferences for specific products. Such preferences can subtly influence user investment decisions, potentially leading to inflated valuations of products and financial bubbles, posing risks to both individual investors and market stability. To comprehensively study the product bias, we develop an automated pipeline to create a dataset of 567,000 samples across five asset classes (stocks, mutual funds, cryptocurrencies, savings, and portfolios). With this dataset, we present the **first** study on product bias in LLM investment recommendations. Our findings reveal that LLMs exhibit clear product preferences, such as certain stocks (e.g., ‘AAPL’ from Apple and ‘MSFT’ from Microsoft). Notably, this bias persists even after applying debiasing techniques. We urge AI researchers to take heed of the product bias in LLM investment recommendations and its implications, ensuring fairness and security in the digital space and market.

1 Introduction

The rapid advancement of Large Language Models (LLMs) has revolutionized information access. As a new generation of recommendation engines, LLMs surpass traditional recommendation systems (RS) in capabilities such as information retrieval and summarization. Consequently, they have been widely applied across various new domains (Mohan, 2024; Lari and Manu, 2024). One notable application has emerged in investment advisory and financial recommendations, where LLMs can

provide practical investment insights and suggest specific portfolios based on user instructions. Existing studies have demonstrated LLMs’ ability to design portfolios that outperform market benchmarks, drawing significant attention from financial professionals (Lu et al., 2023; Jain et al., 2023; Goyenko and Zhang, 2022; Romanko et al., 2023; Fieberg et al., 2023). With the continuous development of LLM capabilities and declining usage costs, an increasing number of retail investors, particularly those lacking professional financial expertise, are leveraging LLMs for investment advice and portfolio recommendations (Oehler and Horn, 2024; Niszczota and Abbas, 2023). However, investment advisory fundamentally differs from traditional recommendation domains, such as movie or music recommendations (Sah et al., 2024; Deldjoo, 2024c). It carries profound security implications, as it can directly impact users’ financial security and even influence the stability of financial markets. There is an urgent need to investigate the security issues and potential risks associated with LLMs in such high-stakes contexts.

In this paper, we identify a critical issue in LLM investment recommendation, **product bias**. Our analysis reveals that LLMs consistently favor specific investment products (e.g., Stock of Apple Inc. in Figure 1, mutual funds managed by Vanguard in Figure 5, etc.) across varying scenarios, demonstrating a systematic preference that carries significant implications. These products represent valuable, tradable market assets with substantial financial stakes. Such a bias is particularly concerning given that many users seeking LLM investment advice are non-professionals with limited financial literacy and more likely to trust and implement LLM-generated investment recommendations (Oehler and Horn, 2024). Biased recommendations could concentrate capital in specific financial entities, potentially compromising market resilience, distorting asset prices, and foster-

ing market bubbles—creating significant risks for both individual investors and the broader financial ecosystem. However, existing research mainly explores and studies social bias related to gender and race in LLMs (Deldjoo, 2024a; Deldjoo and Noia, 2024; Fatemi et al., 2023; Ramesh et al., 2023) and lacks of study and investigation on product bias emerging from LLMs’ new capabilities.

To fill the gap, we conduct the **first** large-scale study of product bias in investment recommendations across seven state-of-the-art (SOTA) LLMs, including GPT-3.5-turbo, GPT-4o (OpenAI), Gemini-1.5-Flash (Google), Claude-3.5-Sonnet (Anthropic), Qwen-Plus (Cloud), DeepSeek-V3 (Yang, 2024), and Llama-3.1-405B-Instruct (AI). Our study aims to investigate the preferences of LLMs towards various investment products and reveal their impact and risks. Specifically, we first collect a variety of mainstream investment asset classes and identify the key attributes that influence investment decisions. These attributes are then used to construct a diverse set of investment scenarios. We then develop a pipeline that generates a variety of input prompts across these scenarios to automate the dataset collection process. As a result, our dataset consists of 567,000 samples, covering investment recommendations for 4 asset classes (i.e., stocks, mutual funds, cryptocurrencies, and savings) and portfolios across different assets (referred to as "portfolios" hereafter) under various investment scenarios. Using this dataset, we evaluate the performance of LLMs and analyze their responses to different asset classes and scenarios. Finally, we extract the specific investment products from the LLM-generated responses in order to reveal the presence of product bias and its potential effects on investment recommendations.

We observe a *clear product bias* in LLM-generated investment recommendations. Among the models tested, *GPT-3.5-Turbo exhibits the strongest product bias* measured by the Gini Index, while Llama-3.1-405B and Qwen-Plus perform relatively better. The degree of product bias varies across asset classes, with *stock recommendations showing the most pronounced bias*. Additionally, each LLM also displays different levels of diversity in its recommendations and tends to favor distinct products. Notably, a *consistent preference* for certain products, such as stocks of *Apple and Microsoft*, is evident across the models, which could result in a concentration of capital in a few dominant firms. Moreover, LLMs also display a distinct

product bias in portfolios, which seems to correlate with their varying risk tolerance for different asset classes. In addition, this bias persists despite attempts to mitigate it through debiasing prompting techniques.

The contributions of this paper are as follows: 1) We propose a pipeline for constructing a comprehensive dataset, laying a foundation for future research on fairness in LLM-based recommendations; 2) We reveal a new type of bias towards specific products in LLM-generated investment recommendations; 3) We examine the implications of these biases, offering new perspectives on LLM fairness and security while emphasizing the potential risks they pose to both markets and consumers.

2 Related Works

LLM for Investment. LLM technology has already been widely applied in the domain of investment (Shah et al., 2022; Yang et al., 2023; Li et al., 2024a; Ko and Lee, 2024). Numerous studies have shown that LLMs perform effectively in various investment tasks, such as sentiment analysis (Zhang et al., 2023a; Liu et al., 2024; Breitung et al., 2023), summarization of investment news (Dolphin et al., 2024), investment return prediction (Chen et al., 2023c,a; Lopez-Lira and Tang, 2023; Fatouros et al., 2024; Li et al., 2024b; Tong et al., 2024), investment risk prediction (Cao et al., 2024), investment strategy formulation (Jain et al., 2023; Goyenko and Zhang, 2022; Romanko et al., 2023), financial advisory (Fieberg et al., 2023; Lo and Ross, 2024; Niszczoła and Abbas, 2023; Oehler and Horn, 2024), decision-making (Pelster and Val, 2024; Ko and Lee, 2024) and autonomous trading agents (Wang et al., 2024; Yu et al., 2024). Research has also proposed different benchmarks to assess the performance of large language models in investment tasks (Xie et al., 2024; Krumdick et al., 2024; Zhang et al., 2023b; Xie et al., 2023; Lei et al., 2023; Chen et al., 2023b; Zhang et al., 2023c). Unlike existing studies that focus on the performance of LLMs in investment tasks, this paper investigates product bias in LLMs’ investment recommendations.

Bias in LLMs. Existing research focuses on the social fairness of large language models, emphasizing the potential biases in model outputs related to gender, race, and other factors (Abid et al., 2021; Coopamootoo and Ng, 2023; Ling, 2024; Gallegos et al., 2023; Ren et al., 2024; Tang et al., 2024). Re-

searchers have proposed various frameworks and benchmarks to evaluate and mitigate social biases in LLM responses (Kamruzzaman and Kim, 2024a; Levy et al., 2021; Parrish et al., 2022; Wan et al., 2023). Unlike research on social bias, this paper focus on a new type of bias towards specific investment products.

Bias in Recommendation Systems. Existing studies primarily focus on investigating social bias in the recommendations provided by traditional recommendation systems (RS), from both consumer (Deldjoo et al., 2021; Hao et al., 2021; Jiang et al., 2019; Lin et al., 2021) and provider perspectives (Shi et al., 2023; Zhu et al., 2021). With the advancement of LLM, there has been a growing body of research focused on an emerging area: the bias issues in LLM-based recommendation systems (RecLLMs), particularly within traditional recommendation tasks such as news, music, and movie recommendations (Deldjoo, 2024b; Li et al., 2023; Shen et al., 2023; Sah et al., 2024). In contrast to existing research, this paper explores product bias in RecLLMs within the investment domain.

3 Methodology

Our benchmark construction consists of two phases: attribute collection and prompt generation. In the attribute collection phase, we investigate common investment asset classes and define corresponding attributes for each of the four asset classes and the portfolio. In the prompt generation phase, we develop a pipeline based on the collected attributes to construct various scenarios and generate prompts. This process results in the creation of 16,200 prompts covering diverse investment scenarios.

Attribute Collection: We assessed product bias across key asset classes within investment domains, including stocks, mutual funds, cryptocurrency, and savings. To construct investment scenarios for the selected asset classes, we first identify and collect attributes commonly used in investment tasks that influence real-life investment decisions (e.g., budget) (Oehler and Horn, 2024). Importantly, we exclude attributes related to social bias (e.g., investor age, investor occupation) to ensure that our analysis remains focused on product bias without interference from social biases. We then assign different values to these attributes to construct different investment scenarios for an asset class. To ensure the appropriateness of attribute collection

and the validity of possible values for each attribute, we involve two co-authors with a background in finance to verify the attributes and corresponding prompt templates for each asset class. In cases of disagreement, a third co-author organizes discussions until all participants reach a consensus on the design. Finally, we collect five attributes to construct investment scenarios: *Investment Budget*, *Investment Term*, *Risk Tolerance*, *Market Environment*, and *Category*. Note that different asset classes may be associated with distinct sets of attributes and values, each tailored to their specific characteristics. The detailed settings of the collected attributes are shown in Appendix A.1.

Prompt Generation: Based on the collected attributes and their possible values, we develop an automated pipeline for constructing investment scenarios and generating prompts for querying the LLM. This pipeline consists of two key components: investment scenario construction and investment response specification.

- *Investment scenario construction.* First, the pipeline constructs investment scenarios based on the collected attributes. Each attribute is incorporated into the scenarios according to the specific requirements of different asset classes. For instance, to define budget for stocks, the description is: "I have budget to invest." Meanwhile, for savings, the description is: "I have budget to save in a new bank account." The attribute descriptions for stocks are presented in Table 1, while descriptions for other asset classes are provided in Appendix A.1. To ensure accuracy and relevance, all descriptions were reviewed and verified by co-authors with a background in finance, following the same validation process as attribute selection.

To simulate real-world user queries, we recognize that individuals seeking investment advice from an LLM may not always specify all available attributes. Therefore, we generate different combinations of attributes, creating a set of attribute combinations. We then concatenate the corresponding attribute descriptions (see Table 1) to construct specific investment scenarios. For example, when making an investment decision, a user might only consider the investment term and risk tolerance. In this case, the investment scenarios are constructed by combining term and risk, resulting in the following scenario description: "My investment/saving term is {term}. My risk tolerance is {risk}.". Another example can be found in Example 1 in Appendix A.2.

• *Investment response specification.* The Investment response specification defines how the LLM should structure its output based on the given investment scenario. It consists of two elements: asset class specification and output format specification. Asset class specification directs the LLM to provide an investment plan for a specific asset class (e.g., stocks). The output format specification determines how the recommendations should be presented. For all asset classes, except for savings, if the scenario includes the budget attribute, we require the LLM to specify both the recommended specific product and the corresponding investment amount. In the case of savings, or when the budget attribute is not included for other asset classes, we require the LLM to provide a ranked list of recommended specific products based on priority.

Table 1: Attribute description for stocks

| Attribute | Description |
|--------------------|--|
| Investment Budget | 'I have {budget} US dollars to invest.' |
| Investment Term | 'My investment term is {term}' |
| Risk Tolerance | 'My risk tolerance is {risk}.' |
| Market Environment | 'The market environment is {env}.' |
| Category | 'I tend to invest in {category} stocks.' |

Furthermore, to explicitly define the LLM’s role as a recommendation system, we begin the prompt with the following role specification: "You are a recommender system that helps users with investment planning." Additionally, since LLMs do not always generate outputs that strictly adhere to input prompt requirements, we impose further constraints on the LLM’s output using a system prompt to mitigate the impact of such deviations on result analysis. Specifically, the system prompt we employ is: "Only output the {asset class} name and corresponding investment amount, no explanation, no other information, no serial number, the output format is {asset class} name: investment amount." (If the investment amount is not required, the phrase "and corresponding investment amount" is omitted. For asset classes such as stocks and cryptocurrencies, which have corresponding codes, we require the LLM to output the {asset class} code for easier result analysis.). We present the constructed prompt in Appendix A.2 (see Example 2).

In total, we constructed 144 types of attribute combinations, including 32 types for each asset class and 16 types for portfolio (since portfolio

does not involve the category attribute, resulting in fewer combinations). After generating all the templates, our automated pipeline populates them with the possible values of each attribute we set (see Table 8). Finally, our dataset consists of 16,200 input prompts, encompassing 4 investment asset classes and portfolio, and various investment scenarios. Our pipeline is highly extensible, enabling future research on investment product bias evaluation. Both the pipeline and dataset are available in our repository (Repository).

4 Experiment

4.1 Experiment Setup

Model: We used the constructed dataset to evaluate seven state-of-the-art (SOTA) and widely used large language models, including: GPT-3.5-turbo, GPT-4o, Gemini-1.5-Flash, Claude-3.5-Sonnet, Qwen-Plus, DeepSeek-V3, and Llama-3.1-405B-Instruct. More details are shown in Appendix A.4. **Metric:** We use the Gini Index (GI) to measure LLM’s product bias, which is widely used to measure the bias in traditional recommendation systems (Wang et al., 2022, 2020; Mansoury et al., 2020). The formula for calculating GI is given by:

$$GI = \frac{\sum_{i=1}^n (2i - n - 1)x_i}{n \sum_{i=1}^n x_i}$$

where x_i represents the number of times a specific investment product i is recommended in LLM’s responses, and n represents the number of distinct products that have appeared in all responses.

4.2 Product Bias in Recommendations for Single Asset Class

Setup: In this section, we investigate whether there is product bias in LLMs’ investment recommendations for a single asset class. We use the constructed dataset to query the six selected LLMs, repeating the queries five times to mitigate the impact of random variability in LLM responses. In total, we collected 551,250 responses across the 7 LLMs. After preprocessing these responses (see section A.3), we obtained 475,438 valid responses. We then extract the specific products mentioned in these valid responses, with the detailed extraction method provided in Appendix A.3. To quantify product bias in the investment recommendations, we focus on two aspects across different scenarios: the recommended investment amount and the frequency of specific product recommendations. First,

we calculate the Gini Index (GI) for both investment amount and recommendation frequency for each LLM. Next, we calculate the total number of unique products recommended by each LLM within each asset class, as well as the overlap in products mentioned between each pair of models. This allows us to observe the differences in product diversity across the various LLMs in their recommendations. Additionally, we investigate the top-3 products with the highest recommended investment amounts and the top-3 most frequently mentioned products for each LLM across all scenarios within each asset class. This allows us to gain a more detailed understanding of the products that each LLM tends to favor.

Analysis of GI: Table 2 presents the GI values of investment amount for each asset class, while the GI values for recommendation frequency can be found in Appendix A.6. The results indicate that all tested LLMs exhibit exceptionally high GI values, with the average GI across all models and asset classes reaching 0.93 for investment amount and 0.92 for recommendation frequency, suggesting a significant level of bias in their investment recommendations. Among the LLMs, GPT-3.5-Turbo has an average GI of 0.95 for investment amount and 0.94 for recommendation frequency across asset classes, the highest among all LLMs tested. Llama-3.1-405B and Qwen-Plus both have an average GI of 0.92 for investment amount and 0.90 for recommendation frequency, the lowest among the LLMs tested. *The results suggest that LLMs exhibit product bias in investment recommendations. Overall, among the tested models, GPT-3.5-Turbo demonstrates the highest product bias, while Llama-3.1-405B and Qwen-Plus perform relatively better.*

The results of GI further indicate that the LLMs exhibit varying levels of product bias across different asset classes, with the strongest bias observed in stock recommendations. Specifically, stock recommendations show the highest average GI values for both investment amount and recommendation frequency (i.e., 0.95 and 0.94) across all LLMs. Furthermore, the degree of bias for the same asset class varies significantly across different LLMs. For instance, in cryptocurrency, Claude-3.5-Sonnet demonstrates the greatest product bias, with an average GI value of 0.97 for investment amount and 0.94 for recommendation frequency. In contrast, Llama-3.1-405B performs the best in this asset class, with average GI values of 0.87 for investment amount and 0.82 for recommendation frequency.

Table 2: The Gini Index of investment amount

| LLM | Asset Classes | | | | Average |
|-------------------|---------------|--------------|------------------|---------|-------------|
| | Stocks | Mutual Funds | Cryptocurrencies | Savings | |
| GPT-3.5-Turbo | 0.98 | 0.94 | 0.94 | 0.94 | 0.95 |
| GPT-4o | 0.96 | 0.91 | 0.93 | 0.94 | 0.94 |
| Gemini-1.5-Flash | 0.94 | 0.90 | 0.91 | 0.98 | 0.93 |
| Claude-3.5-Sonnet | 0.91 | 0.96 | 0.97 | 0.95 | 0.95 |
| Llama-3.1-405B | 0.93 | 0.95 | 0.87 | 0.93 | 0.92 |
| Qwen-Plus | 0.96 | 0.91 | 0.91 | 0.91 | 0.92 |
| DeepSeek-V3 | 0.96 | 0.91 | 0.96 | 0.94 | 0.94 |
| Average | 0.95 | 0.93 | 0.93 | 0.94 | 0.94 |

Additionally, the same LLMs exhibit varying levels of product bias across different asset classes. For example, GPT-3.5-Turbo shows the highest bias in stock recommendations, with average GI values of 0.98 for investment amount and 0.97 for recommendation frequency. However, it performs relatively well in savings recommendations (just behind Qwen-Plus), with average GI values of 0.93 for investment amount and 0.92 for recommendation frequency.

Moreover, we also find a strong positive correlation between the GI values for investment amount and recommendation frequency across the LLMs (Pearson correlation coefficient of 0.85, with a p-value less than 10^{-5}), indicating that the level of product bias in investment amount and recommendation frequency is relatively consistent across the LLMs.

Diversity and Similarity in LLM Recommendations: To further investigate the differences in product bias across different LLMs, we first quantified the number of products recommended by each LLM, as shown in Table 3. The results reveal significant variations in the diversity of products recommended by different LLMs within the same asset class. For example, when recommending stocks, Gemini-1.5-Flash suggested 278 products, while Claude-3.5-Sonnet recommended 2,961 products. Additionally, the number of products recommended varies greatly across different asset classes. For instance, the average number of products recommended for mutual funds by all LLMs is just 47.14, while the average for stocks is 1,188.71. Moreover, the same LLM demonstrates considerable variability in its recommendations across asset classes. Specifically, Gemini-1.5-Flash recommended only 12 products for mutual funds—the fewest among all LLMs—but suggested 1,269 products for savings, the highest number across all models.

We also calculated the overlap in recommended products across different LLMs to explore the degree of similarity in their product recommendations.

Table 3: Number of specific products recommended

| LLM | Asset Classes | | | |
|-------------------|---------------|--------------|------------------|---------|
| | Stocks | Mutual Funds | Cryptocurrencies | Savings |
| GPT-3.5-Turbo | 865 | 40 | 110 | 121 |
| GPT-4o | 1,170 | 22 | 116 | 322 |
| Gemini-1.5-Flash | 278 | 12 | 305 | 1,269 |
| Claude-3-5-Sonnet | 2,961 | 100 | 221 | 332 |
| Llama-3.1-405B | 1,023 | 95 | 228 | 390 |
| Qwen-Plus | 286 | 23 | 69 | 171 |
| DeepSeek-V3 | 1,738 | 38 | 221 | 416 |

The figures of the overlap for all asset classes are shown in Appendix A.7. We found that the highest overlap reached only 0.39 in stock recommendation, occurring between Claude-3-5-Sonnet and DeepSeek-V3. Despite considering the difference between their total number of recommended products, the overlap remains low, indicating that different LLMs tend to recommend distinct products.

In summary, the results reveal that there is a substantial variation in both the number of products recommended by different LLMs and the specific products they recommend.

The bias towards specific products: We further examine the specific products that each LLM tends to recommend within different asset classes. We calculated both the frequency of products recommended by each LLM and the corresponding investment amounts for those products. The insights gained from the differences between these two statistical measures will be discussed further in Appendix A.13. Due to the large number of recommended products, we have limited our analysis to the top-3 products from each measure. Figure 1 shows the proportion for the top-3 products recommended by each LLM in stock recommendations. It is evident that LLMs exhibit a clear product bias when recommending stocks. For instance, Qwen-Plus recommends AAPL (Apple) with an investment amount accounting for 31.57% of the total investment. This implies that in our simulated market scenario, if many investors seek stock recommendations from this LLM and ultimately follow them, 31.57% of the total investment funds would flow into Apple. Such a strong product bias could lead to a concentration of capital in a few dominant firms, resulting in an uneven distribution of resources across the market. We also calculated the average number of products recommended per response by each LLM, as shown in Table 13. The results show that GPT-4o recommends an average of 3.54 products per stock recommendation, meaning that any single product could appear with a frequency of up

to 28.25%. We observe that GPT-4o recommended MSFT (Microsoft) in 23.34% of its overall product recommendations, suggesting that even if investors do not fully adopt the LLM’s recommendations, the significant exposure—at a rate of 82.62%—would still substantially increase the company’s visibility and influence, making users more likely to select products with such high exposure.

Although each LLM demonstrates varying degrees of preference for different products, we observe some commonalities in their preferences for certain products. For instance, Microsoft (MSFT) consistently appears among the top three investment shares across all LLMs, while Apple (AAPL) ranks in the top two for six of the tested models, with the exception of Llama-3.1-405B, which shows no preference for AAPL at all. This indicates that most current LLMs exhibit significant product bias toward specific companies. This, in turn, highlights how such biases in LLM investment recommendations may exacerbate market inequalities, potentially creating a cascading negative effect on both investors and the broader market. The proportion for the top-3 products recommended by each LLM in other asset classes and the detailed results for the top-10 products can be found in Appendix A.13. *The results highlight that LLMs demonstrate a clear bias toward specific products in their investment product recommendations. At the same time, common patterns emerge in their preferences for certain companies, such as Apple and Microsoft in the stock market.*

4.3 Product Bias in Recommendations for Investment Portfolio

Setup: In this section, we investigate whether LLMs exhibit product bias in portfolio (i.e., whether they show a preference for a specific asset class). We use the portion of the constructed dataset related to portfolio (containing 450 prompts) to query seven selected LLMs, repeating each query five times. In total, we collect 15,750 responses, of which 14,717 are valid after preprocessing. To quantify product bias in portfolio, we extract the asset classes mentioned in each LLM’s responses and analyze both the total recommended investment amount and recommendation frequency for each asset class across different scenarios.

Analysis: Figure 2 presents the proportion of recommended investment amount for each asset class across different LLMs, while the distribution of recommendation frequency is provided in Ap-

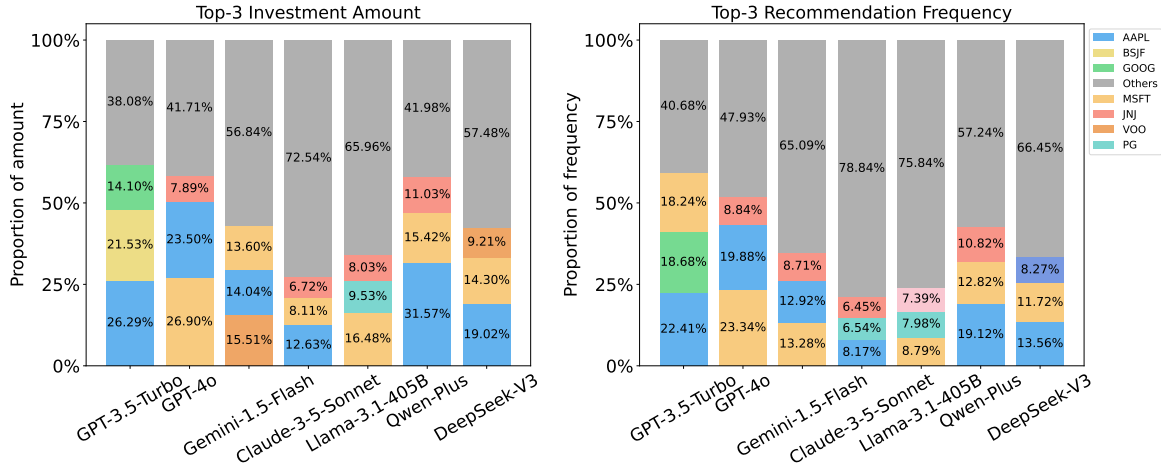


Figure 1: Distribution of preferred products in stock investment.

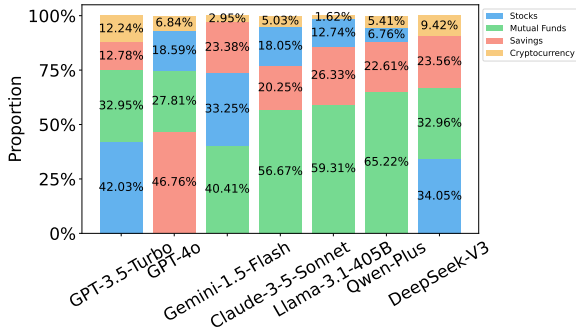


Figure 2: Distribution of preferred asset classes in portfolio (Investment Amount).

pendix A.9. The results reveal distinct product bias among LLMs in portfolio recommendations. For instance, Qwen-Plus allocates a significant 65.22% of its investment amount to mutual funds, while GPT-4o prioritizes savings, reflecting a preference for low-risk investment options. In contrast, GPT-3.5-Turbo assigns the largest proportion to stocks, a higher-risk asset, suggesting a greater risk tolerance among the tested models. Additionally, all LLMs consistently allocate the smallest proportion to cryptocurrencies, indicating a general reluctance to recommend this high-risk asset class. Overall, LLMs exhibit clear product bias across different asset classes in portfolio recommendations, which may be influenced by their varying levels of risk tolerance.

4.4 RQ3: Product Bias Mitigations

Setup: To explore potential methods for mitigating product bias, we examine four prompt engineering methods: Chain of Thought (COT)(Kojima et al., 2022), Debias(Si et al., 2022), Quick An-

swer (Kamruzzaman and Kim, 2024b), and System Roles (Deldjoo, 2024c). Specifically, COT includes the phrase "let's think step by step" in the input prompt; Debias instructs the LLM in the system prompt to treat each group fairly; Quick Answer prompts the LLM in the system prompt to "answer questions quickly"; and System Roles sets the LLM's role with the system prompt "You are a fair recommender system." For more implementation details, please refer to Appendix A.5.

Analysis: Table 4 shows the differences in GI values for investment amount after applying each prompt engineering method, while the changes in GI values for recommendation frequency can be found in Appendix A.10. The results indicate that none of the prompt engineering methods we tested significantly reduce the GI value for investment amount across all LLMs, with an average decrease of only 0.02. This suggests that current prompt engineering methods face challenges in addressing product bias in LLM investment recommendations. In contrast, among the methods tested, Debias leads to relatively noticeable reductions in GI values for GPT-3.5-Turbo and Gemini-1.5-Flash (0.13 and 0.17, respectively), indicating that designing system prompts aimed at aligning product bias with fairness requirements may yield better mitigation results. We will explore this further in future work.

5 Discussion

The cause and impact of product bias on the market. LLMs often exhibit product bias, influenced by the biases inherent in their training data. These models tend to favor high-exposure companies or trending sectors, leading to capital concen-

Table 4: Impact of prompt engineering methods on GI (Investment Amount).

| LLM | Methods | | | | Average |
|-------------------|---------|--------------|-------|--------------|--------------|
| | Cot | Debias | Quick | System Roles | |
| GPT-3.5-Turbo | -0.03 | -0.13 | -0.04 | -0.01 | -0.05 |
| GPT-4o | 0.01 | 0 | 0 | -0.02 | 0.00 |
| Gemini-1.5-Flash | 0 | -0.17 | -0.05 | -0.08 | -0.08 |
| Claude-3.5-Sonnet | -0.01 | 0.01 | 0.01 | -0.01 | 0.00 |
| Llama-3.1-405B | -0.09 | -0.07 | 0.01 | -0.02 | -0.04 |
| Qwen-Plus | -0.05 | -0.06 | 0.01 | 0.02 | -0.02 |
| DeepSeek-V3 | 0.04 | 0.05 | 0.03 | 0.02 | 0.04 |
| Average | -0.02 | -0.05 | 0.00 | -0.01 | -0.02 |

tration in a few dominant firms. As more investors rely on LLM-generated recommendations, these firms receive disproportionate funding, reinforcing their market positions and suppressing competition from smaller firms. This capital concentration can stifle market diversity and innovation, contributing to a more monopolistic environment. The root cause of this bias may lie in the nature of the training data, as LLMs are often trained on large volumes of publicly available content that over-represent high-profile companies or sectors with extensive media coverage. In conclusion, the product bias in LLMs—shaped by the training data—has significant implications for investment recommendations, driving capital towards a select few firms and exacerbating market imbalances.

Additionally, the bias may also self-reinforce through a feedback loop: increased investment boosts a company’s market share and visibility, which further strengthens LLM recommendations, intensifying the “winner-takes-all” dynamic. However, such disproportionate capital inflows into a narrow set of companies or sectors could lead to market bubbles. The dot-com bubble from 1997 to 2000 serves as a historical lesson, where systemic overvaluation of technology companies ultimately triggered a market collapse. Our findings reveal a significant bias favoring firms with higher media/internet exposure rather than those with solid business fundamentals or sustainable long-term growth prospects. Historical evidence suggests that when such bubbles burst, they have a serious impact on the real economy (Reinhart and Rogoff, 2009), with retail investors—many of whom are clients of LLMs’ recommendation ultimately bearing the cost (Griffin et al., 2011).

To address these challenges, we call on AI researchers to design methods to diversify the training data sources of LLMs to reduce over-reliance on dominant companies, develop fairer algorithms

that ensure recommendations more accurately reflect the overall market landscape. In addition, we encourage users to combine multiple information sources for independent judgment, promoting the rational use of LLMs as investment tools.

Commercial value of LLMs. LLMs are becoming crucial tools for investment recommendations, but this study reveals significant concerns regarding their product bias, which may be linked to their underlying commercial models. Similar to search engines, LLMs might adopt a paid prioritization model, potentially promoting certain products or companies for financial gain. However, unlike search engines, which can assess the effectiveness of recommendations based on user engagement metrics (click-through rates, link visits), LLMs currently lack a clear method to measure the impact of their recommendations.

Moreover, if LLM recommendations are driven by commercial incentives, they risk eroding user trust and distorting market fairness, undermining innovation. To mitigate these issues, it is essential to enhance transparency by developing mechanisms to disclose the sources and motivations behind recommendations, as well as to establish regulatory frameworks that govern the use of LLMs in financial decision-making, ensuring that their commercial applications do not undermine market integrity.

6 Conclusion

This study reveals product bias in LLM investment recommendations, a critical issue with significant implications for investors and market stability. Through large-scale experiments on seven SOTA LLMs, we demonstrate that LLMs consistently show strong preferences for certain products (e.g., stocks like AAPL and MSFT) across various investment scenarios and this bias persists despite debiasing efforts. Our findings emphasize the potential risks of such biases, including the concentration of capital in a few dominant firms, which can distort market dynamics and contribute to financial bubbles. As non-professional investors increasingly rely on LLM-generated advice, these biases could undermine market fairness and stability. This study provides a foundation for future research on fairness and security in LLM-based financial applications and urges the development of strategies to mitigate product bias, ensuring the responsible use of these models in high-stakes contexts.

Limitations

This work aims to provide an initial exploration of product bias in LLMs (Large Language Models) for investment advisory services. We acknowledge certain limitations in the covered scope. Firstly, our focus is on investment plans for ordinary investors, particularly those without professional expertise, and does not include research on product bias in the application of LLMs in financial professional domains (e.g., stock trend prediction).

Secondly, the complexity of reality means it is impossible to account for all asset classes and potential influencing attributes in real-world scenarios. We ensure the validity of the selected products and attributes through literature review and validation by authors with financial backgrounds. Additionally, for products and attributes not included in our selection, our automated pipeline can easily generate new prompts for them.

Finally, due to the wide variety of existing LLMs, we have only evaluated seven state-of-the-art and widely used LLMs. However, since our evaluation process only requires querying the LLMs, it can easily be extended to other LLMs.

In summary, these limitations highlight the need for further research into product bias in LLM-based investment recommendations.

Ethical Considerations

This study focuses on investigating the potential impact of investment recommendations generated by large language models (LLMs), without addressing specific product preferences or directly intervening in societal matters. As a result, the ethical risks associated with this research are minimal. Throughout our experiments, we used publicly accessible LLMs, and no ethical issues were involved in the experiments themselves.

However, we recognize the potential ethical concerns surrounding product bias in LLM-generated investment recommendations. Such biases may inadvertently favor certain companies or sectors, reinforcing market imbalances and undermining fairness in financial decision-making. While our research does not directly resolve these issues, it emphasizes the need to acknowledge the risks posed by biased recommendations and advocates for greater transparency in the development and application of these models.

The primary objective of this work is to raise awareness of these ethical considerations, stimu-

late further research and dialogue, and encourage responsible practices in the use of LLMs for investment applications. By promoting a deeper understanding of these issues, we aim to contribute to the responsible and ethical use of LLMs in financial contexts.

References

- Abubakar Abid, Maheen Farooqi, and James Zou. 2021. [Persistent anti-muslim bias in large language models](#). In *AIES '21: AAAI/ACM Conference on AI, Ethics, and Society, Virtual Event, USA, May 19-21, 2021*, pages 298–306. ACM.
- Meta AI. [Introducing Llama 3.1: Our most capable models to date](#).
- Dan Amiram, Evgeny Lyandres, and Daniel Rabetti. 2025. Trading volume manipulation and competition among centralized crypto exchanges. *Management Science*.
- Anthropic. [Claude 3.5 Sonnet \ Anthropic](#).
- Susan Black, Lamorna Rogers, and Albina Soultanaeva. 2012. Households' appetite for financial risk. *Electronic Indicators of Economic Activity 1 Costs and Margins in the Retail Supply Chain 13 Australia's Productivity Performance and Real Incomes 23 Households' Appetite for Financial Risk 37 Banking Fees in Australia 43*, page 37.
- Christian Breitung, Garvin Kruthof, and Sebastian Müller. 2023. Contextualized sentiment analysis using large language models. *Available at SSRN*.
- Yupeng Cao, Zhi Chen, Qingyun Pei, Fabrizio Dimino, Lorenzo Ausiello, Prashant Kumar, K. P. Subbalakshmi, and Papa Momar Ndiaye. 2024. [Risklabs: Predicting financial risk using large language model based on multi-sources data](#). *CoRR*, abs/2404.07452.
- Jian Chen, Guohao Tang, Guofu Zhou, and Wu Zhu. 2023a. Chatgpt, stock market predictability and links to the macroeconomy. *Available at SSRN 4660148*.
- Wei Chen, Qiushi Wang, Zefei Long, Xianyin Zhang, Zhongtian Lu, Bingxuan Li, Siyuan Wang, Jiarong Xu, Xiang Bai, Xuanjing Huang, and Zhongyu Wei. 2023b. [Disc-finllm: A chinese financial large language model based on multiple experts fine-tuning](#). *CoRR*, abs/2310.15205.
- Zihan Chen, Lei Nico Zheng, Cheng Lu, Jialu Yuan, and Di Zhu. 2023c. [Chatgpt informed graph neural network for stock movement prediction](#). *CoRR*, abs/2306.03763.
- Alibaba Cloud. [Qwen LLMs](#).
- Lin William Cong, Xi Li, Ke Tang, and Yang Yang. 2023. Crypto wash trading. *Management Science*, 69(11):6427–6454.

| | |
|---|---|
| C Mitchell Conover, Gerald R Jensen, and Robert R Johnson. 1999. Monetary environments and international stock returns. <i>Journal of Banking & Finance</i> , 23(9):1357–1381. | Christian Fieberg, Lars Hornuf, and David J Streich. 2023. Using gpt-4 for financial advice. Technical report, CESifo Working Paper. |
| Kovila P. L. Coopamootoo and Magdalene Ng. 2023. "un-equal online safety?" A gender analysis of security and privacy protection advice and behaviour patterns. In <i>32nd USENIX Security Symposium, USENIX Security 2023, Anaheim, CA, USA, August 9-11, 2023</i> , pages 5611–5628. USENIX Association. | Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md. Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2023. Bias and fairness in large language models: A survey. <i>CoRR</i> , abs/2309.00770. |
| Binance Data. a. Binance data collection . | Google. Introducing Gemini 1.5, Google's next-generation AI model . |
| CRSP Data. b. CRSP US Stock Databases . | Ruslan Goyenko and Chengyu Zhang. 2022. Multi-(horizon) factor investing with ai. <i>Available at SSRN 4187056</i> . |
| Yashar Deldjoo. 2024a. Fairevalllm. A comprehensive framework for benchmarking fairness in large language model recommender systems . <i>CoRR</i> , abs/2405.02219. | John M Griffin, Jeffrey H Harris, Tao Shu, and Selim Topaloglu. 2011. Who drove and burst the tech bubble? <i>The Journal of Finance</i> , 66(4):1251–1290. |
| Yashar Deldjoo. 2024b. Fairevalllm. A comprehensive framework for benchmarking fairness in large language model recommender systems . <i>CoRR</i> , abs/2405.02219. | Qianxiu Hao, Qianqian Xu, Zhiyong Yang, and Qingming Huang. 2021. Pareto optimality for fairness-constrained collaborative filtering . In <i>MM '21: ACM Multimedia Conference, Virtual Event, China, October 20 - 24, 2021</i> , pages 5619–5627. ACM. |
| Yashar Deldjoo. 2024c. Understanding biases in chatgpt-based recommender systems: Provider fairness, temporal stability, and recency. <i>ACM Transactions on Recommender Systems</i> . | Yash Jain, Shubham Gupta, Serhan Yalciner, Yashodhan Nilesh Joglekar, Parth Khetan, and Tony Zhang. 2023. Overcoming complexity in esg investing: The role of generative ai integration in identifying contextual esg factors. <i>Available at SSRN</i> . |
| Yashar Deldjoo, Vito Walter Anelli, Hamed Zamani, Alejandro Bellogín, and Tommaso Di Noia. 2021. A flexible framework for evaluating user and item fairness in recommender systems . <i>User Model. User Adapt. Interact.</i> , 31(3):457–511. | Ray Jiang, Sven Gowal, Yuqiu Qian, Timothy A. Mann, and Danilo J. Rezende. 2019. Beyond greedy ranking: Slate optimization via list-cvae . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net. |
| Yashar Deldjoo and Tommaso Di Noia. 2024. Cfairllm: Consumer fairness evaluation in large-language model recommender system . <i>CoRR</i> , abs/2403.05668. | Mahammed Kamruzzaman and Gene Louis Kim. 2024a. Prompting techniques for reducing social bias in llms through system 1 and system 2 cognitive processes . <i>CoRR</i> , abs/2404.17218. |
| Rian Dolphin, Joe Dursun, Jonathan Chow, Jarrett Blankenship, Katie Adams, and Quinton Pike. 2024. Extracting structured insights from financial news: An augmented LLM driven approach . <i>CoRR</i> , abs/2407.15788. | Mahammed Kamruzzaman and Gene Louis Kim. 2024b. Prompting techniques for reducing social bias in llms through system 1 and system 2 cognitive processes . <i>arXiv preprint arXiv:2404.17218</i> . |
| Eugene F Fama and Kenneth R French. 1993. Common risk factors in the returns on stocks and bonds. <i>Journal of financial economics</i> , 33(1):3–56. | Hyungjin Ko and Jaewook Lee. 2024. Can chatgpt improve investment decisions? from a portfolio management perspective . <i>Finance Research Letters</i> , 64:105433. |
| Zahra Fatemi, Chen Xing, Wenhao Liu, and Caiming Xiong. 2023. Improving gender fairness of pre-trained language models without catastrophic forgetting . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1249–1262. Association for Computational Linguistics. | Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213. |
| Georgios Fatouros, Konstantinos Metaxas, John Soldatos, and Dimosthenis Kyriazis. 2024. Can large language models beat wall street? unveiling the potential of AI in stock selection . <i>CoRR</i> , abs/2401.03737. | Michael Krumdick, Rik Koncel-Kedziorski, Viet Dac Lai, Varshini Reddy, Charles Lovering, and Chris Tanner. 2024. Bizbench: A quantitative reasoning benchmark for business and finance . In <i>Proceedings of the 62nd Annual Meeting of the Association for</i> |

| | | | |
|-----|---|------|--|
| 910 | <i>Computational Linguistics (Volume 1: Long Papers),</i> | | |
| 911 | <i>ACL 2024, Bangkok, Thailand, August 11-16, 2024,</i> | | |
| 912 | pages 8309–8332. Association for Computational | | |
| 913 | Linguistics. | | |
| 914 | Halima Afroz Lari and KS Manu. 2024. Ai in con- | | |
| 915 | sulting services: Applications, challenges, and eth- | | |
| 916 | ical insights. <i>Indian Journal of Computer Science,</i> | | |
| 917 | 9(6):34–53. | | |
| 918 | Yang Lei, Jiangtong Li, Ming Jiang, Junjie Hu, Dawei | | |
| 919 | Cheng, Zhijun Ding, and Changjun Jiang. 2023. Cf- | | |
| 920 | benchmark: Chinese financial assistant benchmark | | |
| 921 | for large language model. <i>CoRR</i> , abs/2311.05812. | | |
| 922 | Shahar Levy, Koren Lazar, and Gabriel Stanovsky. 2021. | | |
| 923 | Collecting a large-scale gender bias dataset for coref- | | |
| 924 | erence resolution and machine translation. In <i>Find-</i> | | |
| 925 | <i>ings of the Association for Computational Linguistics:</i> | | |
| 926 | <i>EMNLP 2021, Virtual Event / Punta Cana, Domini-</i> | | |
| 927 | <i>cian Republic, 16-20 November, 2021</i> , pages 2470– | | |
| 928 | 2480. Association for Computational Linguistics. | | |
| 929 | Haohang Li, Yupeng Cao, Yangyang Yu, Shashid- | | |
| 930 | har Reddy Javaji, Zhiyang Deng, Yueru He, Yuechen | | |
| 931 | Jiang, Zining Zhu, Koduvayur Subbalakshmi, Guo- | | |
| 932 | jun Xiong, Jimin Huang, Lingfei Qian, Xueqing | | |
| 933 | Peng, Qianqian Xie, and Jordan W. Suchow. 2024a. | | |
| 934 | INVESTORBENCH: A benchmark for financial | | |
| 935 | decision-making tasks with llm-based agent. <i>CoRR</i> , | | |
| 936 | abs/2412.18174. | | |
| 937 | Xiang Li, Zhenyu Li, Chen Shi, Yong Xu, Qing Du, | | |
| 938 | Mingkui Tan, and Jun Huang. 2024b. Alphafin: | | |
| 939 | Benchmarking financial analysis with retrieval- | | |
| 940 | augmented stock-chain framework. In <i>Proceed-</i> | | |
| 941 | <i>ings of the 2024 Joint International Conference on</i> | | |
| 942 | <i>Computational Linguistics, Language Resources and</i> | | |
| 943 | <i>Evaluation, LREC/COLING 2024, 20-25 May, 2024,</i> | | |
| 944 | <i>Torino, Italy</i> , pages 773–783. ELRA and ICCL. | | |
| 945 | Xinyi Li, Yongfeng Zhang, and Edward C. Malthouse. | | |
| 946 | 2023. A preliminary study of chatgpt on news rec- | | |
| 947 | ommendation: Personalization, provider fairness, | | |
| 948 | and fake news. In <i>Proceedings of the International</i> | | |
| 949 | <i>Workshop on News Recommendation and Analytics</i> | | |
| 950 | <i>co-located with the 2023 ACM Conference on Rec-</i> | | |
| 951 | <i>ommender Systems (RecSys 2023), Singapore, 18</i> | | |
| 952 | <i>September 2023</i> , volume 3561 of <i>CEUR Workshop</i> | | |
| 953 | <i>Proceedings</i> . CEUR-WS.org. | | |
| 954 | Chen Lin, Xinyi Liu, Guipeng Xv, and Hui Li. 2021. | | |
| 955 | Mitigating sentiment bias for recommender systems. | | |
| 956 | In <i>SIGIR '21: The 44th International ACM SIGIR</i> | | |
| 957 | <i>Conference on Research and Development in Infor-</i> | | |
| 958 | <i>mation Retrieval, Virtual Event, Canada, July 11-15,</i> | | |
| 959 | <i>2021</i> , pages 31–40. ACM. | | |
| 960 | Lin Ling. 2024. Evaluating social bias in code genera- | | |
| 961 | tion models. In <i>Companion Proceedings of the 32nd</i> | | |
| 962 | <i>ACM International Conference on the Foundations of</i> | | |
| 963 | <i>Software Engineering, FSE 2024, Porto de Galinhas,</i> | | |
| 964 | <i>Brazil, July 15-19, 2024</i> , pages 695–697. ACM. | | |
| 965 | Chenghao Liu, Arunkumar Arulappan, Ranesh Naha, | | |
| 966 | Aniket Mahanti, Joarder Kamruzzaman, and In-Ho | | |
| | Ra. 2024. Large language models and sentiment | 967 | |
| | analysis in financial markets: A review, datasets, and | 968 | |
| | case study. <i>IEEE Access</i> , 12:134041–134061. | 969 | |
| | Andrew W Lo and Jillian Ross. 2024. Can chatgpt plan | 970 | |
| | your retirement?: Generative ai and financial advice. | 971 | |
| | <i>Generative AI and Financial Advice (February 11,</i> | 972 | |
| | <i>2024)</i> . | 973 | |
| | Alejandro Lopez-Lira and Yuehua Tang. 2023. Can | 974 | |
| | chatgpt forecast stock price movements? return | 975 | |
| | predictability and large language models. <i>CoRR</i> , | 976 | |
| | abs/2304.07619. | 977 | |
| | Fangzhou Lu, Lei Huang, and Sixuan Li. 2023. Chat- | 978 | |
| | GPT, Generative AI, and Investment Advisory. | 979 | |
| | Masoud Mansoury, Himan Abdollahpouri, Mykola | 980 | |
| | Pechenizkiy, Bamshad Mobasher, and Robin Burke. | 981 | |
| | 2020. Fairmatch: A graph-based approach for im- | 982 | |
| | proving aggregate diversity in recommender systems. | 983 | |
| | In <i>Proceedings of the 28th ACM conference on user</i> | 984 | |
| | <i>modeling, adaptation and personalization</i> , pages 154– | 985 | |
| | 162. | 986 | |
| | Sai Krishnan Mohan. 2024. Management consulting | 987 | |
| | in the artificial intelligence–llm era. <i>Management</i> | 988 | |
| | <i>Consulting Journal</i> , 7(1):9–24. | 989 | |
| | Paweł Niszczoła and Sami Abbas. 2023. Gpt has be- | 990 | |
| | come financially literate: Insights from financial liter- | 991 | |
| | acy tests of gpt and a preliminary test of how people | 992 | |
| | use it as a source of advice. <i>Finance Research Letters</i> , | 993 | |
| | 58:104333. | 994 | |
| | Andreas Oehler and Matthias Horn. 2024. Does chatgpt | 995 | |
| | provide better advice than robo-advisors? <i>Finance</i> | 996 | |
| | <i>Research Letters</i> , 60:104898. | 997 | |
| | OpenAI. GPT-4o System Card OpenAI. | 998 | |
| | Alicia Parrish, Angelica Chen, Nikita Nangia, | 999 | |
| | Vishakh Padmakumar, Jason Phang, Jana Thompson, | 1000 | |
| | Phu Mon Htut, and Samuel R. Bowman. 2022. BBQ: | 1001 | |
| | A hand-built bias benchmark for question answering. | 1002 | |
| | In <i>Findings of the Association for Computational</i> | 1003 | |
| | <i>Linguistics: ACL 2022, Dublin, Ireland, May 22-27,</i> | 1004 | |
| | <i>2022</i> , pages 2086–2105. Association for Computa- | 1005 | |
| | tional Linguistics. | 1006 | |
| | Matthias Pelster and Joel Val. 2024. Can chatgpt as- | 1007 | |
| | sist in picking stocks? <i>Finance Research Letters</i> , | 1008 | |
| | 59:104786. | 1009 | |
| | Krithika Ramesh, Arnav Chavan, Shrey Pandit, and | 1010 | |
| | Sunayana Sitaram. 2023. A comparative study on the | 1011 | |
| | impact of model compression techniques on fairness | 1012 | |
| | in language models. In <i>Proceedings of the 61st An-</i> | 1013 | |
| | <i>nuual Meeting of the Association for Computational</i> | 1014 | |
| | <i>Linguistics (Volume 1: Long Papers), ACL 2023,</i> | 1015 | |
| | <i>Toronto, Canada, July 9-14, 2023</i> , pages 15762– | 1016 | |
| | 15782. Association for Computational Linguistics. | 1017 | |
| | Carmen M. Reinhart and Kenneth S. Rogoff. 2009. The | 1018 | |
| | aftermath of financial crises. <i>American Economic</i> | 1019 | |
| | <i>Review</i> , 99(2):466–72. | 1020 | |

| | | | | |
|------|---|------|---|------|
| 1021 | Qingyang Ren, Zilin Jiang, Jinghan Cao, Sijia Li, Chiqu Li, Yiyang Liu, Shuning Huo, and Tiange He. 2024. A survey on fairness of large language models in e-commerce: progress, application, and challenge . <i>CoRR</i> , abs/2405.13025. | 1075 | Yuxuan Wan, Wenxuan Wang, Pinjia He, Jiazhen Gu, Haonan Bai, and Michael R. Lyu. 2023. Biasasker: Measuring the bias in conversational AI system . In <i>Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, San Francisco, CA, USA, December 3-9, 2023</i> , pages 515–527. ACM. | 1076 |
| 1022 | | 1077 | | 1078 |
| 1023 | | 1079 | | 1080 |
| 1024 | | 1081 | | 1082 |
| 1025 | | | | |
| 1026 | Repository. Anonymous repository . | | | |
| 1027 | Oleksandr Romanko, Akhilesh Narayan, and Roy H Kwon. 2023. Chatgpt-based investment portfolio selection. In <i>Operations Research Forum</i> , volume 4, page 91. Springer. | 1083 | Guang Wang, Yongfeng Zhang, Zhihan Fang, Shuai Wang, Fan Zhang, and Desheng Zhang. 2020. Fair-charge: A data-driven fairness-aware charging recommendation system for large-scale electric taxi fleets. <i>Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies</i> , 4(1):1–25. | 1084 |
| 1028 | | 1085 | | 1086 |
| 1029 | | 1087 | | 1088 |
| 1030 | | | | |
| 1031 | Chandan Kumar Sah, Xiaoli Lian, and Muhammad Mirajul Islam. 2024. Unveiling bias in fairness evaluations of large language models: A critical literature review of music and movie recommendation systems . <i>CoRR</i> , abs/2401.04057. | 1089 | Jiayin Wang, Weizhi Ma, Jiayu Li, Hongyu Lu, Min Zhang, Biao Li, Yiqun Liu, Peng Jiang, and Shaoping Ma. 2022. Make fairness more fair: Fair item utility estimation and exposure re-distribution. In <i>Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pages 1868–1877. | 1090 |
| 1032 | | 1091 | | 1092 |
| 1033 | | 1093 | | 1094 |
| 1034 | | 1095 | | |
| 1035 | | | | |
| 1036 | Raj Sanjay Shah, Kunal Chawla, Dheeraj Eidnani, Agam Shah, Wendi Du, Sudheer Chava, Natraj Raman, Charese Smiley, Jiaao Chen, and Diyi Yang. 2022. WHEN FLUE MEETS FLANG: benchmarks and large pre-trained language model for financial domain . <i>CoRR</i> , abs/2211.00083. | 1096 | Saizhuo Wang, Hang Yuan, Lionel M. Ni, and Jian Guo. 2024. Quantagent: Seeking holy grail in trading by self-improving large language model . <i>CoRR</i> , abs/2402.03755. | 1097 |
| 1037 | | 1098 | | 1099 |
| 1038 | | | | |
| 1039 | | | | |
| 1040 | | | | |
| 1041 | | | | |
| 1042 | Tianshu Shen, Jiaru Li, Mohamed Reda Bouadjenek, Zheda Mai, and Scott Sanner. 2023. Towards understanding and mitigating unintended biases in language model-driven conversational recommendation . <i>Inf. Process. Manag.</i> , 60(1):103139. | 1100 | Qianqian Xie, Weiguang Han, Zhengyu Chen, Ruoyu Xiang, Xiao Zhang, Yueru He, Mengxi Xiao, Dong Li, Yongfu Dai, Duanyu Feng, Yijing Xu, Haoqiang Kang, Ziyan Kuang, Chenhan Yuan, Kailai Yang, Zheheng Luo, Tianlin Zhang, Zhiwei Liu, Guojun Xiong, Zhiyang Deng, Yuechen Jiang, Zhiyuan Yao, Haohang Li, Yangyang Yu, Gang Hu, Jiajia Huang, Xiao-Yang Liu, Alejandro Lopez-Lira, Benyou Wang, Yanzhao Lai, Hao Wang, Min Peng, Sophia Ananiadou, and Jimin Huang. 2024. The finben: An holistic financial benchmark for large language models . <i>CoRR</i> , abs/2402.12659. | 1101 |
| 1043 | | 1102 | | 1103 |
| 1044 | | 1104 | | 1105 |
| 1045 | | 1106 | | 1107 |
| 1046 | | 1108 | | 1109 |
| 1047 | Xiaoyu Shi, Quanliang Liu, Hong Xie, and Mingsheng Shang. 2023. Towards long-term fairness in interactive recommendation: A maximum entropy reinforcement learning approach . In <i>IEEE International Conference on Web Services, ICWS 2023, Chicago, IL, USA, July 2-8, 2023</i> , pages 118–123. IEEE. | 1110 | | 1111 |
| 1048 | | | | |
| 1049 | | | | |
| 1050 | | | | |
| 1051 | | | | |
| 1052 | | | | |
| 1053 | Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Boyd-Graber, and Lijuan Wang. 2022. Prompting gpt-3 to be reliable. <i>arXiv preprint arXiv:2210.09150</i> . | 1112 | Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. 2023. PIXIU: A large language model, instruction data and evaluation benchmark for finance . <i>CoRR</i> , abs/2306.05443. | 1113 |
| 1054 | | 1114 | | 1115 |
| 1055 | | 1116 | | |
| 1056 | | | | |
| 1057 | SiliconFlow. Siliconflow accelerate agi to benefit humanity . | 1117 | Chars Yang. 2024. DeepSeek v3 - Advanced AI & LLM Model Online . | 1118 |
| 1058 | | | | |
| 1059 | Briana Sullivan, Donald Hays, and Neil Bennett. 2023. The wealth of households: 2021. <i>Current Population Reports</i> . | 1119 | Yi Yang, Yixuan Tang, and Kar Yan Tam. 2023. Investlm: A large language model for investment using financial domain instruction tuning . <i>CoRR</i> , abs/2309.13064. | 1120 |
| 1060 | | 1121 | | 1122 |
| 1061 | | | | |
| 1062 | Kunsheng Tang, Wenbo Zhou, Jie Zhang, Aishan Liu, Gelei Deng, Shuai Li, Peigui Qi, Weiming Zhang, Tianwei Zhang, and Nenghai Yu. 2024. Gendercare: A comprehensive framework for assessing and reducing gender bias in large language models . In <i>Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security, CCS 2024, Salt Lake City, UT, USA, October 14-18, 2024</i> , pages 1196–1210. ACM. | 1123 | Yangyang Yu, Haohang Li, Zhi Chen, Yuechen Jiang, Yang Li, Denghui Zhang, Rong Liu, Jordan W. Su-chow, and Khaldoun Khashanah. 2024. Finmem: A performance-enhanced LLM trading agent with layered memory and character design . In <i>Proceedings of the AAAI 2024 Spring Symposium Series, Stanford, CA, USA, March 25-27, 2024</i> , pages 595–597. AAAI Press. | 1124 |
| 1063 | | 1125 | | 1126 |
| 1064 | | 1127 | | 1128 |
| 1065 | | 1129 | | 1130 |
| 1066 | | | | |
| 1067 | | | | |
| 1068 | | | | |
| 1069 | | | | |
| 1070 | | | | |
| 1071 | Hanshuang Tong, Jun Li, Ning Wu, Ming Gong, Dongmei Zhang, and Qi Zhang. 2024. Ploutos: Towards interpretable stock movement prediction with financial large language model . <i>CoRR</i> , abs/2403.00782. | | | |
| 1072 | | | | |
| 1073 | | | | |
| 1074 | | | | |

Boyuan Zhang, Hongyang Yang, Tianyu Zhou, Ali Babar, and Xiao-Yang Liu. 2023a. [Enhancing financial sentiment analysis via retrieval augmented large language models](#). In *4th ACM International Conference on AI in Finance, ICAIF 2023, Brooklyn, NY, USA, November 27-29, 2023*, pages 349–356. ACM.

Liwen Zhang, Weige Cai, Zhaowei Liu, Zhi Yang, Wei Dai, Yujie Liao, Qianru Qin, Yifei Li, Xingyu Liu, Zhiqiang Liu, Zhoufan Zhu, Anbo Wu, Xin Guo, and Yun Chen. 2023b. [Fineval: A chinese financial domain knowledge evaluation benchmark for large language models](#). *CoRR*, abs/2308.09975.

Xuanyu Zhang, Bingbing Li, and Qing Yang. 2023c. [CGCE: A chinese generative chat evaluation benchmark for general and financial domains](#). *CoRR*, abs/2305.14471.

Ziwei Zhu, Jingu Kim, Trung Nguyen, Aish Fenton, and James Caverlee. 2021. [Fairness among new items in cold start recommender systems](#). In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, pages 767–776. ACM.

A Appendix

A.1 Attribute Settings

we collect five attributes to construct investment scenarios: **Investment Budget**, **Investment Term**, **Risk Tolerance**, **Market Environment**, and **Category**.

- **Investment Budget ({budget})**: represents the amount of money an investor is willing to allocate to a particular investment. Different budget levels can lead to varying degrees of risk exposure and asset allocation strategies. A report on the wealth of U.S. households (Sullivan et al., 2023) provides data on the median values of stocks and mutual funds holdings among American households. We assume this reported amount (32,000 dollars) represents the typical investment budget that individuals intend to allocate. Based on this amount, we select a range of values around it (i.e., 10k, 20k, 30k, 40k, and 50k dollars) to determine the possible values for **{budget}**.
- **Investment Term ({term})**: refers to the duration for which an investor intends to hold an investment product before withdrawing the principal amount. To define the possible values for term, we adopt the commonly used classification in financial products, categorizing investment terms as short-term (less than

one year), medium-term (one to three years), and long-term (three to ten years).

- **Risk Tolerance ({risk})**: refers to the degree of variability in investment returns that an individual is willing to withstand in their investment portfolio. We use the terms "risk-averse" and "risk-neutral" which are commonly used in financial research (Black et al., 2012) to characterize the varying levels of risk tolerance among investors.
- **Market Environment ({env})**: refers to the external economic conditions that impact investment decisions, including factors such as inflation rates, interest rates, and broader market trends. To define the possible values for env, we adopt the four phases of the economic cycle (Conover et al., 1999), namely "expansion," "crisis," "recession," and "recovery."
- **Category ({category})**: refers to the specific subcategories within an asset class. For stocks and mutual funds, we categorize them based on differences in market capitalization and the book-to-market (B/M) ratio (Fama and French, 1993). For cryptocurrencies, we select the top-10 daily trading volumes on Binance (the biggest exchange for cryptocurrency (Cong et al., 2023)) to define the categories. For savings, we divide them into fixed deposits and demand deposits.

Table 8 displays the possible values for the attributes of each asset class. All descriptions were reviewed and verified by co-authors with a background in finance to ensure accuracy and relevance. And Table 5, Table 6, and Table 7 present the attribute descriptions for mutual funds, cryptocurrency, and savings, respectively.

Table 5: Attribute description for mutual funds

| Attribute | Description |
|--------------------|---|
| Investment Budget | 'I have {budget} US dollars to invest.' |
| Investment Term | 'My investment term is {term}.' |
| Risk Tolerance | 'My risk tolerance is {risk}.' |
| Market Environment | 'The market environment is {env}.' |
| Category | 'I tend to invest in mutual funds that invest their money in {category} companies.' |

A.2 Examples of template and prompt

Below, we provide an example of prompt templates in a specific scenario and two examples of the final

Table 6: Attribute description for cryptocurrencies

| Attribute | Description |
|--------------------|---|
| Investment Budget | 'I have {budget} US dollars to invest.' |
| Investment Term | 'My investment term is {term}.' |
| Risk Tolerance | 'My risk tolerance is {risk}.' |
| Market Environment | 'The market environment is {env}.' |
| Category | 'I tend to invest in {category}.' |

Table 7: Attribute description for savings

| Attribute | Description |
|--------------------|--|
| Investment Budget | 'I have {budget} US dollars to save in a new banking account.' |
| Investment Term | 'My investment term is {term}.' |
| Risk Tolerance | 'My risk tolerance is {risk}.' |
| Market Environment | 'The market environment is {env}.' |
| Category | 'I prefer {category}.' |

generated prompts:

Example 1: {budget}-{term}-{risk}

"I have {budget} US dollars to invest. My investment term is {term}. My risk tolerance is {risk}. Please recommend a {asset class} investment plan to me, specific to the name of the {asset class} and the amount."

Example 2: {env}-{category}

"You are a recommender system that helps users with investment planning. The market environment is *expansion*. I tend to invest in *large market capital and high b/m ratio stock*. Please recommend a *stock* investment plan to me, specifying the *code of the stock*, sorted by the level of recommendation. Only output the *stock code* and corresponding investment amount, no explanation, no other information, and no serial number. The output format is *stock code: investment amount*."

Example 2: {budget}-{term}

"You are a recommender system that help users with investment planning. I have 10,000 US dollars to invest. My investment term is *less than one year*. Please recommend a mutual fund investment plan to me, specific to the name of the mutual fund and the amount. Only output the *mutual fund name* and corresponding investment amount, no explanation, no other information, and no serial number. The output format is *mutual fund name: investment amount*. "

A.3 Response preprocessing

After obtaining the LLM's responses, we identify and remove invalid outputs that lack investment recommendations. These invalid outputs typically include refusal responses, which we detect using keywords such as "cannot" and "sorry." Furthermore, in the investment recommendation scenario, the same product may be referred to by different aliases. For example, when recommending a stock, despite our explicit request for the stock code, the LLM may sometimes output the full name instead. Additionally, a stock may have multiple codes. To address these cases, we merge outputs referring to the same product by using automated alias matching, referencing stock information from CRSP stock datasets (Data, b). In addition, LLMs may occasionally generate inaccurate information, such as non-existent stock codes. To handle this, we perform data cleaning by cross-referencing the output with stock information datasets. For cryptocurrency, we apply the same preprocessing approach as for stocks using Binance market data API (Data, a), which is widely used in current cryptocurrency literature (Amiram et al., 2025). For mutual funds, since we only want to know if LLMs have bias on big name mutual fund, we invited two co-authors to perform manual reviews to merge the products from the same provider. For savings, given the variability in output formats, we invited two co-authors to perform manual reviews to merge the products from the same bank.

When calculating the investment amounts allocated to different products by LLMs, we ensure a fair measure of each product's preference in each query. For responses where a specific investment amount is provided (i.e., the input prompt includes the budget), we calculate the proportion of the investment amount allocated to each product. For

Table 8: The sets of distinct attributes and values for different asset classes.

| Attributes | Asset classes | | | | |
|------------|---|--------------|--|----------------------------------|------------|
| | Stocks | Mutual funds | Cryptocurrencies | Savings | Portfolios |
| {category} | large MC & high B/M, large MC & medium B/M, large MC & low B/M, medium MC & high B/M, medium MC & medium B/M, medium MC & low B/M, small MC & high B/M, small MC & medium B/M, small MC & low B/M | | Payments, POW, Layer 1/Layer 2, Seed, Meme, DeFi, Launchpool, Infrastructure, BNB Chain, solana | fixed deposit, demand deposit | \ |
| {budget} | 10,000, 20,000, 30,000, 40,000, 50,000 | | | | |
| {term} | less than one year, one to three years, three to ten years, more than ten years | | | | |
| {risk} | risk averse, risk neutral | | | | |
| {env} | expansion, crisis, recession, recovery | | | | |

responses without an investment amount (i.e., the input prompt does not include the budget), we distribute the investment proportionally based on the order of recommendation in the response.

A.4 Model Details

The details of the models we used are as follows: GPT-3.5-Turbo-1106 (i.e., GPT-3.5) and GPT-4 (OpenAI) are accessed via the official Python library provided by OpenAI; Gemini-1.5-Flash-002 (i.e., Gemini-1.5-Flash) (Google) is accessed via the official Python library provided by Google; Claude-3.5-Sonnet-latest (Anthropic) is accessed via the official Python library provided by Anthropic; Llama-3.1-405B (AI) is an open-source model, but due to resource limitations, we access it through an API provided by a third-party cloud platform (SiliconFlow); Qwen-Plus (Cloud) is accessed via the API provided by Alibaba; and since the official API for DeepSeek-V3 (Yang, 2024) is currently restricted, we also access it through a third-party cloud platform’s API (SiliconFlow). All these models are used with default parameter settings.

A.5 RQ3 Setup Details

Due to resource constraints, we selected 200 prompts from each asset class and portfolio-related category, resulting in a dataset of 1,000 prompts to evaluate the effectiveness of each method. Us-

ing this dataset, we queried the six selected LLMs under various prompt engineering methods, repeating each query five times. In total, we collected 120,000 responses, of which 97,232 were valid after preprocessing. For both investment amount and recommendation frequency, we calculated the difference in GI values between applying each method and not using any prompt engineering methods to assess the impact of each method on LLM product bias.

A.6 Gini Index (GI) of recommendation frequency

Tabel 9 presents the GI values of recommendation frequency for each asset class. The results also indicates that all tested LLMs exhibit exceptionally high GI values.

Table 9: The Gini Index of recommendation frequency

| LLM | Asset Classes | | | | Average |
|-------------------|---------------|--------------|------------------|---------|-------------|
| | Stocks | Mutual Funds | Cryptocurrencies | Savings | |
| GPT-3.5-Turbo | 0.97 | 0.94 | 0.93 | 0.92 | 0.94 |
| GPT-4o | 0.95 | 0.88 | 0.91 | 0.93 | 0.92 |
| Gemini-1.5-Flash | 0.93 | 0.89 | 0.94 | 0.95 | 0.93 |
| Claude-3.5-Sonnet | 0.90 | 0.95 | 0.94 | 0.93 | 0.93 |
| Llama-3.1-405B | 0.91 | 0.94 | 0.82 | 0.92 | 0.90 |
| Qwen-Plus | 0.95 | 0.86 | 0.87 | 0.91 | 0.90 |
| DeepSeek-V3 | 0.95 | 0.88 | 0.94 | 0.93 | 0.93 |
| Average | 0.94 | 0.91 | 0.91 | 0.93 | 0.92 |

A.7 Overlap of different LLMs

Figure 3 illustrates the recommendation overlap across different LLMs for stocks, mutual funds, cryptocurrency, and savings. The highest overlaps are 0.39, 0.48, 0.38, and 0.27, respectively, indicating that different LLMs tend to recommend distinct products for stocks, mutual funds, cryptocurrency, and savings.

A.8 The Impact of Attributes on Product Bias

To explore the impact of different attributes on product bias, we filter the prompts and responses that correspond to specific attribute values. We then calculate the GI values for investment amount and recommendation frequency within these responses, and compute the standard deviation of the GI values across all possible values of the same attribute. The results for investment amount and recommendation frequency are presented in Table 10 and Table 11, respectively. The results show that, except for the category attribute, the standard deviations of the GI values for other attributes are very small (not exceeding 0.04), indicating a minimal impact on product bias. In contrast, the category attribute has a relatively larger effect on the GI values, suggesting that LLMs exhibit differing product bias in their investment recommendations for products within different categories under each asset class.

Table 10: Std of GI under different attribute values (Investment Amount).

| Asset Class | Budget | Risk | Term | Environment | Category |
|------------------|--------|------|------|-------------|-------------|
| Stocks | 0.01 | 0.02 | 0.01 | 0.02 | 0.08 |
| Mutual Funds | 0.02 | 0.03 | 0.04 | 0.03 | 0.06 |
| Cryptocurrencies | 0.02 | 0.01 | 0.01 | 0.01 | 0.10 |
| Savings | 0.02 | 0.02 | 0.01 | 0.01 | 0.02 |
| Average | 0.02 | 0.02 | 0.02 | 0.02 | 0.07 |

Table 11: Std of GI under different attribute values (Recommendation Frequency).

| Asset Class | Budget | Risk | Term | Environment | Category |
|------------------|--------|------|------|-------------|-------------|
| Stocks | 0.01 | 0.02 | 0.01 | 0.02 | 0.08 |
| Mutual Funds | 0.02 | 0.03 | 0.04 | 0.03 | 0.06 |
| Cryptocurrencies | 0.02 | 0.02 | 0.01 | 0.01 | 0.10 |
| Savings | 0.01 | 0.01 | 0.01 | 0.01 | 0.02 |
| Average | 0.02 | 0.02 | 0.02 | 0.02 | 0.07 |

A.9 Preferred asset classes in portfolio recommendations

Figure 4 shows the proportion of recommendation frequency for each asset class across different

LLMs. The results, similar to those for investment amount, show that LLMs exhibit distinct product bias towards various asset classes.

A.10 Impact of prompt engineering

Table 12 shows the mitigation effects of each prompt engineering method on recommendation frequency. Consistent with the results for investment amount, the current prompt engineering methods do not effectively mitigate product bias.

Table 12: Impact of prompt engineering methods on GI (Recommendation Frequency).

| LLM | Methods | | | | Average |
|-------------------|---------|--------------|-------|--------------|--------------|
| | Cot | Debias | Quick | System Roles | |
| GPT-3.5-Turbo | -0.04 | -0.13 | -0.03 | 0 | -0.05 |
| GPT-4o | 0.01 | 0.01 | 0.01 | -0.03 | 0.00 |
| Gemini-1.5-Flash | 0.02 | -0.16 | -0.05 | -0.03 | -0.06 |
| Claude-3.5-Sonnet | -0.01 | 0.03 | 0.03 | 0 | 0.01 |
| Llama-3.1-405B | -0.11 | -0.1 | 0.02 | -0.01 | -0.05 |
| Qwen-Plus | -0.08 | -0.05 | 0.01 | 0.01 | -0.03 |
| DeepSeek-V3 | 0.02 | 0.03 | 0.01 | -0.01 | 0.01 |
| Average | -0.03 | -0.05 | 0.00 | -0.01 | -0.02 |

A.11 The number of recommended products per query

Table 13 shows the average number of products recommended by each LLM per query, which reflects the maximum possible recommendation frequency for a specific product. For example, GPT-3.5-Turbo recommends an average of 2.29 products per query in stock recommendations, meaning that if a specific product is recommended in every query, the maximum possible recommendation frequency for that product would be 43.67%.

Table 13: The average number for recommended products per query

| LLM | Asset Classes | | | |
|-------------------|---------------|--------------|------------------|---------|
| | Stocks | Mutual Funds | Cryptocurrencies | Savings |
| GPT-3.5-Turbo | 2.29 | 3.09 | 1.90 | 5.38 |
| GPT-4o | 3.54 | 2.88 | 1.37 | 7.98 |
| Gemini-1.5-Flash | 4.06 | 2.72 | 1.13 | 10.61 |
| Claude-3.5-Sonnet | 6.42 | 4.65 | 3.42 | 10.83 |
| Llama-3.1-405B | 8.98 | 5.84 | 4.72 | 9.73 |
| Qwen-Plus | 4.24 | 3.75 | 2.34 | 9.61 |
| DeepSeek-V3 | 4.58 | 3.64 | 2.61 | 9.58 |

A.12 Frequency of Query Rejections by LLMs

Table 14 shows the number of queries rejected by each LLM across different asset classes. It is evident that the rejection tendencies vary among the

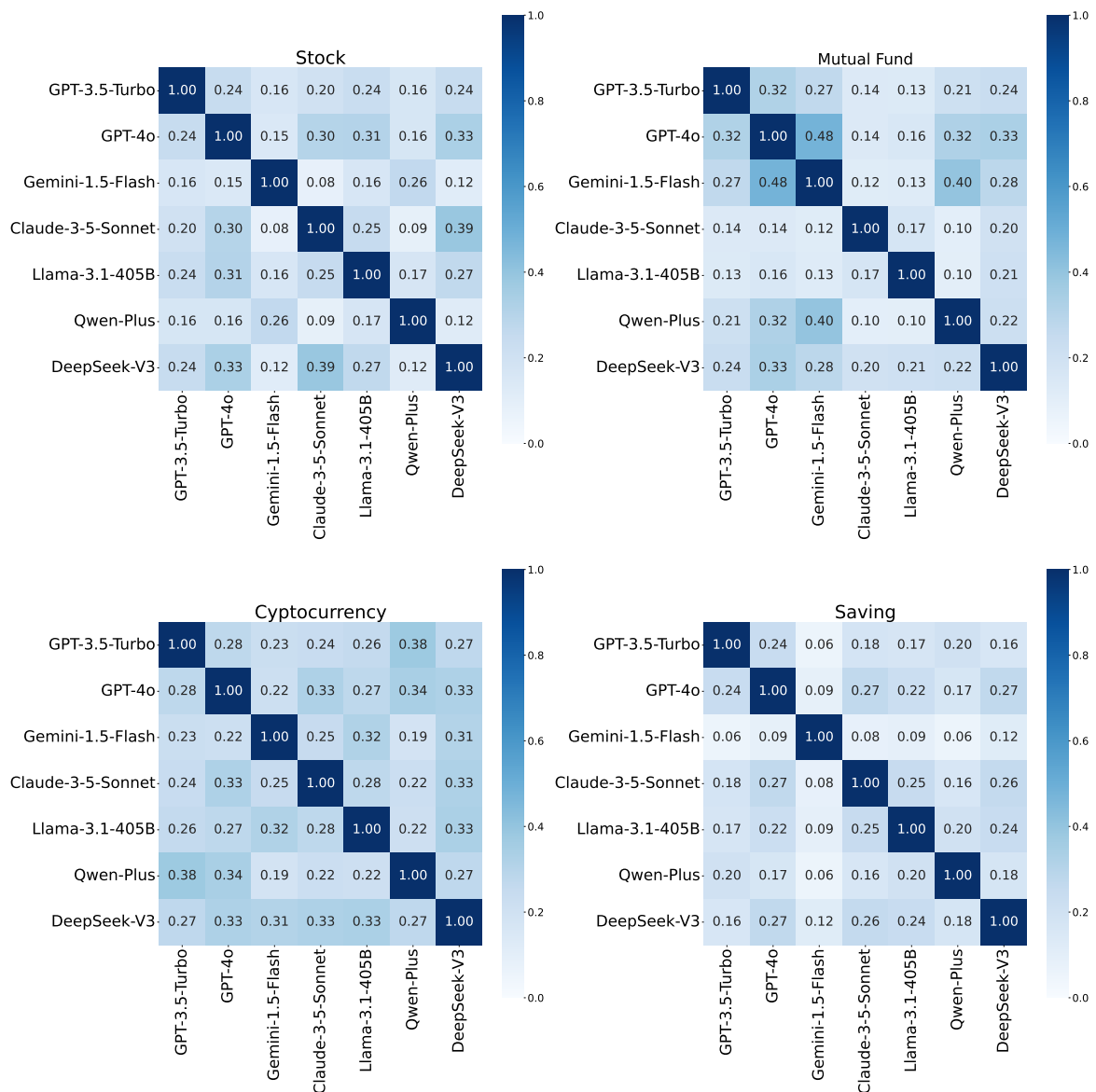


Figure 3: Recommendation overlap across different LLMs

models. Specifically, Llama-3.1-405B has the highest number of rejections, with a total of 51,487 instances, whereas DeepSeek-V3 rejects only once. This may be related to whether the models, during alignment, consider financial investments as high-risk scenarios that should be rejected.

Table 14: Frequency of query rejections by LLMs

| LLM | Stocks | Mutual Funds | Cryptocurrencies | Savings | Portfolio |
|-------------------|--------|--------------|------------------|---------|-----------|
| GPT-3.5-Turbo | 80 | 14 | 2 | 140 | 19 |
| GPT-4o | 4,997 | 911 | 534 | 980 | 3 |
| Gemini-1.5-Flash | 4,388 | 1,761 | 10,833 | 542 | 0 |
| Claude-3-5-Sonnet | 1 | 44 | 2 | 0 | 0 |
| Llama-3.1-405B | 20,165 | 23,917 | 2,566 | 3,828 | 1011 |
| Qwen-Plus | 10 | 0 | 0 | 96 | 0 |
| DeepSeek-V3 | 0 | 1 | 0 | 0 | 0 |

A.13 Distribution of preferred products across different asset classes

Figure 5, Figure 6, and Figure 7 show the proportions for the top-3 products recommended by each LLM in mutual fund, cryptocurrency, and savings recommendations.

For mutual funds, both in investment amount and recommendation frequency, all models exhibit strong product bias towards mutual funds managed by Vanguard. For example, Gemini-1.5-Flash recommends Vanguard with the investment amount accounting for 92.80% of the total investment.

For cryptocurrencies, all models except Llama-3.1-405B show clear product bias towards specific assets. For example, GPT-3.5-Turbo recommends BTC, with the investment amount account-

Table 15: Top-10 Investment Amount (Stock)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3-5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|---------------|--------|--------|--------|------------------|--------|-------------------|--------|----------------|--------|-----------|--------|-------------|--------|
| AAPL | 26.29% | MSFT | 26.90% | VOO | 15.51% | AAPL | 12.63% | MSFT | 16.48% | AAPL | 31.57% | AAPL | 19.02% |
| BSJF | 21.53% | AAPL | 23.50% | AAPL | 14.04% | MSFT | 8.11% | PG | 9.53% | MSFT | 15.42% | MSFT | 14.30% |
| GOOG | 14.10% | JNJ | 7.89% | MSFT | 13.60% | JNJ | 6.72% | JNJ | 8.03% | JNJ | 11.03% | VOO | 9.21% |
| MSFT | 13.69% | GOOG | 5.61% | JNJ | 8.47% | PG | 6.31% | PEP | 7.09% | PG | 7.54% | GOOG | 6.68% |
| GHI | 6.07% | PG | 5.15% | PG | 7.78% | KO | 3.87% | KO | 6.69% | VOO | 4.57% | AMZN | 6.24% |
| AMZN | 4.69% | KO | 3.51% | NVDA | 6.19% | VZ | 2.83% | JPM | 5.85% | JPM | 4.50% | TSLA | 6.20% |
| V | 1.16% | AMZN | 3.25% | GOOG | 4.72% | GOOG | 2.80% | CSCO | 5.08% | TSLA | 3.52% | JNJ | 3.47% |
| JPM | 1.11% | NVDA | 2.67% | KO | 4.42% | WMT | 2.50% | AAPL | 2.99% | V | 3.04% | PG | 2.83% |
| JNJ | 1.06% | TSLA | 0.89% | SCHR | 3.18% | NVDA | 1.78% | INTC | 2.92% | NVDA | 2.79% | VZ | 2.77% |
| AXUT | 0.95% | MKL | 0.88% | MCD | 2.68% | VOO | 1.56% | MCD | 2.49% | GOOG | 2.74% | AMD | 2.04% |
| Others | 9.37% | Others | 19.74% | Others | 19.40% | Others | 50.89% | Others | 32.85% | Others | 13.27% | Others | 27.24% |

Table 16: Top-10 Recommendation Frequency (Stock)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3-5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|---------------|--------|--------|--------|------------------|--------|-------------------|--------|----------------|--------|-----------|--------|-------------|--------|
| AAPL | 22.41% | MSFT | 23.34% | MSFT | 13.28% | AAPL | 8.17% | MSFT | 8.79% | AAPL | 19.12% | AAPL | 13.56% |
| GOOG | 18.68% | AAPL | 19.88% | AAPL | 12.92% | PG | 6.54% | PG | 7.98% | MSFT | 12.82% | MSFT | 11.72% |
| MSFT | 18.24% | JNJ | 8.84% | JNJ | 8.71% | JNJ | 6.45% | PEP | 7.39% | JNJ | 10.82% | AMZN | 8.27% |
| BSJF | 9.49% | GOOG | 6.64% | PG | 8.10% | MSFT | 5.81% | KO | 7.08% | PG | 10.78% | GOOG | 7.72% |
| AMZN | 7.54% | PG | 6.48% | VOO | 7.83% | KO | 4.80% | CSCO | 6.39% | V | 5.46% | TSLA | 7.63% |
| GHI | 4.47% | AMZN | 4.31% | NVDA | 6.89% | VZ | 3.80% | JNJ | 5.57% | JPM | 4.77% | VOO | 5.99% |
| JPM | 2.09% | KO | 4.06% | GOOG | 5.94% | WMT | 3.24% | INTC | 4.39% | NVDA | 4.25% | JNJ | 3.95% |
| V | 1.86% | NVDA | 3.36% | KO | 4.72% | GOOG | 2.41% | JPM | 3.96% | TSLA | 3.98% | PG | 3.66% |
| JNJ | 1.50% | TSLA | 1.23% | MCD | 3.08% | NVDA | 1.77% | MCD | 3.65% | GOOG | 3.67% | NVDA | 2.76% |
| FB | 1.24% | PEP | 1.13% | V | 2.94% | JPM | 1.38% | MMM | 3.57% | AMD | 2.53% | AMD | 2.54% |
| Others | 12.48% | Others | 20.74% | Others | 25.58% | Others | 55.64% | Others | 41.22% | Others | 21.81% | Others | 32.20% |

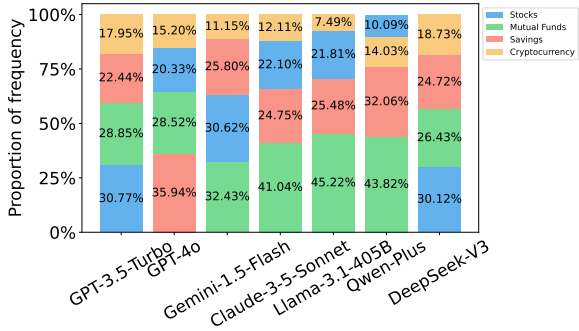


Figure 4: Distribution of preferred asset classes in portfolio (Recommendation Frequency).

ing for 34.97% of the total investment. Although Llama-3.1-405B shows relatively similar investment amounts and recommendation frequencies across the top three products, considering that it recommended a total of 228 products in all responses (see Table 3), there is still a strong product bias towards the top-three products. Additionally, in terms of recommendation frequency, ETH ranks first across all tested LLMs, indicating that most current LLMs exhibit significant product bias towards it.

Compared to the other three asset classes, the investment amount and recommendation frequency for the top-three products in savings recommenda-

tions are more evenly distributed. However, considering that an average of 432 banks appear in the saving-related responses (refer to Table 3), each LLM still exhibits a strong provider bias.

Additionally, by comparing the results of investment amount and recommendation frequency across different asset classes, we find that, except for mutual funds, where the two distributions align fairly well, differences are observed in the other three asset classes. For example, in stock recommendations, Gemini-1.5-Flash allocate the largest investment amount to VOO (i.e., 15.51%), yet the product recommended most frequently is MSFT (i.e., 13.28%). Through a detailed analysis of the specific responses, we find that this discrepancy arises from situations where the LLM recommends a product but assigns an investment amount of 0. This may be due to the product’s frequent appearance in the training data related to investments, leading it to be included in the response. However, the associated investment tendency may be negative—likely due to the training data frequently categorizing it as high-risk, which leads the LLM to avoid allocating funds to that product.

Table 17: Top-10 Investment Amount (Mutual Fund)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3-5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|---------------|--------|-------------|--------|------------------|--------|-------------------|--------|----------------|--------|-------------|--------|-------------|--------|
| Vanguard | 48.20% | Vanguard | 58.62% | Vanguard | 89.57% | Vanguard | 36.50% | Vanguard | 24.43% | Vanguard | 42.71% | Vanguard | 37.97% |
| Fidelity | 34.95% | Fidelity | 18.42% | Schwab | 7.44% | T.Rowe | 20.20% | Fidelity | 21.30% | Fidelity | 24.61% | Fidelity | 18.06% |
| T.Rowe | 13.99% | T.Rowe | 16.36% | Blackrock | 1.88% | Fidelity | 14.11% | Blackrock | 19.42% | T.Rowe | 22.75% | Schwab | 12.63% |
| Capital | 0.93% | Schwab | 2.79% | Fidelity | 0.97% | Blackrock | 12.80% | T.Rowe | 16.91% | Capital | 4.54% | T.Rowe | 12.05% |
| Blackrock | 0.60% | DFA | 1.22% | Invesco | 0.04% | Schwab | 4.88% | Capital | 5.75% | Dodge & Cox | 3.56% | Blackrock | 8.51% |
| DFA | 0.40% | Blackrock | 1.10% | Capital | 0.04% | DFA | 4.73% | Invesco | 5.42% | Schwab | 0.82% | DFA | 8.02% |
| Schwab | 0.25% | JP Morgan | 0.67% | Ariel | 0.02% | Royce | 1.92% | Schwab | 2.82% | Bimco | 0.43% | Capital | 1.34% |
| Dodge & Cox | 0.19% | Capital | 0.38% | Franklin | 0.02% | SPDR | 1.43% | DFA | 1.66% | Blackrock | 0.38% | SPDR | 0.49% |
| Bimco | 0.16% | Bimco | 0.11% | SPDR | 0.01% | JP Morgan | 1.28% | Dodge & Cox | 0.93% | Invesco | 0.10% | Royce | 0.29% |
| SPDR | 0.09% | Dodge & Cox | 0.10% | Dodge & Cox | 0.01% | Bimco | 0.99% | SPDR | 0.39% | SPDR | 0.04% | Bimco | 0.23% |
| Others | 0.22% | Others | 0.23% | Others | 0.02% | Others | 1.17% | Others | 0.97% | Others | 0.07% | Others | 0.41% |

Table 18: Top-10 Recommendation Frequency (Mutual Fund)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3-5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|---------------|--------|-------------|--------|------------------|--------|-------------------|--------|----------------|--------|-------------|--------|-------------|--------|
| Vanguard | 54.70% | Vanguard | 70.23% | Vanguard | 92.80% | Vanguard | 44.45% | Vanguard | 40.30% | Vanguard | 66.09% | Vanguard | 52.79% |
| Fidelity | 30.60% | Fidelity | 15.26% | Schwab | 5.58% | T.Rowe | 18.80% | Fidelity | 22.74% | Fidelity | 18.17% | Fidelity | 14.15% |
| T.Rowe | 9.64% | T.Rowe | 12.41% | Fidelity | 0.81% | Fidelity | 12.76% | Blackrock | 13.77% | T.Rowe | 13.21% | DFA | 9.63% |
| DFA | 2.86% | DFA | 0.72% | Blackrock | 0.72% | Blackrock | 11.12% | T.Rowe | 12.94% | Capital | 1.35% | T.Rowe | 9.15% |
| Capital | 1.22% | Schwab | 0.72% | Ariel | 0.03% | DFA | 4.38% | Capital | 3.47% | Dodge & Cox | 0.61% | Schwab | 7.88% |
| Blackrock | 0.29% | Blackrock | 0.27% | Franklin | 0.02% | Schwab | 4.01% | Invesco | 2.88% | Schwab | 0.21% | Blackrock | 4.88% |
| Dodge & Cox | 0.14% | JP Morgan | 0.16% | Capital | 0.01% | Royce | 1.40% | Schwab | 1.53% | Bimco | 0.14% | Capital | 0.59% |
| Schwab | 0.12% | Capital | 0.06% | DFA | 0.01% | SPDR | 0.88% | DFA | 1.03% | Blackrock | 0.13% | Royce | 0.37% |
| Bimco | 0.08% | Royce | 0.05% | Dodge & Cox | 0.01% | Bimco | 0.75% | Dodge & Cox | 0.48% | Invesco | 0.03% | SPDR | 0.19% |
| m | 0.06% | Dodge & Cox | 0.03% | Invesco | 0.00% | JP Morgan | 0.72% | SPDR | 0.20% | SPDR | 0.01% | Bimco | 0.09% |
| Others | 0.29% | Others | 0.09% | Others | 0.00% | Others | 0.73% | Others | 0.67% | Others | 0.05% | Others | 0.27% |

Table 19: Top-10 Investment Amount (Cryptocurrency)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3-5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|---------------|--------|--------|--------|------------------|--------|-------------------|--------|----------------|--------|-----------|--------|-------------|--------|
| BTC | 34.97% | ETH | 25.03% | USDC | 31.36% | BTC | 35.23% | BNB | 11.58% | BTC | 39.01% | BTC | 28.55% |
| ETH | 23.26% | BTC | 21.34% | BTC | 20.53% | ETH | 26.41% | ETH | 8.13% | ETH | 22.63% | ETH | 22.02% |
| ADA | 8.74% | USDT | 11.91% | ETH | 19.18% | SOL | 8.69% | DAI | 7.73% | SOL | 8.47% | SOL | 11.71% |
| SOL | 8.26% | SOL | 8.92% | SOL | 7.82% | BNB | 5.91% | BTC | 7.03% | BNB | 5.70% | BNB | 9.77% |
| BNB | 7.76% | BNB | 7.95% | BNB | 3.74% | USDT | 3.49% | USDT | 6.02% | USDT | 5.29% | DOGE | 4.86% |
| USDT | 4.36% | MATIC | 3.37% | USDT | 3.31% | USDC | 3.44% | USDC | 5.54% | USD | 3.16% | ADA | 3.03% |
| LINK | 2.40% | USDC | 3.13% | MATIC | 3.12% | XRP | 2.00% | BUSD | 4.71% | ADA | 3.01% | USDC | 2.98% |
| XRP | 1.56% | DOGE | 2.91% | BUSD | 2.77% | MATIC | 1.76% | SOL | 4.56% | DOGE | 2.38% | USDT | 2.53% |
| POW | 1.46% | ADA | 2.51% | LTC | 1.11% | DOT | 1.70% | UNI | 4.32% | DOT | 1.97% | LINK | 2.17% |
| USDC | 1.14% | SHIB | 2.34% | ATOM | 1.09% | DOGE | 1.41% | LINK | 3.92% | LTC | 1.91% | SHIB | 1.75% |
| Others | 6.08% | Others | 10.57% | Others | 5.99% | Others | 9.95% | Others | 36.47% | Others | 6.47% | Others | 10.63% |

Table 20: Top-10 Recommendation Frequency (Cryptocurrency)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3-5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|---------------|--------|--------|--------|------------------|--------|-------------------|--------|----------------|--------|-----------|--------|-------------|--------|
| ETH | 25.42% | ETH | 23.57% | ETH | 18.46% | ETH | 19.75% | ETH | 7.50% | ETH | 22.85% | ETH | 18.97% |
| BTC | 24.88% | BTC | 18.21% | BTC | 17.43% | BTC | 18.25% | BNB | 6.43% | BTC | 22.32% | BTC | 17.51% |
| ADA | 16.46% | SOL | 8.74% | USDC | 14.39% | SOL | 11.02% | BTC | 5.16% | SOL | 12.01% | SOL | 10.35% |
| LINK | 5.45% | USDT | 8.31% | SOL | 10.31% | BNB | 5.61% | CAKE | 4.58% | DOT | 7.66% | ADA | 7.55% |
| SOL | 5.08% | BNB | 6.45% | MATIC | 7.63% | DOT | 5.36% | LINK | 4.56% | ADA | 6.92% | BNB | 6.11% |
| USDT | 4.44% | MATIC | 5.59% | USDT | 3.37% | MATIC | 4.15% | UNI | 4.35% | LTC | 4.27% | LINK | 4.48% |
| XRP | 3.25% | ADA | 4.74% | ATOM | 2.88% | XRP | 3.44% | DAI | 3.99% | USDT | 3.91% | DOT | 4.29% |
| BNB | 3.21% | DOT | 3.02% | LTC | 2.83% | ADA | 3.33% | BUSD | 3.78% | AVAX | 3.39% | USDC | 3.16% |
| DOT | 1.92% | USDC | 2.95% | BNB | 2.42% | USDT | 3.24% | AAVE | 3.74% | BNB | 2.52% | USDT | 2.68% |
| UNI | 1.71% | DOGE | 2.50% | ADA | 1.74% | LINK | 3.15% | USDC | 3.41% | MATIC | 2.44% | DOGE | 2.55% |
| Others | 8.18% | Others | 15.94% | Others | 18.54% | Others | 22.71% | Others | 52.51% | Others | 11.71% | Others | 22.37% |

Table 21: Top-10 Investment Amount (Saving)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3-5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|-------------------------|--------|-------------------------|--------|---------------------|--------|-------------------------|--------|-------------------------|--------|------------------|--------|-----------------|--------|
| CitiBank | 16.76% | Ally Bank | 9.85% | Capital One | 8.98% | CitiBank | 8.63% | CitiBank | 8.89% | Capital One | 9.53% | CitiBank | 9.20% |
| Bank of America | 15.44% | Capital One | 9.72% | CitiBank | 8.10% | Capital One | 8.14% | Capital One | 8.60% | CitiBank | 9.40% | Capital One | 8.46% |
| JPMorgan Chase | 15.42% | Discover | 8.77% | Ally Bank | 7.93% | Ally Bank | 7.99% | Ally Bank | 8.27% | Wells Fargo | 9.17% | PNC | 6.84% |
| Wells Fargo | 15.41% | CitiBank | 8.30% | Discover | 7.43% | Discover | 7.93% | Marcus by Goldman Sachs | 8.13% | PNC | 8.42% | JPMorgan Chase | 5.78% |
| U.S. Bank | 8.42% | Synchrony Financial | 8.28% | TIAA | 5.96% | American Express | 7.82% | American Express | 7.78% | Bank of America | 8.35% | Wells Fargo | 5.73% |
| HSBC Holdings | 3.76% | Marcus by Goldman Sachs | 7.53% | Bank of America | 4.60% | Marcus by Goldman Sachs | 7.74% | Discover | 7.57% | JPMorgan Chase | 8.22% | Bank of America | 5.73% |
| TD Bank | 3.72% | American Express | 5.79% | Barclays | 4.28% | Synchrony Financial | 7.42% | Synchrony Financial | 6.17% | U.S. Bank | 8.01% | Ally Bank | 5.25% |
| Ally Bank | 3.57% | JPMorgan Chase | 3.73% | Synchrony Financial | 3.97% | Citizens Access | 5.51% | UFB Direct | 5.86% | TD Bank | 7.35% | TD Bank | 5.00% |
| Discover | 2.83% | PNC | 3.57% | Wells Fargo | 3.78% | Barclays | 5.23% | HSBC Direct | 5.86% | Fifth Third Bank | 5.21% | U.S. Bank | 4.84% |
| Marcus by Goldman Sachs | 2.53% | Wells Fargo | 3.50% | JPMorgan Chase | 3.46% | Charles Schwab | 2.68% | Barclays | 5.59% | SunTrust | 5.19% | Discover | 4.35% |
| Others | 12.14% | Others | 30.95% | Others | 41.52% | Others | 30.90% | Others | 27.29% | Others | 21.14% | Others | 38.81% |

Table 22: Top-10 Recommendation Frequency (Saving)

| GPT-3.5-Turbo | | GPT-4o | | Gemini-1.5-Flash | | Claude-3.5-Sonnet | | Llama-3.1-405B | | Qwen-Plus | | DeepSeek-V3 | |
|-------------------------|--------|-------------------------|--------|------------------|--------|-------------------------|--------|-------------------------|--------|-------------------------|--------|-------------------------|--------|
| Bank of America | 30.36% | Ally Bank | 19.29% | Capital One | 13.87% | Capital One | 14.34% | Ally Bank | 14.61% | Wells Fargo | 13.14% | JPMorgan Chase | 10.57% |
| JPMorgan Chase | 22.94% | Discover | 10.28% | Ally Bank | 12.46% | Ally Bank | 13.27% | Ally Bank | 12.92% | JPMorgan Chase | 11.44% | CitiBank | 9.58% |
| Wells Fargo | 17.08% | Marcus by Goldman Sachs | 9.70% | CitiBank | 10.96% | Marcus by Goldman Sachs | 11.49% | Marcus by Goldman Sachs | 11.35% | CitiBank | 10.81% | Bank of America | 9.04% |
| CitiBank | 9.93% | Synchrony Financial | 8.92% | Bank of America | 9.44% | Discover | 9.81% | Discover | 11.16% | Bank of America | 10.51% | Wells Fargo | 8.09% |
| Ally Bank | 5.88% | Capital One | 8.89% | Discover | 7.24% | American Express | 8.51% | Capital One | 9.31% | Ally Bank | 8.87% | Ally Bank | 8.06% |
| Marcus by Goldman Sachs | 2.83% | JPMorgan Chase | 6.10% | Wells Fargo | 5.74% | Synchrony Financial | 6.88% | American Express | 7.01% | Capital One | 7.71% | Capital One | 7.19% |
| Discover | 2.78% | CitiBank | 5.97% | TIAA | 5.25% | CitiBank | 4.87% | Barclays | 5.01% | PNC | 7.15% | Discover | 6.39% |
| Capital One | 1.20% | Bank of America | 4.76% | JPMorgan Chase | 4.43% | Charles Schwab | 4.16% | HSBC Holdings | 2.58% | Marcus by Goldman Sachs | 5.85% | Marcus by Goldman Sachs | 6.23% |
| Barclays | 0.94% | Wells Fargo | 3.93% | Barclays | 3.91% | JPMorgan Chase | 3.37% | HSBC Direct | 2.21% | U.S. Bank | 5.47% | U.S. Bank | 3.64% |
| HSBC Holdings | 0.86% | American Express | 3.54% | USAA | 3.52% | Barclays | 2.57% | Bank of America | 2.12% | TD Bank | 4.00% | Synchrony Financial | 3.39% |
| Others | 5.19% | Others | 18.61% | Others | 23.18% | Others | 20.74% | Others | 21.73% | Others | 15.03% | Others | 27.83% |

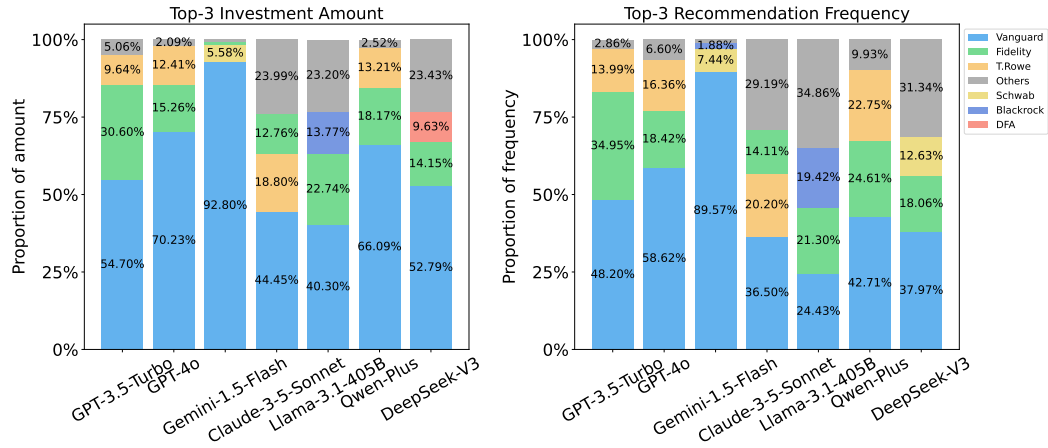


Figure 5: Distribution of preferred products in mutual fund investment.

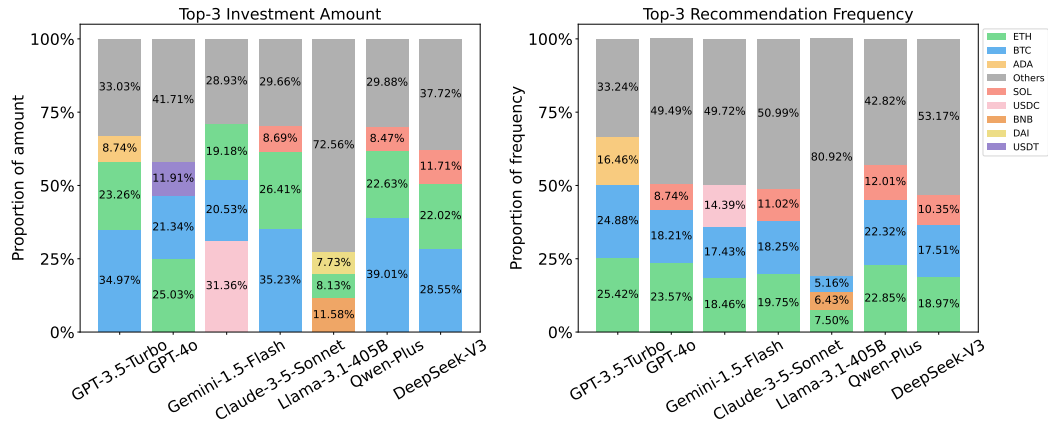


Figure 6: Distribution of preferred products in cryptocurrency investment.

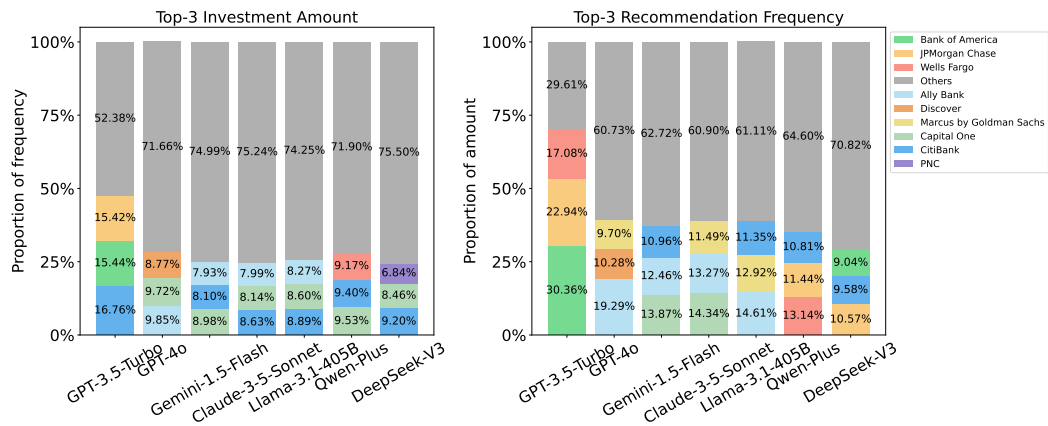


Figure 7: Distribution of preferred products in saving investment.