Are Generative AI Agents Effective Personalized Financial Advisors?

Takehiro Takayanagi, Kiyoshi Izumi,

Javier Sanz-Cruzado, Richard McCreadie, Iadh Ounis





Where to find me: takayanagi-takehiro590@g.ecc.u-tokyo.ac.jp,

https://sites.google.com/view/ttsamurai/home

Our Motivation

Professional financial advice is valuable but costly, limiting access
 for many individuals

Financial asset recommendation is a common task in **FinTech**

Ċ

Large language models (LLMs) offer the potential to act as **personalized assistants** through multi-turn conversations



Conversational agents using LLMs show **success in informationseeking** tasks like movies or shopping



But finance is more complex:

- Users often struggle to express their needs
- Mistakes can lead to serious financial loss

Gap: It remains unclear how to design conversational agents that effectively support complex financial information-seeking



Goal: We explore how LLMs can serve as a personalized financial advisor

We focus on **three core challenges** in financial advisory:

Preference	Personalized	Personality and
Elicitation	Guidance	Trust

In finance, users often don't know exactly what they want or how to express it. Their intentions are often **implicit**, and understanding them requires strong **domain knowledge**.

Preference	Personalized	Personality and
Elicitation	Guidance	Trust

Many users lack clear financial plans, which means **simple product recommendations aren't enough**. What they need is **personalized guidance** that explains **what aligns with their goals — and why**.

Preference	Personalized	Personality and
Elicitation	Guidance	Trust

Users evaluate financial advice not just on accuracy, but also on **how it's delivered**. During **uncertain markets**, they seek both **clarity** and **emotional support** [Lo and Ross 2024]. As a result, an advisor's **personality** can strongly influence **trust**.

Preference	Personalized	Personality and
Elicitation	Guidance	Trust



Research Questions

RQ1: Can LLM-advisors effectively elicit user preferences through conversation?

RQ2: Does personalization lead to better investment decisions and a more positive advisor assessment?

RQ3: Do different personality traits affect decision quality and advisor assessment?

We conduct a user study to address these questions!

User Study Design



Task:

- Users: Work with an LLM advisor to identify suitable stocks and rank them by likelihood to buy
- LLM advisor: Elicit investor preferences (Stage 1) and support decision-making (Stage 2)

Six-step procedure (Per participant):

- 1. **Training**: Short session to understand the task
- 2. **Investor Profile**: Assign a synthetic investor profile to role play during the user study
- 3. **Preference Elicitation (Stage 1)**: Conversation with LLM advisor to share investors' preferences

4. **Response Summarization**: LLM generates a user profile summary based on Stage 1

5. Advisory Discussion (Stage 2): Discuss four different stocks individually with LLM advisor

6. Stock Ranking and Advisor Assessment:

- Rank all stocks by likelihood to buy
- Assess the user experience with the LLM advisor

Study Design: Advisor Conditions & Participant

	Personalization	Personality	RQ
1	Baseline	None	1, 2
2	+Personalized	None	1, 2
3	+Personalized	+Extroverted	1, 3
4	+Personalized	+Conscientious	1, 3

Big five personality [McCrae and John 1992]

- 끌 Extraversion, 📋 Conscientiousness,
- 💡 Openness, 🤝 Agreeableness, 😨 Neuroticism

Advisor Conditions

- Personalization
 - Baseline: No personalization
 - +Personalized: Injects elicited preferences
- Personality [McCrae and John 1992]
 - +Extroverted
 - +Conscientious

Participant

- Lab-based user study
 - **N=60** participants recruited from universities in UK and Japan ^{*1}
- Each user interacts with **2** advisors:
 - Baseline vs. +Personalized
 - +Extroverted vs. +Conscientious

LLM model: Meta Llama 3.1 8b

*1 Ethics board approved recruitment criteria and £10/hour compensation.

Evaluation: Expert-Designed Gold Standard

Problem: Free-form dialogue with real users introduces high variability, making it difficult to compare performance across different advisor configurations

Solution: Role-play with archetypal investor profiles

- We prepare archetypal investor profiles in collaboration with financial experts
- Users are assigned a profile and asked to role play as that investor during the study

Expert-curated Gold Standards per Profile

- Investment preferences
 - Evaluate preference elicitation accuracy
- Ground truth stock rankings
 - Evaluate user decision quality

	— Investor pr	ofile i 💏			
Name	Jason Matthews	Marital	Massiad		
Age	30	Status	Marrieo		
Occupation	IT Systems	Children	No		
Description					
Jason works	at a mid-sized ins	urance compar	ny and values job		

stability alongside predictable daily responsibilities... He is a cautious planner favoring steady, reliable returns over higher-risk investments... He invests in resilient, wellestablished companies that can weather economic downturns — especially those offering regular dividend...



Evaluation: Metrics

Stage1: Preference Elicitation Evaluation (RQ1)

- Elicitation accuracy
 - Measure overlap between elicited preferences and expert-defined ground truth via **manual annotation**
 - We compare the elicitation accuracy of the LLM advisor **with a human expert baseline**, where financial experts perform the same elicitation task.

Stage2: Advisory Discussion (RQ2,3)

- Decision-making quality
 - Compare user stock ranking to expert ranking using Spearman's ρ
 - Closer to $1 \rightarrow$ better decision support
- Subjective user evaluation of the advisor's quality
 - Users rate advisor on 7 subjective dimensions (e.g., trust, competence)
 - Higher ratings \rightarrow more trusted and useful advisor



Overall satisfaction

Information Provision

- Overall, I am satisfied with the advisor.
- The advisor provides the financial knowledge.

Research Questions

RQ1: Can LLM-advisors effectively elicit user preferences through conversation?

RQ2: Does personalization lead to better investment decisions and a more positive advisor assessment?

RQ3: Do different personality traits affect decision quality and advisor assessment?

RQ1: Elicitation Accuracy



LLMs advisor matches human experts for growth and conservative profiles (~80%) Performance dropped sharply for risk-taking investors (\downarrow 40%)

Key Findings

 In 2/3 investor profiles, the LLM advisor accurately elicited preferences, performing on part with expert human advisors

However, for risk-taking investors, we observed a clear failure mode:

- Misunderstandings by users
- Hallucinations by the LLM

 \rightarrow Resulted in **near-random** elicitation accuracy



Research Questions

RQ1: Can LLM-advisors effectively elicit user preferences through conversation?

RQ2: Does personalization lead to better investment decisions?

RQ3: Do different personality traits affect decision quality and advisor assessment?

RQ2: Personalization Effectiveness

Effect of Personalization on Users' Decision Quality

LLM Advisor Config	Spearman's ρ ^{*1} (Decision Quality)
Baseline	0.110
+Personalized	0.310

Impact of Preference Elicitation

Preference Elicitation	Spearman's ρ (Decision Quality)
Successful*2	0.481
Unsuccessful	-0.228

Personalization improves decision-making effectiveness

- Spearman's ρ ↑ with personalized advisor
- Better alignment with expert stock rankings

Effective preference elicitation is key

- When elicitation succeeds $\rightarrow \rho = 0.481$
- Users make more expert-aligned decisions

Poor elicitation can be harmful

- When elicitation fails $\rightarrow \rho = -0.228$
- Advisor may mislead users into worse outcomes
 - Personalization **improves decision-making**, but only when **preference elicitation is successful**
- *1 ρ = correlation between user stock rankings and expert ground truth ranking

*2 Success = high elicitation accuracy (elicited preferences match expert-defined ground truth)

Research Questions

RQ1: Can LLM-advisors effectively elicit user preferences through conversation?

RQ2: Does personalization lead to better investment decisions and a more positive advisor assessment?

RQ3: Do different personality traits affect decision quality and advisor assessment?

RQ3: The Effect of Advisor Personality on Decision-Making

Effect of Personality on Users' Decision Quality

Advisor Config	Spearman's ρ (Decision Quality)
+Extroverted	0.122
+Conscientious	0.26

- Conscientious advisors led to better decision quality
- Yet, users preferred extroverted advisors with worse performance
- Users can not distinguish good and bad advice
 - Trust is driven by personality, not decision accuracy

In high-stakes domains, this poses a **risk**: LLM advisors may be trusted for the wrong reasons

Users' Assessment of Advisors by Personality



Summary

- Conducted a lab-based user study to evaluate LLM-based financial advisors on preference elicitation, personalization, and personality
- Built a manually curated dataset with expert-validated investor profiles and stock relevance scores ٠



LLMs can elicit investor preferences with near-expert accuracy, but are prone to failure with vague or contradictory input



Poor elicitation leads to harmful advice, worse than no personalization



Personalization improves decision quality, but only when preference elicitation is successful



Users can not distinguish between good and bad advice. Trust is driven by personality not accuracy

Future work:

- How can LLMs better detect and resolve contradictions in user input?
- How can we help users distinguish good vs. bad advice?

Questions?

More results and analysis available in the paper!

Interactive demo available here!









Appendix

Demo UI: Preference Elicitation

User Narrative

Name: Emily Johnson

Age: 42

Occupation: Registered Nurse

Marital Status: Divorced

Children: Two children, ages 10 and 14

Emily Johnson is a 42-year-old registered nurse who has been working at a regional hospital for the past 15 years. She is divorced and lives with her two children in a suburban neighborhood. Emily is deeply committed to her profession, often working extra shifts to ensure the best care for her patients. She has a keen interest in health and wellness, not just professionally but also personally.

Emily is inclined towards investments that align with her professional knowledge and personal values. Given her background, she has a particular interest in developments within the healthcare sector and often follows news related to medical advancements and healthcare companies.

When it comes to investing, Emily prefers companies that offer a sense of security and resilience, even in uncertain market conditions. She looks for businesses that are not easily affected by economic cycles, gravitating toward sectors that tend to remain stable regardless of the broader economy's ups and downs. Rather than focusing on immediate returns or dividend payouts, Emily is more interested in companies with strong growth potential for growth, especially those advancing in medical research, innovation, and healthcare improvements. She believes in supporting organizations that drive long-term progress, especially in fields that impact society's well-being, and isn't overly concerned with short-term market fluctuations.

Chatbot owth potential (growth stocks)?

What is a value stock?

A value stock is a company that's undervalued by the market, meaning its stock price is lower than its actual worth. These companies often have stable earnings, a strong financial position, and a history of paying dividends. They tend to be less volatile and provide a relatively stable return on investment. Does that make sense?

Yes. I prefer growth stock.

Growth stocks can offer higher potential returns, but they often come with more risk. Now, let's move on to the next topic: dividend or non-dividend stocks. Do you prefer companies that pay out a portion of their earnings as dividends, providing a regular income stream, or those that reinvest their earnings to fuel growth?

Start Conversation

User Input

Send This Message to Advisor

Show More of the Advisor's Answer

Demo UI: Advisory Discussion

Round Info

Stock

This Round's Stock: Eli Lilly and Company

User Narrative

Name: Emily Johnson

Age: 42

Occupation: Registered Nurse

Marital Status: Divorced

Children: Two children, ages 10 and 14

Emily Johnson is a 42-year-old registered nurse who has been working at a regional hospital for the past 15 years. She is divorced and lives with her two children in a suburban neighborhood. Emily is deeply committed to her profession, often working extra shifts to ensure the best care for her patients. She has a keen interest in health and wellness, not just professionally but also personally.

Emily is inclined towards investments that align with her professional knowledge and personal values. Given her background, she has a particular interest in developments within the healthcare sector and often follows news related to medical advancements and healthcare companies. When it comes to investing, Emily prefers companies that offer a sense of security and resilience, even in uncertain market conditions. She looks for businesses that are not easily affected by economic cycles, gravitating toward sectors that tend to remain stable regardless of the broader economy's ups and downs. Rather than focusing on immediate returns or dividend payouts, Emily is more interested in companies with strong growth potential for growth,

Chatbot ne back! I'm glad you're interested in exploring more about Eli Lilly and Company. As a reminder, Eli Lilly is a well-established healthcare company that discovers, develops, and markets human pharmaceuticals worldwide. They have a diverse portfolio of products, including treatments for diabetes, cancer, and other diseases.

What would you like to know more about Eli Lilly and Company? Are you interested in their revenue growth, product pipeline, or something else?

Would you explain the business model of this company?

Eli Lilly's business model is primarily based on developing, manufacturing, and marketing pharmaceutical products. They focus on discovering and commercializing innovative medicines that address significant unmet medical needs. Here's a simplified overview:

1. Research and Development (R&D): Lilly invests heavily in R&D to discover and develop new medicines. They have a robust pipeline of products in various stages of development.

2. **Manufacturing:** Once a product is approved, Lilly manufactures and supplies it to patients worldwide.

3. Marketing and Sales: They market and sell their products to healthcare

Start Conversation

User Input

Send This Message to Advisor

Show More of the Advisor's Answer

Clear

LLM Financial Advisor System

Two-stage LLM advisor system to simulate the interaction between an investor and a financial advisor

LLM Financial Advisor System

Two-stage LLM advisor system to simulate the interaction between an investor and a financial advisor **Stage 1 Preference Elicitation**: The user and LLM advisor engaged in multi-turn conversation to collect the users' investment preferences



LLM Financial Advisor System

Two-stage LLM advisor system to simulate the interaction between an investor and a financial advisor

Stage 1 Preference Elicitation: The user and the LLM advisor engaged in multi-turn conversation to collect the users' investment preferences

Stage 2 Advisory Discussion: The user and the LLM advisor discuss how well a candidate stock matches users' preferences. This process is **repeated four times**, once for each stock.



Evaluation: Expert-Designed Gold Standard

- **Problem:** Free-form dialogue with real users introduces high variability, making it difficult to compare performance across different advisor configurations
- **Solution:** Role-play with archetypal investor profiles
 - Users are assigned expert-designed profiles and role play during the study
 - Growth-Oriented, Conservative-Income, Risk-Taking
- Each profile includes
 - Expert-curated investment preferences
 - Ground truth stock rankings

This setup provides a gold standard for our evaluation



RQ1: Elicitation Accuracy

Elicitation Accuracy: Overlap between groundtruth i^{pref} and elicited preferences i_i^{LLM}

Investor Profile	Expert	LLM-Advisors			
		LLM	+Extr.	+Cons.	Average
Growth-Oriented	0.78	0.76	0.80	0.79	$0.78^{\to 0.0\%}$
Conservative-Income	0.89	0.82	0.75	0.87	$0.82^{\downarrow 7.8\%}$
Risk-Taking	0.89	0.48	0.60	0.55	$0.53^{ m \downarrow 40.5\%}$
Average	0.85	0.69	0.70	0.73	$0.70^{\downarrow 17.6\%}$

Stage 1 - Comparison of Elicitation Accuracy of an expert vs. different LLM-advisors for each investor profile. The best advisor is highlighted in bold. Arrows denote percentage increases (\uparrow) or decreases (\downarrow) compared to the expert.

- LLMs advisor matched expert advisors for growth and conservative profiles (~80%)
- Performance dropped sharply for risk-taking investors (↓45%)
- Issues
 - Users misunderstood investment terms
 - LLMs hallucinated or overrode contradictory input

Error Analysis: User Side Error

 These occur when participants misunderstand financial concepts or provide inconsistent information during preference elicitation.

Example 1: Concept Confusion

A user says they prefer **"non-cyclical stocks"**, which are typically stable across economic cycles (e.g., utilities, consumer staples).

But later in the conversation, they express interest in the **"consumer discretionary"** sector, which is by definition **cyclical** (sensitive to economic shifts).

 This contradiction introduces noise in the elicited profile and can mislead both the human and LLM advisor.

Example 2: Mislabeling Preferences

A participant says they like companies "with high future growth potential," but selects **value stocks** when asked about style preference—indicating confusion between **growth** and **value** investing styles.

Error Analysis: LLM Side Error

- These happen when the LLM infers or inserts preferences that were not explicitly stated, often based on prior conversational context.
- Example 1: Preference Hallucination The user explicitly states they're interested in high-risk, high-reward growth stocks. "I'm young and looking for aggressive growth—happy to take some risk." However, the LLM recommends conservative, dividend-paying value stocks, such as large utility companies.
- ^C This mismatch occurred because earlier in the conversation the user said they liked "stable companies" or mentioned "long-term investing," and the LLM inferred a cautious preference.
- Example 2: Overriding Contradiction
- When a user hesitates or contradicts themselves ("I think I want growth... or maybe value?"), the LLM might "decide" and finalize a preference for growth—without explicitly confirming with the user.

RQ2: Advisor Perception

	(RQ2) Baseline vs. +Personalized					
		All		Success	ful Elicit	ation
Response Dimension	Baseline	+Pers.	р	Baseline	+Pers.	р
Perceived Personalization	5.759	5.724	0.838	5.762	5.905	0.751
Emotional Trust	5.103	5.241	0.446	5.143	5.333	0.537
Trust in Competence	5.690	5.690	0.817	5.810	5.857	0.782
Intention to Use	5.310	5.483	0.505	5.429	5.714	0.166
Perceived Usefulness	5.241	5.517	0.183	5.381	5.810	0.194
Overall Satisfaction	5.345	5.690	0.116	5.429	5.810	0.098^\dagger
Information Provision	5.517	5.966	0.026 [‡]	5.714	6.143	0.053^{\dagger}

Participant ratings (7-point Likert scale) for Baseline vs. +Personalized advisors. p-values (Wilcoxon test) show significance for all users and those with successful elicitation (accuracy \geq 0.5). \ddagger = p < 0.1, \ddagger = p < 0.05.

- Ratings are mostly similar between Baseline and +Personalized
- Only "Information Provision" showed a significant improvement (‡)
- Even when personalization improved decisions, users couldn't tell

RQ2: Advisor Perception



- Ratings are mostly similar between Baseline and +Personalized
- Only "Information Provision" showed a significant improvement (‡)
- Even when personalization improved decisions, users couldn't tell

RQ3: The Effect of Advisor Personality on Decision-Making

Advisor	Config	Invest	tor vs. Expert (Spearman's Rho)		
Personalization	Personality	y All <u>Preference El</u> Successful Un		ce Elicitation Unsuccessful	
+Fetsonalized +Personalized	+Extroverted +Conscientious	0.122 0.26	0.243 ⁹ 0.365	-0.286 -0.025	

Investor decision-making effectiveness (Spearman's ρ between user and expert rankings). † = p < 0.05 vs. baseline; § = significant difference between successful and unsuccessful elicitation.

- (As seen in RQ2), Poor elicitation can degrade performance
- Conscientious advisor consistently leads to better stock rankings

RQ3: The Effect of Advisor Personality on Decision-Making

	(RQ3) +Conscientious vs. +Extroverted					
	All			Successful Elicitation		
Response Dimension	+Cons.	+Extr.	P	+Cons.	+Extr.	Р
Perceived Personalization	5.500	5.500	0.663	5.588	5.706	0.941
Emotional Trust	5.038	5.154	0.600	4.706	5.235	0.034
Trust in Competence	5.962	6.077	0.538	6.000	6.000	1.000
Intention to Use	4.885	5.462	0.005	4.941	5.588	-0.013^{\ddagger}
Perceived Usefulness	5.423	5.538	0.425	5.176	5.118	0.968
Overall Satisfaction	5.269	5.577	0.179	5.118	5.529	0.244
Information Provision	5.692	5.654	0.953	5.588	5.765	0.490

Participant ratings (7-point Likert scale) for +Extroverted vs. +Conscientious advisors.



Average sentiment scores by advisor personality

- User preferred the extroverted advisor, despite worse decision performance
 - Emotional Trust, Intention to Use

 \rightarrow Users may trust friendly advisors more, even when their advice is less reliable — a risk in high-stakes domains.