Kyoto University, Graduate School of Economics Discussion Paper Series



Comparing Risk Preferences and Loss Aversion in Humans and AI: A Persona-Based Approach with Fine-Tuning

Ryota Iwamoto, Takunori Ishihara, and Takanori Ida

Discussion Paper No. E-25-006

Graduate School of Economics Kyoto University Yoshida-Hommachi, Sakyo-ku Kyoto City, 606-8501, Japan

JULY, 2025

Comparing Risk Preferences and Loss Aversion in Humans and AI: A Persona-Based Approach with Fine-Tuning

Ryota Iwamoto¹, Takunori Ishihara², and Takanori Ida³

Abstract:

This study empirically investigates the differences in risk preferences and loss aversion between humans and generative AI. We conduct a nationwide online survey of 4,838 individuals and generate AI responses under identical conditions by using personas constructed from demographic attributes. The results show that in gain domains, both humans and the AI select risk-averse options and exhibit similar preference patterns. However, in loss domains, AI shows a stronger risk-loving tendency and responds more sharply to individual attributes such as gender, age, and income. We retrain the AI by fine-tuning it based on human choice data. After fine-tuning, the AI's preference distribution moves closer to that of humans, with loss-related decisions showing the greatest improvement. Using Wasserstein distance, we also confirm that fine-tuning reduces the behavioral gap between AI and humans.

JEL Classification: D91, C91

Keywords: bias, loss aversion, risk preference, generative AI, persona, fine-tuning, Wasserstein distance

¹ Graduate School of Public Policy, Kyoto University, <u>iwamoto.ryota.82a@st.kyoto-u.ac.jp</u>

² Institute for International Academic Research, Kyoto University of Advanced Science, <u>ishihara.takunori@kuas.ac.jp</u>

³ Graduate School of Economics, Kyoto University, <u>ida@econ.kyoto-u.ac.jp</u>

1. Introduction

As artificial intelligence (AI) has increasingly automated judgment processes, concerns have arisen regarding its legitimacy and explainability. Human-in-the-loop (HITL), a framework that integrates human and AI in economic decision making (Rahwan, 2018; Rahwan et al., 2019), offers both theoretical and practical responses to these concerns. The comparison of bias structures between humans and AI plays a central role in understanding the value of HITL. AI often inherits biases from its training data. It tends to reproduce the structural biases embedded in past economic behavior and institutions. Meanwhile, as behavioral economics reveals, human decision making also reflects systematic deviations due to emotion-driven cognitive biases. Understanding these distinct bias structures and creating designs that allow mutual correction between humans and AI are essential for achieving rationality. Therefore, when designing HITL-based decision systems, the commonalities and differences between algorithmic and cognitive biases, from both theoretical and empirical perspectives, must be examined to provide a theoretical foundation for developing reliable decision-support systems that leverage the complementary strengths of human judgment and AI reasoning.

To advance this line of analysis, systematic biases that arise in human decision making must be first clarified. Humans often deviate from the predictions of standard utility-maximization models. In particular, under conditions of risk or uncertainty, existing studies observe behavioral patterns such as the certainty effect and loss aversion. Kahneman and Tversky (1979) organize these tendencies into a theoretical framework known as prospect theory, which is the basis for studies that document various forms of irrational judgments. Tversky and Kahneman (1992) estimate the degree of loss aversion at approximately $\lambda = 2.25$ and suggest that people experience losses about twice as intensely as equivalent gains and therefore tend to avoid them. Falk et al. (2018) conducts a cross-national survey of approximately 80,000 individuals in 76 countries. They find that loss aversion varies not only across cultures but also according to individual attributes, such as age, gender, and education. Brown et al. (2024) conduct a meta-analysis of 150 studies in economics and psychology, reporting a median loss aversion coefficient of approximately $\lambda = 1.69$. These findings offer a foundational understanding of human behavior that informs the design of HITL systems, particularly regarding the role of human agents within such frameworks.

In recent years, the rise of generative AI, exemplified by models such as ChatGPT, has sparked growing interest in whether the cognitive biases commonly observed in humans also appear in AI. As generative AI models are trained on large-scale human-generated language data, these models may inherently acquire humanlike tendencies and biases. If generative AI can replicate such biases, it can serve as a practical alternative in fields such as social surveys, economic experiments, and policy simulations. However, concerns remain regarding the amplification of biases and the unpredictability of AI-generated outputs. Thus, understanding the behavioral patterns and limitations of AI responses

has become a central research challenge. Against this backdrop, more studies have compared the judgment patterns of humans and AI to identify similarities and differences.

An important question is whether assigning specific personas to generative AI leads to more humanlike outputs. A persona refers to a set of hypothetical individual attributes, such as age, gender, education, income, and cultural background, assigned to an AI agent to generate responses with greater consistency and personality. Many previous studies embed such attribute information directly into the prompt, reporting improved output coherence. Park et al. (2024) use detailed interview data from 1,052 American respondents to test how well AI agents replicate actual human behavior when assigned attributes such as age and political orientation. Their results indicated a matching rate of up to 85%. Jia et al. (2024) assign demographic profiles to GPT models and measure their risk preferences and loss aversion. They find that differences in gender and educational attainment significantly affect the estimated loss aversion coefficients and probability weighting. In summary, the findings suggest that generative AI can replicate humanlike behavioral patterns and provide valuable insights for future empirical and applied research.

Several studies reveal that generative AI shows consistent patterns in gain-and-loss scenarios. That is, while generative AI exhibits risk-averse behavior similar to that of humans, it tends to show weaker loss aversion. Qiu et al. (2024) report that the GPT-4 has a loss aversion coefficient of 1.09, compared to 2.56 for actual consumers, suggesting that the model underestimates losses relative to gains. Jia et al. (2024) find that although generative AI generally behaves in a risk-averse manner, it tends to outweigh the probabilities of rare events. Meanwhile, model-level comparisons reveal several important differences. Ross et al. (2024) demonstrate that GPT-4 generates the most stable and humanlike preference structures compared to other models, such as Claude and Gemini, which exhibit greater variability in their outputs. Chen et al. (2025) highlight the sensitivity of AI preferences to phrasing prompts. Even when presented with identical choice tasks, generative AI produces different responses depending on contextual cues, indicating a lack of internal consistency and a high dependence on prompt design.

Fine-tuning is the process of retraining a pretrained foundation model on additional data related to specific tasks or domains, thereby adjusting its output tendencies and performance to meet objectives. Although many existing studies do not explicitly fine-tune models using human data, they are concerned that reinforcement learning from human feedback (RLHF), often implemented as part of HITL, may significantly affect the behavior of the model. Santurkar et al. (2023) show that large language models (LLMs) tend to produce outputs biased toward liberal, highly educated, and high-income groups, and underrepresent the views of older adults, low-income individuals, and religious communities. The RLHF process may reinforce specific value systems and amplify the underlying biases. Hagendorff et al. (2023) note that although RLHF improves the models' ability to avoid intuitive errors and generate accurate responses, it risks entrenching normative assumptions within the model. In contrast, Park et al.

(2024) adopt a different approach by embedding actual human response data into a prompt to construct individual AI agents. Although this does not constitute formal fine-tuning, it effectively injects human knowledge into the model and functions as an implicit calibration method. These findings suggest that generative AI exhibits humanlike behavioral tendencies when conditioned on social attributes and contextual framing but still lacks consistent judgment across various settings. Although RLHF may enhance coherence and response quality, empirical evaluations of its effects remain limited. Therefore, assessing how well generative AI replicates human decision-making biases before and after fine-tuning represents a key research agenda in both behavioral economics and applied AI.

This study compares the similarities and differences in decision-making biases, specifically risk preferences and loss aversion, between humans and generative AI (GPT-4). Specifically, it analyzes how individual attributes influence these biases across the two entities. Our study designs a unified experimental framework to analyze the presence and direction of biases using equivalent tasks for both humans and AI. We conduct a nationwide online survey of 4,838 Japanese adults and explore the relationship between individual attributes such as age, gender, household income, and their response patterns. Based on these demographic profiles, we construct virtual personas and present the same set of questions to the GPT-4 under multiple temperature settings. We then quantitatively compare the AI-generated responses with human preferences to evaluate both the divergence and degree of alignment in their underlying decision structures.

At this point, *loss aversion* is broadly defined as a property of the value function wherein the slope in the loss domain is steeper than that in the gain domain, typically resulting in risk-averse behavior in the gain domain and risk-seeking behavior in the loss domain⁴. Our empirical analysis reveals that, in the gain domain, humans and GPT-40 exhibit broadly similar patterns of risk aversion. In contrast, in the loss domain, GPT-40 demonstrates significantly stronger risk-seeking tendencies than humans, indicating a divergence in risk attitudes that may stem from differences in the degree or functional form of loss aversion. Moreover, the results suggest that GPT-40 is more sensitive to demographic cues in its decision-making, with individual attributes such as age and income exerting a stronger influence on AI-generated responses than on human choices.

We also fine-tune the AI model using web-based survey data on human loss aversion. Our results show that the AI's choice patterns and marginal effects for demographic attributes become more aligned with human behavior. In particular, in the loss domain, the fine-tuned AI exhibits preference structures that more closely resemble those of humans, positioning it between the pre-trained AI and actual human

⁴ However, our study does not directly estimate differences in the slope of the value function between the gain and loss domains. Rather, we focus on a key behavioral implication of loss aversion—namely, the reversal of risk preferences between these domains—as observed in both human and AI agents. It should be noted, however, that the stronger risk-seeking behavior observed in the loss domain may be attributable to differences in the curvature of the value function.

respondents in terms of behavioral traits. However, in some cases, signs of overfitting or divergence are observed, particularly in terms of accuracy and distributional fit. We also evaluate the distributional similarity using the Wasserstein distance. The fine-tuned AI displays a statistically closer alignment with human choice distributions than the pre-trained model. These results suggest that fine-tuning contributes to the convergence of behavioral distributions between AI and humans.

The remainder of this paper is organized as follows. Section 2 reviews previous studies that compare generative AI with human behavior. Section 3 describes the survey design and data collection procedures. Section 4 introduces the estimation methods for behavioral biases and compares the decision patterns of humans and pretrained AI using the collected data. Section 5 describes the fine-tuning process based on human choice patterns and compares the behavior of the fine-tuned AI with that of humans. Section 6 compares the bias tendencies of humans, pre-trained AI, and fine-tuned AI and summarizes their interrelationships. Section 7 discusses the main findings. Lastly, Section 8 concludes the study and outlines future research directions.

2. Related Literature

We review previous studies that compare decision-making patterns and cognitive biases between generative AI and humans⁵. Table 1 lists the studies based on the following seven dimensions: 1) type of bias or experimental task examined; 2) version and type of generative AI used; 3) target population and sample size; 4) extent to which LLMs replicate average human behavior; 5) replicability of human responses when conditioned on socioeconomic attributes; 6) whether the model incorporates human feedback, such as fine-tuning; and 7) metrics used to evaluate differences in response distributions between humans and LLMs.

⁵ Several studies examine the research applications of LLMs, outlining both their potential and limitations. Demszky et al. (2023) provide a systematic overview of how LLMs can be applied in psychology and highlight the risks of biased outputs when cultural biases or underrepresented groups are not adequately captured in the training data. Hagendorff (2023) introduces the concept of "machine psychology," proposing a framework that treats LLMs as virtual participants in psychological experiments and draws attention to the ethical concerns associated with this approach. Sarstedt et al. (2024) review studies comparing so-called "silicon samples" with human respondents to assess the applicability of LLMs in consumer behavior and marketing research. They find that LLMs can achieve high replicability in tasks such as brand evaluation and framing effects but also note their limitations in replicating preference formation and behavioral decision-making tasks. These reviews offer valuable guidance for incorporating LLMs into social science research and provide critical perspectives on ethical governance and future application areas.

< Table 1: Prior Studies on Decision-Making and Cognitive Biases: Generative AI vs. Humans >

We introduce studies that compare human decision-making with that of large LLMs through economic experiments. We first focus on game-theoretic experiments that involve social preferences such as altruism and strategic interactions. Although these studies do not directly compare AI behavior with human data, they evaluate LLM behavior in economic games by allowing multiple models to interact with one another (e.g., Akata et al., 2025; Guo, 2023; Kitadai et al., 2024; Tsuchihashi, 2023; Lorè and Heydari, 2024). Akata et al. (2025) conduct repeated games using GPT-3, GPT-3.5, and GPT-4 in scenarios such as the Prisoner's Dilemma and the Battle of the Sexes, each with different payoff matrices. They show that GPT-4 exhibits humanlike cooperative behavior but frequently struggles to maintain implicit coordination. However, coordination success rates significantly improve when a "cooperative persona" is introduced in the prompt. Guo (2023) assigns GPT-4 personas such as "fair" and "selfish" and repeatedly runs ultimatum games and the Prisoner's Dilemma. They show that the "fair" persona generates offer and rejection thresholds close to human averages, while the "selfish" persona leads to low offers and high acceptance rates. Kitadai et al. (2024) conduct simulations of the ultimatum game using GPT-3.5 and GPT-4. They demonstrate that enhancing the reasoning abilities of GPT-based generative agents yields outcomes closer to theoretical predictions than to actual human experimental data. Lorè and Heydari (2024) present four two-player games under five contextual framing conditions and compare the behavior of GPT-3.5, GPT-4, and LLaMa-2. They reveal that GPT-4 tends to select strategies consistent with game-theoretic rationality depending on the game structure. In contrast, GPT-3.5 and LLaMa-2 exhibit greater sensitivity to contextual framing. Tsuchihashi (2023) examines GPT-3.5's bidding behavior in sealed-bid auctions. They find that in first-price auctions (FPA), the model overbids similarly to humans, whereas in second-price auctions (SPA), it submits truthful bids in line with theory. Furthermore, when given a "student persona," GPT-3.5 places more theoretically consistent bids in FPA and underbids in SPA, suggesting that persona conditioning affects strategic behavior.

Several studies compare the behavior of LLMs with human decisions through experiments on social preferences (e.g., Horton, 2023; Brookins and DeBacker, 2024; Mei et al., 2024; Xie et al., 2024). Horton (2023) conceptualizes LLMs as Homo Silicus, a virtual economic agent, and presents GPT-3 with standard behavioral economics tasks, such as fairness evaluations, status quo bias, and responses to minimum wage policies. They show that the model exhibits biases consistent with those observed in previous human experiments. Moreover, its decision patterns vary depending on the assigned persona, such as selfish, fair-minded, or efficiency-oriented. Brookins and DeBacker (2024) repeatedly run dictator and prisoner's dilemma games using GPT-3.5 to assess its tendency toward fairness and

cooperation. They find that the model tends to behave more altruistically than humans, particularly in contexts that emphasize efficiency, in which its behavior aligned more closely with human patterns. Mei et al. (2024) analyze how GPT-3.5 and GPT-4 replicate human behavior by comparing model outputs with a dataset of approximately 90,000 human decisions across six game types. Using a Turing test framework, they find that GPT-4's responses closely resemble the human average. However, the model also displays a centralizing tendency, avoiding extreme responses. Xie et al. (2024) evaluate five models, including GPT-40, using games such as the dictator and ultimatum game. They assess the differences in response distributions between LLMs and humans using the Wasserstein distance. Overall, LLMs demonstrate fairer and more cooperative behavior than humans, with GPT-40 producing the distribution most similar to human responses.

Meanwhile, several experimental studies examine loss aversion and risk preferences in LLMs (Jia et al., 2024; Qiu et al., 2024; Ross et al., 2024; Chen et al., 2025; Macmillan-Scott and Musolesi, 2024). Jia et al. (2024) present choice tasks related to loss aversion and probability weighting to GPT-4, Claude, and Gemini. They show that all models exhibit risk-averse behavior, with GPT-4 demonstrating lossaversion patterns similar to those of humans. They also report substantial variations in output depending on the prompt conditions based on demographic attributes. Ross et al. (2024) comprehensively analyze 12 LLMs, including GPT-4, focusing on biases such as inequality aversion, loss aversion, and time discounting. The models reveal intense guilt and weak envy in terms of social preferences. While they respond rationally to gains, their behavior under losses deviates from rational expectations. Additionally, the LLMs display stronger time discounting than human subjects. Chen et al. (2025) examine 18 cognitive biases, including risk attitudes. They show that GPT-3.5 tends to be risk-seeking for gains and risk-averse for losses, while GPT-4 exhibits more consistent risk aversion and lower susceptibility to framing effects. Macmillan-Scott and Musolesi (2024) evaluate 12 cognitive bias tasks, including the conjunction and gambler's fallacies. They show that GPT-3.5 reproduces many humanlike biases, whereas GPT-4 achieves a higher overall accuracy. Finally, although not based on an economic experiment, Qiu et al. (2024) simulate health insurance plan choices (n = 5.998) using GPT-4 to assess whether the model could replicate human decisions under risk. They show that while the aggregate choice distributions resemble those of humans, individual-level agreement is low. The estimated loss aversion coefficient for GPT-4 (1.09) is substantially lower than the human average (2.56), suggesting that the model systematically underweights losses relative to human decision-makers.

Other studies examine the reproducibility of survey responses using generative AI. Santurkar et al. (2023), Park et al. (2024), and Dominguez-Olmedo et al. (2025) compare AI-generated outputs with actual public opinion and survey data to evaluate the degree of alignment and divergence. In the domain of cognitive biases, some studies investigate the extent to which LLMs replicate cognitive errors and reflective thinking patterns (Aher et al., 2023; Binz and Schulz, 2023; Hagendorff et al., 2023). Additionally, research on political attitudes and voting behavior analyze the reproducibility of AI

outputs under specific political ideologies and demographic conditions (Argyle et al., 2023; Motoki et al., 2024; Bisbee et al., 2024). Lastly, additional studies explore the behavioral alignment between AI and humans in various domains, covering topics such as market research, economic forecasting, environmental awareness, and tourism behavior (Brand et al., 2023; Bybee, 2023; Lee et al., 2024; Li et al., 2024; Xiong et al., 2024).

In summary, these studies demonstrate that generative AI can replicate human behavior and judgment under specific conditions. Moreover, they highlight that such reproducibility heavily depends on model configurations, including demographic prompts, contextual framing, question formats, and the design of personas or instructions.

Recent LLMs, such as GPT-4, exhibit a strong ability to replicate average tendencies and aggregatelevel human judgments with high accuracy. However, they show apparent limitations in reproducing extreme responses or generating outputs based on specific demographic attributes. Issues such as "hyperaccuracy bias" and excessive internal consistency may lead to unnaturally uniform or nonhuman-like outputs. Moreover, LLMs are prone to systematic bias when responding to prompts involving socially marginalized or culturally underrepresented groups.

Many prior studies comparing human and generative AI decision-making do not clarify the evaluation metrics used to assess reproducibility. Notably, only Santurkar et al. (2023) and Xie et al. (2024) employ the Wasserstein distance, a metric adopted in this study. Other studies rely on Kullback–Leibler (KL) divergence to compare distributions. However, the KL divergence may fail to capture differences in the overall shape of the distributions, making it inadequate for identifying substantive disparities between AI-generated and human response patterns. Therefore, using the Wasserstein distance is more appropriate.

In addition, although many studies evaluate the average-level reproducibility of AI responses, few consider individual attributes. Even when generative AI exhibits behavior similar to that of humans, divergent response patterns may emerge under more granular conditions or when evaluations are made at the individual level. Although incorporating human responses into AI training may improve reproducibility, only a few studies have implemented such a feedback mechanism. Given these gaps, this study offers a more rigorous framework for evaluating the similarities and differences between generative AI and human behavior. Compared to prior work, our approach provides advantages in terms of the validity of its analytical methods and precision of reproducibility assessments.

3. Survey Design

3.1. Design of the Web-Based Survey

We conduct a web-based survey (December 2024) through an online survey company, targeting domestic residents aged 20–65. We divide Japan into nine regions and allocate respondents to balance

gender and age groups within each region. Participants are informed in advance that the survey aims to investigate the relationship between consumer attributes and preferences. They would receive compensation for their participation in the study. We obtain responses from 5,040 participants. After excluding 202 incomplete responses, the final sample includes 4,838 participants. The survey includes individual attributes, such as gender, age, and household income. It also contains multiple questions on cognitive biases and psychological characteristics that cover Big Five personality traits, time preferences, ultimatum games, and the trolley problem. In this study, we focus on gender, age, and household income as individual attributes. Regarding psychological characteristics, we focus only on risk preferences and loss aversion based on prospect theory. Table 2 presents descriptive statistics of the respondents' attributes⁶.

< Table 2: Descriptive Statistics >

3.2. Questions on Prospect Theory

In this section, we discuss the questionnaire used to measure risk preferences and loss aversion. As a reference, the following choice tasks concerning gains and losses are presented to students at Stanford University and the University of British Columbia by Tversky and Kahneman (1988):

Problem 1. (n=126)

Assume yourself richer by \$300 than you are today. You are offered a choice between

A. A sure gain of \$100, or

B. A 50% chance to gain \$200 and A 50% chance to lose \$0.

Problem 2. (n=128)

Assume yourself richer by \$500 than you are today. You are offered a choice between

- A. A sure loss of \$100, or
- B. A 50% chance to lose \$200, and A 50% chance to lose \$0.

Participants are presented with two choices (A and B) with equivalent expected values in both the gain and loss domains. In Problem 1, which corresponds to the gain domain, 72% of participants select the sure gain option (A), whereas 28% choose the risky option (B). This indicates a risk-averse preference in the gain domain, in which most participants prefer certainty over risk. However, in Problem 2, which

⁶ Income is categorized into six groups: "Less than JPY 3 million," "JPY 3 million to less than JPY 5 million," "JPY 5 million to less than JPY 7 million," "JPY 7 million to less than JPY 10 million," "JPY 10 million to less than JPY 15 million," and "JPY 15 million or more." Class values are set as follows: JPY 3 million, JPY 4 million, JPY 6 million, JPY 8.5 million, JPY 12.5 million, and JPY 15 million.

corresponds to the loss domain, 36% of the participants select the sure option (A), while 64% choose the risky option (B). This suggests a risk-seeking tendency in the loss domain, in which participants gamble on the possibility of avoiding loss by accepting a higher potential loss.

Another important point is the difference in the choice tendencies between Problems 1 and 2. Although both problems present identical monetary amounts and probability structures from the same wealth reference point, the observed choices differ. This preference asymmetry reflects framing effects consistent with loss aversion⁷. This tendency highlights a behavioral pattern that deviates from the rational decision-making predicted by the standard expected utility theory, showing that human decisions are heavily influenced by frameworks such as reference points. Based on Tversky and Kahneman (1988), we present the following questions:

Question 1. Imagine you receive an additional JPY 30,000 on top of your current wealth and are asked to choose between the options below. Which option would you choose?

Option 1: Receive a guaranteed JPY 10,000

Option 2: A 50% chance of receiving 20,000 yen and a 50% chance of receiving nothing

Question 2. Imagine you receive an additional JPY 50,000 on top of your current wealth and are asked to choose between the options below. Which option would you choose?

Option 1: Lose 10,000 yen with certainty

Option 2: A 50% chance of losing JPY 20,000 and a 50% chance of losing nothing

Since our survey targets domestic residents in Japan, questions are presented in the Japanese language, and monetary amounts are displayed in JPY (with an exchange rate of USD 1 = JPY 100).

3.3. Survey Design for Generative AI

We implement "personas" constructed from actual web-based survey attribute information into generative AI and compare the AI's decision tendencies with those of humans under identical conditions. Specifically, we extract representative patterns of attributes from human data and input them into

⁷ Kahneman and Tversky (1979) reports a similar tendency. The structure of the questions and the results in their study are as follows. In PROBLEM 1, the proportions of participants choosing options A and B are 16% and 84%. In contrast, in PROBLEM 2, the proportions of participants choosing options C and D are 69% and 31%, respectively.

Problem 1.

In addition to whatever you own, you have been given 1,000. You are now asked to choose between: A: (\$1,000, 0.50), and B: (\$500).

Problem 2.

In addition to whatever you own, you have been given 2,000. You are now asked to choose between: C: (\$-1,000, 0.50), and D: (\$-500).

prompts to assign virtual personas to the AI. Then, we present the same questions as in the web-based survey to these personas and collect their responses, enabling a condition-controlled comparative analysis. We use OpenAI's GPT-40 as the generative AI model. Unlike the general ChatGPT interface, GPT-40 is accessible via the OpenAI API, allowing program-based operation. Using the API, we can automate large-scale response generation with Python code, enabling systematic analysis under numerous persona conditions⁸.

The process of collecting responses from generative AI is as follows. First, we extract the three attributes (gender, age, and household income) from web-based survey data and save them in CSV format. Second, we create a Python script that references this CSV file and links it to the OpenAI GPT-40 model to automate the response generation process (see Figure 1). Lines 1–5 of the script import necessary libraries (e.g., pandas, csv, and openai) for reading CSV files and communicating with the OpenAI API. Line 6 sets the API key (authentication information) to enable the use of the OpenAI API. Lines 9–10 specify the path to the CSV file containing the analysis targets. Lines 12–24 define the function "extract_reason_and_answer" to extract "reason" and "answer" from the generated responses, outputting "unknown" if no reasoning is provided. Lines 26–53 define the function "generate_reason_and_answer", which constructs prompts based on each persona's attributes (gender, age, household income) and generates AI responses via the API. Line 55 applies this process to all 4,838 samples on the DataFrame, storing the extracted reasons and answers as separate variables. The subsequent code saves the response data as CSV files for further analysis.

<Figure 1: Python Script for AI Question Presentation and Response Collection>

To control variability in the generated responses, we also adjust the "temperature" parameter, which controls the randomness, and collect responses under three different temperature settings⁹. By adjusting the temperature, we could evaluate the impact of probabilistic fluctuations on AI decision tendencies. Through these procedures, we reproduce 4,838 personas with attributes identical to those of the webbased survey participants on GPT-40 and obtain responses to the two-choice tasks (see Table 3)¹⁰. Then, we apply a logit model to the response data and quantitatively analyze the generative AI's tendencies regarding risk preferences and loss aversion.

⁸ In addition to GPT-40, we also collect responses to the same questions from other chat-based language models, including OpenAI's GPT-3.5, Google's Gemini 2.0 Flash, and DeepSeek's DeepSeek-R1. The results of these comparisons are discussed in Section 7.

⁹ The temperature parameter ranges from 0 (minimum) to 2 (maximum), with a default value of 1. Higher temperatures lead to increased randomness in the generated text. However, when the temperature exceeds 1, the analysis results become unstable in our setting.

¹⁰ Unlike the general ChatGPT interface, responses generated via the OpenAI API are stateless and independent, ensuring response independence comparable to that in web-based surveys.

<Table 3: Example of Responses Using OpenAI API>

3.4. Estimation Methods for Behavioral Biases

We present the estimation method used to clarify the impact of respondent attributes on risk preferences and loss aversion.¹¹ In the questionnaire, respondents are asked to choose between two options: "Option 1: JPY 10,000 with 100%" or "Option 2: JPY 20,000 with 50%." Therefore, the responses can be treated as binary variables. In this study, the binary responses are treated as dependent variables, and the three attributes are treated as independent variables. The estimation is conducted using the maximum likelihood method. Suppose that the error terms associated with the utility of each choice follow a type-I extreme value distribution. In this case, the difference in the error terms follows a logistic distribution; thus, using a logit model is appropriate. The choice probability in the logit model is as follows:

$$P(Y_{i} = 1 | X_{1i}, X_{2i}, X_{3i}) = \frac{1}{1 + e^{-(\beta_{0} + \beta_{1}X_{1i} + \beta_{2}X_{2i} + \beta_{3}X_{3i})}}$$

where i denotes each respondent. The dependent variable Y_i equals 1 if the respondent selects "Option 1: JPY 10,000 with 100%" and 0 if they select "Option 2: JPY 20,000 with 50%." X_{1i} is a dummy variable for the gender of respondent i that takes the value of 1 if the respondent is female, and 0 if the respondent is male. X_{2i} represents the respondent i's age in 10-year increments. X_{3i} indicates the respondent i's household income in 1-million-JPY increments.

4. Comparison of Biases between Humans and Pre-Trained AI¹²

4.1. Overview of Choice Patterns

Table 4 presents the choice results for each question obtained from human respondents and the generative AI. For each question, a binary variable is defined: Option 1 (a certain gain/loss of 10,000 yen) is coded as 1 and Option 2 as 0. The average of this variable represents the average selection rate of Option 1.

<Table 4: Choice Results: Human vs. Pre-trained AI>

¹¹ From this point forward, we refer to both the web-based survey participants and generative AIs with implemented personas as "respondents." The terms "individual" and "answer" are used in the same unified manner.

¹² Generative AI in this section is not yet fine-tuned, unlike what is introduced after Section 5. Therefore,

We first review the results of the web-based survey. For Question 1, 88.4% of the human respondents select Option 1 (a sure gain of JPY 10,000), whereas only 11.6% choose Option 2 (a 50% chance of gaining JPY 20,000 or nothing). This result suggests a strong risk-averse tendency in the gain domain. For Question 2, 57.7% choose Option 1 (a sure loss of JPY 10,000), while 42.3% choose Option 2 (a 50% chance of losing JPY 20,000 or nothing). Although the result indicates some degree of risk-seeking behavior, no clear tendency is found in the loss domain. Then, we examine the response results from the generative AI under the default setting of temperature = 1.0. For Question 1, 91.1% of the respondents select the sure gain option, indicating a clear tendency toward risk aversion in the gain domain. However, for Question 2, only 11.0% choose the sure loss option, whereas the majority opt for the probabilistic loss, indicating a distinct risk-seeking tendency in the loss domain. In summary, the results suggest that humans and generative AI display a high level of risk aversion in the gain domain. In contrast, in the loss domain, while the generative AI exhibits a clear risk-seeking tendency, the human responses do not display a similarly strong pattern. This discrepancy suggests that preferences in the loss domain may differ between humans and AI. When we compare the selection rates for Option 1 in Questions 1 and 2, the generative AI results are 91.1% and 11.0%, indicating stronger loss aversion. For humans, the rates are 88.4% and 57.7%, reflecting relatively weaker loss aversion compared to the AI.

Lastly, we examine how the temperature setting affects the response tendencies of the AI. For Question 1, selection rates for Option 1 remained high across different temperatures (91.1% (1.0), 88.9% (0.5), and 89.8% (0.0)), confirming risk aversion in the gain domain. However, for Question 2, selection rates for Option 1 decrease as the temperatures decrease (11.0% (1.0), 3.0% (0.5), and 0.0% (0.0)), suggesting a stronger risk-seeking tendency at lower temperatures. Notably, at temperature = 0.0, only one respondent (persona is "female, age: 53, income: JPY 3 million") selects Option 1. As this response deviates significantly from the broader trend, the data is excluded from subsequent analyses.

4.2. Analysis Based on Marginal Effects

Tables 5 and 6 present the estimated marginal effects and average predicted probabilities derived from the logit model using datasets from both humans and generative AI across different temperature settings. The marginal effects indicate how the likelihood of $Y_i = 1$ changes with a one-unit increase in each explanatory variable. The average predicted probabilities are calculated by applying the estimated logit model to each sample to obtain the probability of choosing Option 1 and then averaging those probabilities. These values align with the selection rates presented in Table 4, confirming that the logit model accurately predicts actual choice tendencies.

<Table 5: Logit Model Estimation: Humans vs. Pre-trained AI (Gain, 3 Attributes)> <Table 6: Logit Model Estimation Humans vs. Pre-trained AI (Loss, 3 Attributes)> We examine the estimation results for humans. For Question 1, all attributes show statistically significant marginal effects at the 1% level. Specifically, being female increases the probability of choosing Option 1 by 4.1%. A 10-year age increase raises the probability by 1.5%, whereas a JPY 1 million increase in household income reduces the probability by 0.5%. For Question 2, statistically significant marginal effects are found only for age and income, but not for gender. Specifically, a 10-year age increase raises the probability of choosing Option 1 by 2.1% and a JPY 1 million increase in income decreases the probability by 0.8%.

Subsequently, we analyze the generative AI results at a temperature of 1.0. In this setting, statistically significant marginal effects at the 1% level are observed for all three attributes for both questions. For Question 1, being female increases the probability by 3.3%, a 10-year age increase raises the probability by 5.7%, whereas a JPY 1 million increase in income decreases the probability by 2.7%. For Question 2, being female increases the probability by 2.8%, a 10-year age increase raises the probability by 2.3%, whereas a JPY 1 million increase in income decreases the probability by 2.3%, whereas a JPY 1 million increase in income decreases the probability by 0.6%.

We then compare the marginal effects of humans and AI at a temperature of 1.0. Using the delta method, we calculate the standard errors for the marginal effects and conduct statistical tests to determine the differences between the two groups. For Question 1, the z-values for the differences in marginal effects between humans and AI for gender, age, and income are 0.70, -9.90, and 15.56, respectively. The difference in gender is not statistically significant at the 5% level. In contrast, significant differences are observed in terms of age and income. For Question 2, the z-values for the differences in marginal effects for gender, age, and income are -2.04, -0.31, and -0.71, respectively. Unlike in the gain domain, only gender shows a statistically significant difference, while age and income do not show any significant differences. These results suggest that while partial differences in marginal effects between humans and AI are observed, no consistent or statistically significant differences are found.

Finally, we examine how changes in the temperature setting affect the marginal effects of AI. In all temperature conditions, statistically significant marginal effects at the 1% level are observed for all attributes for both questions. For Question 1, the marginal effect of age is consistently larger for AI than for humans and increases with lower temperatures (0.057 (temp=1.0), 0.073 (temp=0.5), and 0.120 (temp=0.0)) compared with the result for humans (0.015). For income, while the negative marginal effects of AI are larger than those of humans (-0.005), a consistent trend is not found across the different temperature settings (-0.027 (temp=1.0), -0.033 (temp=0.5), and -0.031 (temp=0.0)). For gender, no consistent trend is observed for temperature (Human: 0.041; AI: 0.033 (temp = 1.0), 0.026 (temp = 0.5), and 0.086 (temp = 0.0)). For Question 2, no clear temperature-dependent trend is found in the marginal effect of age (Human: 0.021; AI: 0.023 (temp = 1.0), 0.008 (temp = 0.5)). For income, AI shows slightly stronger negative marginal effects than the human result (-0.008); however, no consistent trend across temperature settings is found as well (both temps: -0.006). For gender, no variation is observed across

temperature settings (both temps: 0.028).

4.3. Graphical Analysis of Choice Probabilities

Figures 2 and 3 illustrate the choice probabilities for each attribute in Questions 1 and 2, respectively. Figure 2 shows that under "Gender," the other two attributes are fixed at their average values (age = 44.9 years, income = JPY 5.93 million). The choice probabilities are plotted for the gender dummy variable with values of 0 (male) and 1 (female). Additionally, in "Gender," we observe that the 95% confidence intervals overlap at both points. This indicates no significant difference in choice probabilities between the web-based survey results and AI responses, or among the different temperature conditions.

> <Figure 2: Choice Probabilities: Humans vs. Pre-trained AI (Gain, 3 Attributes)> <Figure 3: Choice Probabilities: Humans vs. Pre-trained AI (Loss, 3 Attributes)>

Meanwhile, in "Age" and "Income," the slope of the black line representing the human web-based survey results is relatively gentle, whereas the lines for the generative AI are steeper. This suggests that AI emphasizes changes in the choice probability more strongly in response to age or income. However, the confidence intervals overlap around the mean values in these graphs, indicating that the differences between humans and AI are relatively small near the average. Figure 3 illustrates a clear difference in the y-intercept across all three attributes. The y-intercepts for human web-based survey results lie around 0.5–0.6, whereas those for the AI results are clustered around 0–0.1. This indicates notable differences in choice probabilities at lower age and income levels. However, the slopes of the lines do not differ significantly between humans and AI, suggesting that incremental changes in the choice probability with increasing age or other attributes are similar.

In addition, by comparing the scale of the y-coordinates in both figures, we can assess the tendency toward loss aversion. In Figure 2, under "Age," the AI's choice probability at age 20 (with temperature = 1.0) is around 0.7. In contrast, in Figure 3, under "Age," the choice probability under the same condition is only about 0.05. This sharp contrast reflects the previously observed strong tendency towards loss aversion in generative AI. Moreover, the figures illustrate that while humans show relatively weaker loss aversion, AI displays this behavior more prominently.

4.4. Analysis of Prediction Accuracy

Table 7 shows the prediction accuracy between the predicted choice probabilities from each model and the actual binary choices. "Predicted value" is defined as a binary classification based on the expected probability $\hat{P}(Y_i = 1 | \mathbf{X})$. If the predicted probability is 0.5 or greater, the respondent is classified as choosing Option 1; otherwise, they are classified as choosing Option 2. Based on this classification, we

calculate the match rate between the model's prediction and the actual observed response to assess the model's fit. In this context, "prediction accuracy" refers to the proportion of cases in which the predicted outcomes from the model match the actual observed outcomes. Specifically, "prediction accuracy" is the sum of the proportions of cases in which both the predicted and observed outcomes are Option 1 (i.e., (1,1)) and those in which both are Option 2 (i.e., (2,2)). In the human web-based survey results, the prediction accuracy for Question 1 is 88.40%. However, for Question 2, the prediction accuracy decreases to 57.59%. In comparison, the prediction accuracy based on the responses generated by the AI ranges from 88.98% to 96.96%, indicating a higher overall level of consistency compared with the web-based survey.

<Table 7: Logit Model Prediction Accuracy: Humans vs. Pre-trained AI>

5. Fine-Tuning

This study also aims to replicate human behavioral biases within generative AI by training it on real human response data and to evaluate the proximity between humans and AI before and after training. Thus, we employ fine-tuning as the training method for the generative AI. Fine-tuning refers to the retraining of a pre-trained model, originally trained on large-scale data, on a new dataset to specialize in a specific task or domain (Devlin et al., 2019). In the case of OpenAI's GPT model, fine-tuning is performed by preparing a dataset in the JSONL format, uploading it to the OpenAI API system, and initiating the training process.

We use the predicted probabilities estimated by the logit model for each respondent as the training data. These probabilities represent $\widehat{P}(Y_i = 1 | \mathbf{X})$, which is the likelihood of choosing Option 1 calculated from the estimated model for each individual. After completing the fine-tuning process, we submit the same questions to the fine-tuned model as we have for the human participants and the pre-trained AI. We estimate the predicted probability using a logit model to clarify the differences in behavioral tendencies among the three groups.

Figure 4 shows an excerpt from the JSONL format dataset used for fine-tuning. In fine-tuning with GPT, each training sample must consist of a sequence of 10 or more dialogues using three roles: "system", "user", and "assistant". The "system" defines the instructions the model should follow. The "user" provides the input or question to the model. The "assistant" provides the model's expected response. In our dataset, we input virtual personal attributes (personas) under "system", questions about the respondent's predicted choice probabilities under "user", and the value of the predicted probabilities under "assistant".

In the fine-tuning, determining the training methods and parameter settings in advance is necessary. For the training methods, we consider two possible approaches: a probability-based method, which directly inputs the value of the predicted probability as reported by the human respondent, and a frequency-based method, which involves inputting question-and-answer data that reflects the frequency of human choices. For the parameters, we consider epoch and learning rate multiplier (LR). Epoch is the number of training iterations. The higher the value, the more the model learns; the risk of overfitting increases. LR is used to adjust the learning speed. A higher value leads to faster and more efficient learning and easier convergence; however, it also increases the risk of overfitting. We conduct preliminary fine-tuning using a sample of 100 respondents to select the appropriate learning method and parameter settings. Then, we instruct the fine-tuned model to generate responses to the two questions and evaluate them using logit analysis. With OpenAI's default settings (epoch = 3, LR = 2), nearly all samples selected Option 1 in Question 1, resulting in an output that may reflect overfitting. This result indicates overfitting. Therefore, we conduct multiple trials with lower epoch and LR values. After comparing the results, we adopt the probability-based method as our training approach, setting the parameters to epoch = 1 and LR = 0.2.

<Figure 4: Dataset for Fine-tuning>

Figure 5 shows the Python script used to train the GPT-4o-2024-08-06 model. Lines 1–3 of the script, "import ...", read JSONL files and communicate with the OpenAI API. Line 9 specifies the path to the JSONL file containing analysis targets. Lines 12–18 read the JSONL training file and prepare it for uploading to the model. Lines 21–31 initiate the fine-tuning of GPT-4o-2024-08-06 with epoch=1 and LR=0.2. Lines 33–36 instruct the output of the progress status during fine-tuning.¹³

<Figure 5: Python Script for Fine-tuning>

We apply the same procedure as the one for the pretrained AI (Section 3.3) to the fine-tuned generative AI. Specifically, we reproduce 4,838 personas on the generative AI, each with the same attributes as those in the web-based survey results and obtained responses to two-choice questions. The resulting responses are analyzed using the logit model to examine the tendencies in risk preferences and loss aversion.

6. Comparison of Biases Among Humans, Pre-Tuned AI, and Fine-Tuned AI

We compare the logit model estimation results by persona based on individual attributes, examining the

¹³ In fine-tuning, OpenAI system randomly shuffles the examples repeatedly, so the order of examples should not make a difference. Namely, each sample contributes independently to the learning.

outcomes for humans, pretrained AI, and fine-tuned AI. All generative AI responses are obtained using the default temperature setting of 1.0.

6.1. Overview of Choice Patterns

Table 8 presents the choice results for each question obtained from human respondents and the pre- and post-trained generative AI. For Question 1, respondents choosing Option 1 (96.1%) significantly exceed those who choose Option 2 (3.9%). This suggests a strong tendency toward risk aversion in the gain domain. Conversely, for Question 2, more respondents select Option 2 (59.6%) over Option 1 (40.4%). This indicates a tendency toward risk-seeking in the loss domain; this tendency is weaker than the risk aversion observed in the gain domain.

<Table 8: Choice Results: Humans, Pre-trained AI, and Fine-tuned AI>

Next, we compare the results of humans, pre-trained AI, and fine-tuned AI. For Question 1, the finetuned AI's selection rate (96.1%) is higher than that of the pre-trained AI (91.1%), which has a higher rate than the human average (88.4%). This suggests overfitting in the form of excessive risk aversion in the gain domain. Similarly, for Question 2, the fine-tuned AI's average selection rate (40.4%) is closer to that of humans (57.7%), indicating that it appropriately learns human tendencies toward risk-seeking in the loss domain. Comparing the selection rates for Option 1 between Questions 1 and 2, the finetuned AI rates for Questions 1 and 2 (96.1% and 40.4%, respectively) are higher than those of the pretrained AI (91.1% and 11.0%, respectively), indicating that the fine-tuned AI has a weaker loss aversion and a pattern similar to that of humans (88.4% and 57.7%, respectively).

6.2. Analysis Based on Marginal Effects

Tables 9 and 10 present the estimated marginal effects and average predicted probabilities derived from the logit model using datasets from the human web-based survey and the pre- and post-trained generative AI. The average predicted probabilities align with the selection rates presented in Table 8, confirming that the logit model accurately predicts actual choice tendencies.

<Table 9: Logit Model Estimation: Humans, Pre-trained AI, and Fine-tuned AI (Gain, 3 Attributes)> <Table 10: Logit Model Estimation: Humans, Pre-trained AI, and Fine-tuned AI (Loss, 3 Attributes)>

First, Question 1 confirms statistically significant marginal effects at the 1% level for all three attributes. Specifically, being female increases the probability of choosing Option 1 by 1.9%. A 10-year age increase raises the probability by 3.0%, and a JPY 1 million increase in income decreases the probability by 1.4%. By contrast, statistically significant marginal effects are observed only for gender

and age in Question 2. Specifically, being female increases the probability of choosing Option 1 by 26.4%, and a 10-year age increase raises the probability by 13.3%. For Question 2, large marginal effects are observed for both gender and age. This indicates that the impacts of gender and age on risk-seeking behavior in the loss domain become stronger through training.

Next, we compare the responses of humans and those of fine-tuned AI. Using the delta method, we calculate standard errors for the marginal effects and conduct statistical tests to determine the differences between them. For Question 1, the z-values for the differences in the marginal effects for gender, age, and income are 2.14, -4.16, and 6.36, respectively; all are statistically significant. For Question 2, the z-values for the differences for gender, age, and income are -15.2, -17.5, and -1.77. Statistically significant differences at the 5% level are found only for gender and age.

6.3. Graphical Analysis of Choice Probabilities

Figures 6 and 7 illustrate the choice probabilities for each attribute for Questions 1 and 2, respectively. In Figure 6, in "Gender," the 95% confidence intervals for humans' and pre-trained AI's lines overlap at both points, indicating no significant difference in the choice probabilities. However, fine-tuned AI's line is positioned above both, suggesting a stronger risk aversion through training.

<Figure 6: Choice Probabilities: Humans, Pre-trained AI, and Fine-tuned AI (Gain, 3 attributes)><Figure 7: Choice Probabilities: Humans, Pre-trained AI, and Fine-tuned AI (Loss, 3 attributes)>

In "Age" and "Income," the slope of the fine-tuned AI's line becomes more moderate through training, approaching that of humans. Nevertheless, near the average values, the fine-tuned AI lies above the 95% confidence intervals for both humans and the pre-trained AI, again indicating a heightened risk aversion.

Figure 7 shows a different trend. In "Gender" and "Age," the fine-tuned AI reflects high positive marginal effects with steep slopes. Thus, when the persona is male or young, its line gets closer to that of the pre-trained AI, whereas for female or older personas, it aligns more closely with human responses. Around the average values, the y-coordinates of the fine-tuned AI fall between those of humans and the pre-trained AI, indicating that it has effectively learned humanlike risk seeking in the loss domain. In "Income," no major slope differences are found among the three lines, but apparent differences in y-intercepts are observed. The fine-tuned AI's line is positioned between the other two and closer to the human line. In addition, by comparing the scale of the y-coordinates in both figures, we can assess the tendency toward loss aversion. Notably, we show that training reduces the asymmetric response pattern in gain and loss domains, leading to behavior more consistent with human-like loss aversion.

6.4. Analysis of Prediction Accuracy

Table 11 shows the prediction accuracy of the choice probabilities for each model. For Question 1, the fine-tuned AI's prediction accuracy is high (96.57%), surpassing that of both humans (88.40%) and the pre-trained AI (92.08%). This indicates that the fine-tuned AI accurately replicates human decision tendencies in the gain domain. In contrast, for Question 2, the prediction accuracy of the fine-tuned AI (69.93%) is higher than that of humans (57.59%) but lower than that of the pre-trained AI (88.98%). This suggests limitations in the ability of fine-tuned AI to replicate human behavior in the loss domain.

<Table 11: Logit Model Prediction Accuracy: Humans, Pre-trained AI, and Fine-tuned AI>

6.5. Comparative Analysis of Choice Probabilities Using the Wasserstein Distance

For the quantitative evaluation of the proximity of decision tendencies between humans and generative AI before and after training, we utilize the Wasserstein distance of order 2, a metric for measuring the distances between probability distributions. The Wasserstein distance of order 2 is defined based on optimal transport theory and reflects the minimum cost required to move mass between two distributions. Specifically, using the squared Euclidean distance as the cost function, it is formalized as follows (Ambrosio et al., 2008):

$$W_2(\mu^1, \mu^2) \coloneqq \min\{\int_{\mathbf{X}^2} d(x_1, x_2)^2 d\mu(x_1, x_2) : \mu \in \mathcal{T}(\mu^1, \mu^2)\}.$$

Here, μ^1 and μ^2 are arbitrary probability distributions over a probability space, $\mathcal{T}(\mu^1, \mu^2)$ denotes the set of joint probability distributions with μ^1 and μ^2 as their marginal distributions, and $d(x_1, x_2)$ represents the Euclidean distance between two points x_1 and x_2 . The right-hand side of the equation represents the minimum transportation cost among all transport plans from the distribution μ^1 to μ^2 , with the cost function defined as $d(x_1, x_2)^2$. Using this distance metric enables a quantitative assessment of the difference between the choice distributions of pre- and post-trained generative AI and those of humans in terms of the distance between their respective choice probability distributions.

Table 12 shows the Wasserstein distances of order 2 between the choice probability distributions of humans and those of the pre-trained and fine-tuned AIs. In calculating the Wasserstein distance, we used the distributions of predicted probabilities estimated by the logit model for humans, as well as for the pre-trained and fine-tuned AI.

<Table 12: Wasserstein Distance of Order 2 between Three Groups>

For Question 1, the distances between humans and the pre-trained AI, between humans and the finetuned AI, and between the pre-trained AI and fine-tuned AI are 0.135, 0.123, and 0.081, respectively. This indicates that the human-AI distance is reduced by approximately 8.9% (0.012 units). Moreover, the fine-tuned AI is 52% (=0.123/0.081-1) closer to the pre-trained AI than to humans. For Question 2, the distances between humans and the pretrained AI are 0.467, 0.263, and 0.353, respectively. This indicates that training reduces the distance between humans and AI by 0.204 (approximately 43.7%). Additionally, the fine-tuned AI is 34% (=1-0.263/0.353) closer to humans than the pre-trained AI. Using the combined data from Questions 1 and 2 (N=9,676), the distances between humans and the fine-tuned AI are 0.343, 0.198, and 0.253, respectively. This shows that training reduces the human-AI distance by 42.3% (=1-0.198/0.313), and that the fine-tuned AI is 28% (=1-0.198/0.253) closer to humans than the pre-trained AI distance by 42.3% (=1-0.198/0.313), and that the fine-tuned AI is 28% (=1-0.198/0.253) closer to humans than the pre-trained AI here-trained AI. In summary, training brings the behavior of the generative AI closer to human decision-making tendencies, particularly in the loss domain, suggesting that the learning process effectively influences the AI's bias structure and enhances its humanlike judgment characteristics.

7. Discussions

7.1. Bias Patterns in Humans and Pre-Trained AI

Based on the analysis presented in Section 4, GPT-40 exhibits a risk-averse preference pattern in the gain domain, closely resembling that of human respondents. The AI demonstrates a high level of consistency in replicating human decision-making behavior. In particular, for Question 1, the selection rate for Option 1 is approximately 90% for both humans and AI, suggesting that the behavior aligns with the predictions of prospect theory. Meanwhile, a comparison of the marginal effects by individual attributes reveals no substantial differences between humans and AI in terms of gender. However, for age and income, the absolute values of the marginal effects are consistently larger for the AI, indicating stronger sensitivity to these attributes. Furthermore, the AI exhibits risk-averse behavior in the loss domain, with a more pronounced asymmetry in preference than humans.

Generative AI tends to replicate the association between demographic attributes and preferences in its training data with excessive precision, suggesting a form of "overfitting to bias." Rather than reflecting the variability and inconsistency inherent in human decisions, the model selects contextually plausible responses, which reinforces stereotypical patterns linked to individual characteristics. When the temperature parameter is set to 0.0, the magnitude of the estimated coefficients increases significantly compared to other settings. The suppression of the output randomness eliminates the uncertainty and ambiguity often observed in human judgment, resulting in more mechanical and normative output patterns. These tendencies suggest that the model's responses heavily rely on the context and attribute conditions. This reliance creates a fundamental divergence from human behavior, particularly in how AI reproduces biased structures.

Meanwhile, generative AI exhibits a significantly greater degree of loss aversion than humans. In the

web-based survey targeting human respondents, 57.7% choose Option 1 in Question 2, indicating a moderate tendency toward risk aversion. This rate exceeds the 36% reported by Tversky and Kahneman (1988), suggesting a general inclination to avoid losses. However, generative AI selects Option 1 in only 11.0% and 3.0% of cases under temperature settings of 1.0 and 0.5, respectively, demonstrating a more pronounced aversion to loss. When the temperature is set to 0.0, all but one sample chooses Option 2.

Moreover, the generative AI reveals a similar pattern of loss aversion. When comparing the selection rates for Option 1 in Questions 1 and 2, Tversky and Kahneman (1988) report rates of 72% and 36%, respectively. In contrast, the rates in our web-based survey are 88% and 58%, respectively. In contrast, the generative AI shows a more polarized pattern, with approximately 90% of the responses for Question 1 and 10% for Question 2. This suggests a stronger aversion to losses than observed among human respondents. Furthermore, the comparison of prediction accuracy highlights this discrepancy (Table 5): 98.04% of the human responses aligned with predicted choice 1. At the same time, the AI consistently produces predicted choice two across all samples. These results indicate a clear divergence between humans and AI in terms of risk preference and loss aversion.

In addition to GPT-4o, we present the same set of questions to the following LLMs: GPT-3.5 by OpenAI, Gemini 2.0 Flash by Google, and DeepSeek-R1 by DeepSeek. The results reveal differences in response patterns across the models. GPT-3.5 demonstrates a limited understanding of risk preferences, whereas Gemini 2.0 Flash's answers frequently fail to capture the intended meaning of the questions accurately. In contrast, DeepSeek-R1 accurately reflects existing findings on risk preferences and loss aversion, producing responses similar to those of GPT-4o. These results suggest that the upgrade from GPT-3.5 to GPT-4o improves alignment with behavioral economic principles, particularly regarding risk- and loss-related choices. Gemini appears only to have limited sensitivity to such features, whereas DeepSeek's level of behavioral consistency might be comparable to that of GPT-4o.

7.2. Bias Patterns in Fine-Tuned AI

The study's results indicate that GPT-40 overfits the human tendency toward risk aversion in the gain domain. For Question 1, the proportions of respondents selecting Option 1 are 88%, 91%, and 96% for humans, the pre-trained AI, and the fine-tuned AI, respectively. These results reflect the risk-averse preferences predicted by prospect theory; however, fine-tuning may intensify this bias. This pattern supports the earlier claim that the model overfits cognitive biases. When comparing the marginal effects by demographic attributes, the gender gap between humans and AI widens after fine-tuning. However, the gaps related to age and income narrow. Nonetheless, statistically significant differences remain between humans and the fine-tuned AI for all three attributes. Using the second-order Wasserstein distance as a metric, we find that fine-tuning reduces the distance between the human and AI choice distributions by nearly 10%, suggesting that the AI's decision pattern moves closer to that of humans.

In the loss domain, fine-tuning reduces the extreme risk seeking of AI and brings its behavior closer to that of humans. Before training, the pretrained model selects Option 1 in Question 2 at a rate of only 11.0%, indicating a strong preference for risk proneness. After fine-tuning, this rate increases to 40.4%, exceeding the 36% reported by Tversky and Kahneman (1988) and approaching the 57.7% observed in our web-based survey. The second-order Wasserstein distance between the AI and human responses also decreases by more than 40%, suggesting that learning progresses more effectively in the loss domain than in the gain domain. Furthermore, classification accuracy improves after training. The model no longer assigns all samples to Prediction 2. Instead, more than 65% of the samples receive Prediction 1. This shift indicates a substantial alignment with human behavior. However, a comparison of the marginal effects by attribute reveals a different pattern. It shows that fine-tuning increases the gap between humans and AI. In particular, the fine-tuned model yields large marginal effects for gender (0.264) and age (0.133), suggesting potential overfitting in the model's response to demographic information.

Additionally, fine-tuning reduces the AI's degree of loss aversion, bringing it closer to human behavior. Comparing the choice rates for Option 1 in Questions 1 and 2, the pre-trained model shows a 90% vs. 10% split. After training, the pattern shifts to 96% vs. 40%. By comparison, Tversky and Kahneman (1988) report a 72% vs. 36% split, whereas the human respondents in this study show an 88% vs. 58% split. These results indicate that fine-tuning weakens the AI's excessive loss aversion and aligns its behavior more closely with human responses.

8. Conclusions

This study examines the decision-making tendencies of humans and generative AI in terms of risk preferences and loss aversion. It focuses on the impact of persona settings, defined by attribute information, on AI responses. Furthermore, we examine whether generative AI can replicate human behavioral patterns by incorporating human decision data and evaluate the extent to which training improves the alignment between AI and human tendencies.

The results show that in the gain domain, both humans and the pre-trained AI exhibit risk-averse preferences. In contrast, in the loss domain, AI has a markedly stronger risk-seeking tendency than humans. Furthermore, the influence of individual attributes such as age and income is more pronounced in the AI responses, suggesting that the model tends to amplify preference-related biases. Meanwhile, our results reveal that training moves the AI's choice patterns closer to the mild loss aversion observed in humans. In particular, training in the loss domain significantly reduces the gap between human and AI responses, indicating that the overall learning process performed reasonably well. That is, in the loss domain, AI exaggerates the influence of individual attributes. However, in the gain domain, AI shows signs of overfitting to risk aversion. These patterns suggest possible improvements in training

approaches.

Despite the study's implications, it has limitations that can be addressed in future research. First, future research should develop more effective fine-tuning methods to replicate human preferences more closely. Second, future research should introduce mechanisms that mitigate the tendency of AI to overemphasize the biases associated with individual attributes. While this study focuses on risk preference and loss aversion as an example of cognitive bias, future research should apply the same framework to other biases, including time preference and present bias.

References

- Aher, G. V., Arriaga, R. I., & Kalai, A. T. (2023). Using large language models to simulate multiple humans and replicate human subject studies. In International Conference on Machine Learning (pp. 337-371). PMLR.
- [2] Akata, E., Schulz, L., Coda-Forno, J., Oh, S. J., Bethge, M., & Schulz, E. (2025). Playing repeated games with large language models. Nature Human Behaviour, 1-11.
- [3] Ambrosio, L., Gigli, N., & Savaré, G. (2008). Gradient flows: in metric spaces and in the space of probability measures. Springer Science & Business Media.
- [4] Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., & Wingate, D. (2023). Out of one, many: Using language models to simulate human samples. Political Analysis, 31(3), 337-351.
- [5] Binz, M., & Schulz, E. (2023). Using cognitive psychology to understand GPT-3. Proceedings of the National Academy of Sciences, 120(6), e2218523120.
- [6] Bisbee, J., Clinton, J. D., Dorff, C., Kenkel, B., & Larson, J. M. (2024). Synthetic replacements for human survey data? the perils of large language models. Political Analysis, 32(4), 401-416.
- [7] Brand, J., Israeli, A., & Ngwe, D. (2023). Using GPT for market research. Harvard business school marketing unit working paper, (23-062).
- [8] Brookins, P., & DeBacker, J. (2024). Playing games with GPT: What can we learn about a large language model from canonical strategic games?. Economics Bulletin, 44(1), 25-37.
- [9] Brown, A. L., Imai, T., Vieider, F. M., & Camerer, C. F. (2024). Meta-analysis of empirical estimates of loss aversion. Journal of Economic Literature, 62(2), 485-516.
- [10] Bybee, L. (2023). Surveying Generative AI's Economic Expectations. arXiv preprint arXiv:2305.02823.
- [11] Chen, Y., Kirshner, S. N., Ovchinnikov, A., Andiappan, M., & Jenkin, T. (2025). A manager and an AI walk into a bar: does ChatGPT make biased decisions like we do?. Manufacturing & Service Operations Management.
- [12] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1, pp. 4171-4186.
- [13] Demszky, D., Yang, D., Yeager, D.S., Bryan, C.J., Clapper, M., Chandhok, S., Eichstaedt, J.C., Hecht, C., Jamieson, J., Johnson, M., Jones, M., Krettek-Cobb, D., Lai, L., JonesMitchell, N., Ong, D.C., Dweck, C.S., Gross, J.J., & Pennebaker, J.W. (2023). Using large language models in psychology. Nature Reviews Psychology, 2(11), 688-701.
- [14] Dominguez-Olmedo, R., Hardt, M., & Mendler-Dünner, C. (2025). Questioning the survey responses of large language models. Advances in Neural Information Processing Systems, 37,

45850-45878.

- [15] Falk, A., Becker, A., Dohmen, T., Enke, B., Huffman, D., & Sunde, U. (2018). Global evidence on economic preferences. The quarterly journal of economics, 133(4), 1645-1692.
- [16] Guo, F. (2023). GPT in game theory experiments. arXiv preprint arXiv:2305.05516.
- [17] Hagendorff, T. (2023). Machine psychology: Investigating emergent capabilities and behavior in large language models using psychological methods. arXiv preprint arXiv:2303.13988, 1.
- [18] Hagendorff, T., Fabi, S., & Kosinski, M. (2023). Humanlike intuitive behavior and reasoning biases emerged in large language models but disappeared in ChatGPT. Nature Computational Science, 3(10), 833-838.
- [19] Horton, J. J. (2023). Large language models as simulated economic agents: What can we learn from homo silicus? (No. w31122). National Bureau of Economic Research.
- [20] Jia, J., Yuan, Z., Pan, J., McNamara, P., & Chen, D. (2024). Decision-making behavior evaluation framework for llms under uncertain context. Advances in Neural Information Processing Systems, 37, 113360-113382.
- [21] Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. Econometrica, 47(2), 363-391.
- [22] Kitadai, A., Lugo, S. D. R., Tsurusaki, Y., Fukasawa, Y., & Nishino, N. (2024). Can AI with high reasoning ability replicate humanlike decision making in economic experiments?. arXiv preprint arXiv:2406.11426.
- [23] Lee, S., Peng, T. Q., Goldberg, M. H., Rosenthal, S. A., Kotcher, J. E., Maibach, E. W., & Leiserowitz, A. (2024). Can large language models estimate public opinion about global warming? An empirical assessment of algorithmic fidelity and bias. PLOS Climate, 3(8), e0000429.
- [24] Li, P., Castelo, N., Katona, Z., & Sarvary, M. (2024). Frontiers: Determining the validity of large language models for automated perceptual analysis. Marketing Science, 43(2), 254-266.
- [25] Lorè, N., & Heydari, B. (2024). Strategic behavior of large language models and the role of game structure versus contextual framing. Scientific Reports, 14(1), 18490.
- [26] Macmillan-Scott, O., & Musolesi, M. (2024). (Ir) rationality and cognitive biases in large language models. Royal Society Open Science, 11(6), 240255.
- [27] Mei, Q., Xie, Y., Yuan, W., & Jackson, M. O. (2024). A Turing test of whether AI chatbots are behaviorally similar to humans. Proceedings of the National Academy of Sciences, 121(9), e2313925121.
- [28] Motoki, F., Pinho Neto, V., & Rodrigues, V. (2024). More human than human: measuring ChatGPT political bias. Public Choice, 198(1), 3-23.
- [29] Park, J.S., Zou, C.Q., Shaw, A., Hill, B.M., Cai, C., Morris, M.R., Willer, R., Liang, P., Bernstein, M.S. (2024). Generative agent simulations of 1,000 people. arXiv preprint, arXiv:2411.10109.
- [30] Qiu, L., Singh, P. V., & Srinivasan, K. (2023). Consumer Risk Preferences Elicitation From Large

Language Models. Available at SSRN 4526072.

- [31] Rahwan (2018). Society-in-the-Loop: Programming the Algorithmic Social Contract. Ethics and Information Technology, 20, 5–14.
- [32] Rahwan et al. (2019). Machine Behaviour. Nature, 568, 477–486.
- [33] Ross, J., Kim, Y., & Lo, A. W. (2024). Llm economicus? mapping the behavioral biases of llms via utility theory. arXiv preprint, arXiv:2408.02784.
- [34] Santurkar, S., Durmus, E., Ladhak, F., Lee, C., Liang, P., & Hashimoto, T. (2023). Whose opinions do language models reflect? In International Conference on Machine Learning, pp. 29971-30004.
- [35] Sarstedt, M., Adler, S. J., Rau, L., & Schmitt, B. (2024). Using large language models to generate silicon samples in consumer and marketing research: Challenges, opportunities, and guidelines. Psychology & Marketing, 41(6), 1254-1270.
- [36] Tsuchihashi, T. (2023). How much do you bid? Answers from ChatGPT in first-price and secondprice auctions. Journal of Digital Life, 3.
- [37] Tversky, A., & Kahneman, D. (1988). Rational choice and the framing of decisions. In D. E. Bell,
 H. Raiffa, & A. Tversky (Eds.). In Decision making: Descriptive, normative, and prescriptive interactions, pp. 167–192.
- [38] Tversky, A., & Kahneman, D. (1992). Advances in prospect theory: Cumulative representation of uncertainty. Journal of Risk and uncertainty, 5, 297-323.
- [39] Xie, Y., Liu, Y., Ma, Z., Shi, L., Wang, X., Yuan, W., ... & Mei, Q. (2024). How Different AI Chatbots Behave? Benchmarking Large Language Models in Behavioral Economics Games. arXiv preprint arXiv:2412.12362.
- [40] Xiong, X., Wong, I. A., Huang, G. I., & Peng, Y. (2024). Understanding AI-generated experiments in tourism: replications using GPT simulations. Journal of Travel Research, 00472875241275945.

Table 1: Prior Studies on Decision-Making and Cognitive Biases: Generative AI vs. Humans

Study ID	Authors	Targeted Bias / Experimental Task	AI Model	Population / Sample Size	Average Reproducibility	Conditional Reproducibility	Human Feedback	Evaluation Metric
Economic	Decision-Making and B	ehavioral Biases						
1	Akata et al. (2025)	Economic experiments (Prisoner's Dilemma, Battle of the Sexes)	GPT-3, 3.5, 4	LLM only, using varied payoff matrices across 10 rounds each (1,224 trials total)	N/A (LLM-only results). The model performs well in the Prisoner's Dilemma but shows lower performance in the Battle of the Sexes.	N/A	N/A	N/A
2	Guo (2023)	Economic experiments (Ultimatum Game, Prisoner's Dilemma)	GPT-4 (gpt4-1106- preview)	LLM only (compared with prior studies): Assigned selfish and fairness-oriented personas (UG: 400 trials, PD: 300 trials)	Similar to findings in prior studies	N/A	Prompt adjustment only	N/A
3	Horton (2023)	Economic experiments (Social Preferences: Charness and Rabin, 2002; Fairness: Kahneman et al., 1986; Status Quo Bias: Samuelson and Zeckhauser, 1988; Minimum Wage: Horton, 2023)	GPT-3 (including davinci-003 and others)	Reproduction of results from prior studies	High similarity with existing results	N/A	RLHF applied.	N/A
4	Tsuchihashi (2023)	Economic experiment (bidding behavior in auctions)	GPT-3.5	LLM Only (compared with prior research): 40 rounds each of FPA and SPA	Similar trend to prior research in FPA: Overbidding in FPA, slightly more accurate bidding in SPA	N/A	Prompt adjustment only	N/A
5	Brookins and DeBacker (2024)	Economic experiments (Dictator Game, Prisoner's Dilemma)	GPT-3.5-turbo	LLM Only (compared with prior studies): Dictator Game: 500 rounds, Prisoner's Dilemma: 1,100 rounds	LLMs tend to give fairer and more cooperative responses than humans.	Evaluated based on risk, temptation, and efficiency. Similar tendencies were observed only in the efficiency metric.	Prompt adjustment only	N/A
6	Kitadai et al. (2024)	Economic experiment (Ultimatum Game)	GPT-3.5 (gpt-3.5- turbo-0613), GPT-4 (gpt-4-1106- preview)	LLM only: For each configuration (prompt type, temperature value, and GPT model version), 1000 agents are generated.	Proposer: Adjusting temperature partially replicates human data; GPT improvements move results closer to the theoretical predictions. Responder: Results approach theory but do not fully match experimental data.	N/A	N/A	N/A
7	Jia et al. (2024)	Economic experiments (loss aversion, probability weighting, risk preference)	ChatGPT-4.0- Turbo, Claude-3- Opus, Gemini-1.0- pro	LLM Only (compared with prior studies)	Somewhat similar tendencies to human choices	Bias varies by attributes (e.g., gender); no direct comparison with human responses.	Prompt adjustment only	N/A
8	Lorè and Heydari (2024)	Economic experiments (strategic behavior in two- player games)	GPT-3.5, GPT-4, LLaMa-2	LLM Only: 4 Games × 5 Contexts × 3 Models × 300 Trials	N/A (LLM results only) GPT-4 shows structure-dependent behavior; GPT-3.5 shows context- dependent behavior	N/A	Prompt adjustment only	N/A
9	Mei et al. (2024)	Economic experiments (altruism, fairness, trust, cooperation)	GPT-3.5-Turbo, GPT-4	Compared with public data from approximately 90,000 human decisions, LLM data: 6 games × 30 rounds.	ChatGPT-4 produces results similar to those of humans.	N/A	Prompt adjustment only	N/A
10	Qiu et al. (2024)	Insurance plan selection	GPT-4-turbo	Western insurance data up to 2006 (5,998 cases)	At the aggregate level, the model exhibits choice patterns similar to those of humans; however, it performs poorly at the individual level.	Psychological parameters estimated from LLM responses, such as loss aversion and the probability weighting coefficient, tend to have smaller	N/A	Flscore

						values compared to those reported in previous human studies.		
11	Ross et al. (2024)	Economic experiments (inequity aversion, loss aversion, and time discounting)	GPT-3.5, GPT-4, Claude 2, and Nine Other LLMs	LLM Only (compared with prior studies)	Behavioral differences from Humans: Inequity aversion: LLMs exhibit intense guilt toward others but weak envy. Loss aversion: they respond rationally to gains but show irrational tendencies toward losses. Temporal discounting: they display a more substantial present bias than humans.	N/A	Prompt adjustment only	N/A
12	Xie et al. (2024)	Economic experiments (altruism, fairness, risk, cooperation)	GPT-40, LLaMa3, Claude 3, and others	LLMs only (compared with prior studies): 6 games × 5 models × 50 trials	The distribution of LLM responses showed patterns similar to those of humans. The models tended to make fair choices and exhibited high rates of cooperation.	N/A	Prompt adjustment only	Wasserstein distance
13	Chen et al. (2025)	18 types of cognitive biases	GPT-3.5-turbo, GPT-4	LLM only (binary classification of presence/absence of bias)	Many biases are replicated. GPT-3.5 tends to avoid losses but prefers risk in gains. GPT-4 exhibits consistent risk aversion across various framing contexts.	N/A	N/A	N/A
14	Macmillan-Scott and Musolesi (2024)	12 cognitive tasks	GPT-3.5, GPT-4, Claude2, Bard, LLaMA	LLM Only (compared with prior studies)	GPT-3.5 exhibits the highest proportion of humanlike biases. GPT-4 shows the most humanlike response patterns overall.	N/A	N/A	N/A
	L				· · ·			
Reproducib	ility of Survey Data					1	1	
Reproducib	bility of Survey Data Santurkar et al. (2023)	Public opinion survey data	GPT models, AI21 Labs models	United States (ATP public opinion surveys, 2017–2021)	The opinions generated by language models exhibit significant discrepancies compared to those of the general U.S. population. They tend to exhibit a specific political bias, leaning toward liberal viewpoints.	N/A	RLHF applied.	Wasserstein distance
Reproducib 15 16	bility of Survey Data Santurkar et al. (2023) Park et al. (2024)	Public opinion survey data Personality assessments and behavioral experiments (15 types of economic and psychological experiments)	GPT models, AI21 Labs models Agent architectures utilizing LLMs	United States (ATP public opinion surveys, 2017–2021) Simulation of 1,052 AI agents	The opinions generated by language models exhibit significant discrepancies compared to those of the general U.S. population. They tend to exhibit a specific political bias, leaning toward liberal viewpoints. Reproduced participants' responses two weeks later with high accuracy (approximately 85%).	N/A Evaluated under conditions based on age, race, and political ideology. Showed a consistent tendency to reduce bias across tasks.	RLHF applied. Prompt adjustment only	Wasserstein distance N/A
Reproducib 15 16 17	Santurkar et al. (2023) Park et al. (2024) Dominguez- Olmedo et al. (2025)	Public opinion survey data Personality assessments and behavioral experiments (15 types of economic and psychological experiments) Reproduction of survey data (order and label bias)	GPT models, AI21 Labs models Agent architectures utilizing LLMs GPT-2 to GPT-4, LLaMA, and 43 other models	United States (ATP public opinion surveys, 2017–2021) Simulation of 1,052 AI agents 25 questions from the 2019 American Community Survey (ACS)	The opinions generated by language models exhibit significant discrepancies compared to those of the general U.S. population. They tend to exhibit a specific political bias, leaning toward liberal viewpoints. Reproduced participants' responses two weeks later with high accuracy (approximately 85%). Overall, the accuracy is low.	N/A Evaluated under conditions based on age, race, and political ideology. Showed a consistent tendency to reduce bias across tasks. N/A	RLHF applied. Prompt adjustment only RLHF applied.	Wasserstein distance N/A KL divergence
Reproducib 15 16 17 Cognitive a	bility of Survey Data Santurkar et al. (2023) Park et al. (2024) Dominguez- Olmedo et al. (2025) nd Psychological Biase	Public opinion survey data Personality assessments and behavioral experiments (15 types of economic and psychological experiments) Reproduction of survey data (order and label bias)	GPT models, AI21 Labs models Agent architectures utilizing LLMs GPT-2 to GPT-4, LLaMA, and 43 other models	United States (ATP public opinion surveys, 2017–2021) Simulation of 1,052 AI agents 25 questions from the 2019 American Community Survey (ACS)	The opinions generated by language models exhibit significant discrepancies compared to those of the general U.S. population. They tend to exhibit a specific political bias, leaning toward liberal viewpoints. Reproduced participants' responses two weeks later with high accuracy (approximately 85%). Overall, the accuracy is low.	N/A Evaluated under conditions based on age, race, and political ideology. Showed a consistent tendency to reduce bias across tasks. N/A	RLHF applied. Prompt adjustment only RLHF applied.	Wasserstein distance N/A KL divergence
Reproducib 15 16 17 Cognitive a 18	bility of Survey Data Santurkar et al. (2023) Park et al. (2024) Dominguez- Olmedo et al. (2025) nd Psychological Biase Aher et al. (2023)	Public opinion survey data Personality assessments and behavioral experiments (15 types of economic and psychological experiments) Reproduction of survey data (order and label bias) s Turing experiments (Ultimatum Game, Garden Path Sentences, Milgram Shock Experiment, and Wisdom of Crowds)	GPT models, AI21 Labs models Agent architectures utilizing LLMs GPT-2 to GPT-4, LLaMA, and 43 other models GPT-3.5, GPT-4, and 8 other models	United States (ATP public opinion surveys, 2017–2021) Simulation of 1,052 AI agents 25 questions from the 2019 American Community Survey (ACS) LLMs only (compared with previous studies): 1,000 virtual subjects per task, several thousand responses in total.	The opinions generated by language models exhibit significant discrepancies compared to those of the general U.S. population. They tend to exhibit a specific political bias, leaning toward liberal viewpoints. Reproduced participants' responses two weeks later with high accuracy (approximately 85%). Overall, the accuracy is low. Models from LM-5 onward show a higher level of resemblance to human responses. However, instances of hyper- accuracy bias are observed, particularly in the wisdom of crowds task.	N/A Evaluated under conditions based on age, race, and political ideology. Showed a consistent tendency to reduce bias across tasks. N/A N/A	RLHF applied. Prompt adjustment only RLHF applied. RLHF applied.	Wasserstein distance N/A KL divergence

		deliberative capacity, and causal reasoning ability)						
20	Hagendorff et al. (2023)	Psychological experiments (CRT, semantic illusion)	GPT-1 to GPT-4	Humans: 455 participants; LLMs: 50 trials per task across 10 model types	Up to GPT-3, responses tend to be intuitive; GPT-4 has a higher accuracy than humans, achieving a correct response rate of 96%.	N/A	Prompt adjustment only	N/A
Political At	titudes and Voting Beh	avior						
21	Argyle et al. (2023)	Political attitudes and voting behavior	GPT-3	ANES: 1,304 individuals + 2,873 individuals	A strong correlation (greater than 0.9) is observed in voting behavior.	Reproducibility is low in specific categories, such as independents.	Prompt adjustment only	Cramér's V, Tetrachoric Correlation
22	Bisbee et al. (2024)	Comparison of ANES data using political persona settings	ChatGPT 3.5 Turbo, ChatGPT 4.0, Falcon-40B	United States (ANES 2016 and 2020)	LLM responses closely resemble human responses but exhibit less variability.	When examined by political attributes, human and LLM responses show similar overall tendencies despite differing variability. However, when conditioned on individual attributes, human and LLM responses tend to diverge in their patterns.	RLHF applied.	N/A
23	Motoki et al. (2024)	Political bias	GPT-3.5	LLM only: 100 responses under Democratic/Republican persona prompts and 100 responses without persona prompts	No comparison with humans. Default responses tend to align with Democratic positions. The political stance changes depending on the specified persona.	N/A	Prompt adjustment only	N/A
Other Studi	es							
24	Brand et al. (2023)	Willingness to pay (WTP) for multiple products	GPT-3.5-turbo- 0125	LLM only (compared with prior studies)	The average WTP patterns are well replicated. However, for new products, consistency with human responses tends to decline.	For specific product categories, the model exhibits patterns similar to those of humans.	Fine-tuning applied.	N/A
25	Bybee (2023)	Economic expectations forecast (based on WSJ articles)	GPT-3.5	United States (300 WSJ articles from 1984 to 2021)	High correlation with existing surveys such as SPF, AAII, and CFO.	N/A	N/A	N/A
26	Lee et al. (2024)	Beliefs about global warming	GPT-3.5-turbo-16k, GPT-4	U.S. national surveys (2017: 1,304 respondents; 2021: 1,006 respondents)	Beliefs about global warming are highly replicable (85%), while replication rates for beliefs about its causes and associated concerns are lower (51% and 48%, respectively).	The opinions of Black respondents tend to be underrepresented.	Prompt adjustment only	F1 score, Cramér's V
27	Li et al. (2024)	Brand recognition and perceived similarity.	GPT-4, GPT-Neo	Evaluations of 21 automobile brands by 530 participants.	High accuracy rate (87.2%).	Consistent patterns across age groups and demographic attributes.	Prompt adjustment only	Triplet matching rate
28	Xiong et al. (2024)	Emotion and belief formation (in the tourism domain).	GPT-3.5-turbo	LLM only (compared with prior research): 16 scenarios \times 100 iterations.	The responses generated by the LLM exhibit trends similar to those observed in human responses.	N/A	Prompt adjustment only	N/A

Individual Attributes	Mean
Female dummy	0.507 [0.500]
Age (10 years)	4.49 [1.33]
Annual income (JPY 1 million)	5.93 [3.10]
Obs.	4,838

 Table 2: Descriptive Statistics

Notes: The values in parentheses in the table indicate standard deviations.

```
1
      import time
      import pandas as pd
2
3
      import openai
4
      import os
5
      import re
6
      openai.api_key = os.getenv("OPENAI_API_KEY")
7
8
      file_path = r"C:***¥***.csv"
9
      df = pd.read_csv(file_path)
10
11
      def extract_reason_and_answer(response_text):
    match_reason = re.search(r"reason[::]¥s*(.*)", response_text)
    match_answer = re.search(r"(answer[::]¥s*[12])", response_text)
12
13
14
15
16
           if not match reason:
17
               lines = response_text.split("¥n")
               reason = lines[0].strip() if lines else "unknown"
18
19
           else:
20
               reason = match_reason.group(1).strip()
21
22
           answer = match_answer.group(1).strip()
23
24
           return reason, answer
25
26
      def generate_reason_and_answer(row):
      question = f"""The following question is asked to someone with the following gender:
{row['gender']}, age: {row['age']}years and annual income: {row['income']}million yen.
Imagine you receive an additional 50,000 yen on top of your current wealth, and are
27
28
      asked to choose between the following options. Which do you choose?
29
           Option 1: Lose 10,000 yen with certainty
          Option 2: A 50% chance of losing 20,000 yen, and a 50% chance of losing nothing Please write your answer in the following exact format:
30
31
           - Reason: [Reason for your choice]
32
33
           - Answer: [Output only 1 or 2]
34
35
           for attempt in range(3):
36
37
                   time.sleep(1)
38
39
                   response = openai.chat.completions.create(
                       model="gpt-40",
40
41
                       messages=[
                           {"role": "system", "content": f"Your attributes are gender:
42
      43
44
                        1.
                       max_tokens=180,
45
46
                       temperature=1.0
                   )
47
48
                   full_response = response.choices[0].message.content.strip()
49
                   print(f" API Response for row {row.name}: {full_response}")
50
51
52
                   reason, answer = extract_reason_and_answer(full_response)
53
                   return pd.Series([reason, answer])
54
55
      df[["Reason", "Answer"]] = df.apply(generate_reason_and_answer, axis=1)
56
      output_file_path = r"C:***¥***.csv"
57
58
      df.to_csv(output_file_path, index=False)
59
      print(f" Generated the reasons and answers and saved them in a new file:
60
      {output_file_path}")
```

Figure 1: Python Script for AI Question Presentation and Response Collection

Notes: The generative AI is assigned the role of "Your attributes are gender: {row['gender']}, age: {row['age']} years, annual income: JPY {row['income']}million." and a question is asked. This diagram illustrates an excerpt of the key elements of the script used. The actual questions are asked in Japanese. In addition, after fine-tuning in Section 5, "gpt-4o" on line 40 is replaced with the name of the fine-tuned model.

	1	1		
Gender	Age	Income	Reason	Answer
Male	26	400	Given my current income and age, I feel more secure in	Answer: 1
			receiving a guaranteed gain than taking a risk.	
			Receiving the guaranteed JPY 10,000 gives me peace	
			of mind.	
Female	39	850	By receiving the guaranteed JPY 10,000, I can increase	Answer: 1
			my funds safely and steadily without any risk.	
Male	38	400	Choosing the guaranteed JPY 10,000 is a less risky	Answer: 1
			option, and with a limited income, I prefer stable and	
			secure gains.	
Male	22	850	Considering my current income and age, I am willing	Answer: 2
			to take some risks and make decisions based on	
			expected value. Since the expected value of Option 2 is	
			JPY 10,000, I choose to take the risk.	
Male	61	400	Receiving JPY 10,000 for sure allows me to earn some	Answer: 1
			extra income while avoiding risk, which is reassuring.	
Female	41	850	I value stability, so I prefer an option with a guaranteed	Answer: 1
			return. I chose the option to receive the guaranteed JPY	
			10,000 because it involves no risk.	
Male	47	300	Given my current income, it feels safer and more	Answer: 1
			beneficial to increase my money in a guaranteed way.	
			I want to prioritize a stable income, so I chose the	
			option that ensures a monthly income of JPY 10,000.	
Male	38	850	Considering my current income and financial situation,	Answer: 1
			I prefer to increase my earnings without taking major	
			risks. Receiving the guaranteed JPY 10,000 helps me	
			maintain financial stability in my life.	
Male	55	300	Since my income is limited, I want the immediate	Answer: 1
			sense of security that comes from receiving the	
			guaranteed JPY 10,000. Avoiding risk and making a	
			safe choice feels more appropriate for me.	

	Table 3:	Examp	le of Res	ponses Us	ing Ope	nAI API
--	----------	-------	-----------	-----------	---------	---------

Notes: Table 3 shows excerpts of responses obtained through the procedure described in Section 3.3 (for Question 1, with the setting temperature = 1.0). In practice, the results are presented in CSV format.

Question 1 Gain	AI temperature 0.0	AI temperature 0.5	AI temperature 1.0	Human
Mean	0.898	0.889	0.911	0.884
Std. Dev.	0.302	0.314	0.285	0.320
Obs.	4,838	4,838	4,838	4,838

Table 4: Choice Results: Humans vs. Pre-trained AI

Question 2 Loss	AI temperature 0.0	AI temperature 0.5	AI temperature 1.0	Human
Mean	0.000	0.030	0.110	0.577
Std. Dev.	0.014	0.172	0.313	0.494
Obs.	4,838	4,838	4,838	4,838

Notes: All instances of "AI" in the table refer to the "pre-trained AI."

Qu	uestion 1 Gain	AI temperature 0.0	AI temperature 0.5	AI temperature 1.0	Human
	Female dummy	0.086^{***} (0.006)	0.026*** (0.007)	0.033*** (0.007)	0.041*** (0.009)
Marginal Effects	Age (10 years)	0.120*** (0.004)	0.073*** (0.003)	0.057*** (0.003)	0.015*** (0.003)
	Annual income (1 million yen)	-0.031*** (0.001)	-0.033*** (0.001)	-0.027*** (0.001)	-0.005*** (0.001)
	obs.	4,838	4,838	4,838	4,838
McFadden R²		0.6355	0.3932	0.3472	0.0140
Predicte	ed Probability	0.898 (0.004)	0.889 (0.005)	0.911 (0.004)	0.884 (0.005)

 Table 5: Logit Model Estimation: Humans vs. Pre-trained AI (Gain, 3 Attributes)

Notes: Rows 2-4 show the marginal effects of each attribute based on the logit model. The numbers in parentheses indicate standard errors. " *** " denotes statistical significance at the 1% level. Row 7 presents the average predicted probabilities based on the estimates from the logit model. Here, all instances of "AI" in the table refer to "pre-trained AI."

Qu	testion 2 Loss	AI temperature 0.0	AI temperature 0.5	AI temperature 1.0	Human
	Female dummy	-	0.028*** (0.006)	0.028*** (0.009)	-0.006 (0.014)
Marginal Effects	Age (10 years)	-	0.008*** (0.002)	0.023*** (0.004)	0.021*** (0.005)
	Annual income (1 million yen)	-	-0.006*** (0.001)	-0.006*** (0.002)	-0.008*** (0.002)
	obs.	4,838	4,838	4,838	4,838
McF	adden R²	-	0.0644	0.0204	0.0038
Predicte	ed Probability	-	0.030*** (0.002)	0.110*** (0.004)	0.577*** (0.007)

Table 6: Logit Model Estimation Humans vs. Pre-trained AI (Loss, 3 Attributes)

Notes: Rows 2-4 show the marginal effects of each attribute based on the logit model. The numbers in parentheses indicate standard errors. " *** " denotes statistical significance at the 1% level. Row 7 presents the average predicted probabilities based on the estimates from the logit model. Here, all instances of "AI" in the table refer to "pre-trained AI."



Figure 2: Choice Probabilities: Humans vs. Pre-trained AI (Gain, 3 Attributes)

Notes: AI_10 refers to the responses generated by the pre-trained AI when queried with a temperature setting of 1.0. The same applies to AI_05 and AI_00. The vertical bars at each point in the figure represent 95% confidence intervals.



Figure 3: Choice Probabilities: Humans vs. Pre-trained AI (Loss, 3 Attributes)

Notes: AI_10 refers to the responses from the pre-trained AI when queried with a temperature setting of 1.0. The same applies to AI_05 and AI_00. The vertical bars at each point in the figure represent 95% confidence intervals.

Question 1	AI	AI	AI	Human
Gain	temperature 0.0	temperature 0.5	temperature 1.0	
Predicted value, Observed value				
(2, 2)	6.39%	4.24%	2.34%	0%
	(n=309)	(n=205)	(n=113)	(n=0)
(2, 1)	1.53%	2.00%	1.34%	0%
	(n=74)	(n=97)	(n=65)	(n=0)
(1, 2)	3.78%	6.84%	6.57%	11.60%
	(n=183)	(n=331)	(n=318)	(n=561)
(1,1)	88.30%	86.92%	89.75%	88.40%
	(n=4,272)	(n=4,205)	(n=4,342)	(n=4,277)
Prediction Accuracy	94.69%	91.15%	92.08%	88.40%

 Table 7: Logit Model Prediction Accuracy: Humans vs. Pre-trained AI

Question 2 Loss	AI temperature 0.0	AI temperature 0.5	AI temperature 1.0	Human
Predicted value,				
Observed value				
(2, 2)		96.96%	88.98%	0.93%
(2, 2)	-	(n=4,691)	(n=4,305)	(n=45)
(2, 1)		3.04%	11.02%	1.03%
(2, 1)	-	(n=147)	(n=533)	(n=50)
(1, 2)		0%	0%	41.38%
(1, 2)	-	(n=0)	(n=0)	(n=2,002)
(1.1)		0%	0%	56.66%
(1,1)	-	(n=0)	(n=0)	(n=2,741)
Prediction Accuracy	-	96.96%	88.98%	57.59%

Notes: All instances of "AI" in the table refer to the "pre-trained AI."

{"messages": [{"role": "system", "content": "You are a 36-year-old man with an annual income of 3 million yen."}, {"role": "user", "content": "Question: Imagine you receive an additional 50,000 yen on top of your current wealth, and are asked to choose between the following options. Which do you choose? Option 1: Lose 10,000 yen with certainty Option 2: A 50% chance of losing 20,000 yen, and a 50% chance of losing nothing. What is the probability of choosing option 1?"}, {"role": "assistant", "content": "0.58"}]

{"messages": [{"role": "system", "content": " You are a 57-year-old woman with an annual income of 3 million yen."}, {"role": "user", "content": " Question: Imagine you receive an additional 50,000 yen on top of your current wealth, and are asked to choose between the following options. Which do you choose? Option 1: Lose 10,000 yen with certainty Option 2: A 50% chance of losing 20,000 yen, and a 50% chance of losing nothing. What is the probability of choosing option 1?"}, {"role": "assistant", "content": "0.62"}]

{"messages": [{"role": "system", "content": " You are a 40-year-old woman with an annual income of 8.5 million yen."}, {"role": "user", "content": " Question: Imagine you receive an additional 50,000 yen on top of your current wealth, and are asked to choose between the following options. Which do you choose? Option 1: Lose 10,000 yen with certainty Option 2: A 50% chance of losing 20,000 yen, and a 50% chance of losing nothing. What is the probability of choosing option 1?"}, {"role":

"assistant", "content": "0.54"}]}

{"messages": [{"role": "system", "content": " You are a 36-year-old man with an annual income of 3 million yen."}, {"role": "user", "content": " Question: Imagine you receive an additional 50,000 yen on top of your current wealth, and are asked to choose between the following options. Which do you choose? Option 1: Lose 10,000 yen with certainty Option 2: A 50% chance of losing 20,000 yen, and a 50% chance of losing nothing. What is the probability of choosing option 1?"}, {"role": "assistant", "content": "0.87"}]

{"messages": [{"role": "system", "content": " You are a 57-year-old woman with an annual income of 3 million yen."}, {"role": "user", "content": " Question: Imagine you receive an additional 50,000 yen on top of your current wealth, and are asked to choose between the following options. Which do you choose? Option 1: Lose 10,000 yen with certainty Option 2: A 50% chance of losing 20,000 yen, and a 50% chance of losing nothing. What is the probability of choosing option 1?"}, {"role": "assistant",

"content": "0.93"}]}

Figure 4: Dataset for fine-tuning

Notes: An excerpt from the JSONL file used for fine-tuning. The model is trained based on the choice probabilities for each question. To proceed with the work while verifying that fine-tuning is being appropriately performed, we create JSONL files containing selection probabilities for every 500 samples and sequentially load them into the model.

```
import openai
 1
    import os
 2
    import json
 3
 4
 5
 6
    openai.api_key = os.getenv("OPENAI_API_KEY")
 7
 8
9
    file_path = r"C:\Users\***\***.jsonl"
10
11
12
    print("Uploading training file...")
    response = openai.files.create(
13
        file=open(file path, "rb"),
14
        purpose="fine-tune"
15
16
    )
    file id = response.id
17
    print(f"File uploaded successfully! File ID: {file_id}")
18
19
20
21
    print("Starting fine-tuning process...")
    fine_tune_job = openai.fine_tuning.jobs.create(
22
23
        training_file=file_id,
        model="gpt-40-2024-08-06",
24
25
        hyperparameters={
26
            "n_epochs": 1,
            "learning rate multiplier": 0.2
27
        }
28
29
        )
30
    fine tune id = fine tune job.id
    print(f"Fine-tuning started! Job ID: {fine_tune_id}")
31
32
33
    print("Checking fine-tuning status...")
34
    status_response = openai.fine_tuning.jobs.retrieve(fine_tune_id)
35
    print("Fine-tuning status:", status_response.status)
36
```

Figure 5: Python Script for Fine-tuning

Notes: Fine-tuning is conducted by loading the dataset shown in Figure 5 into the GPT-40 model with the settings epoch = 1 and learning rate (LR) = 0.2.

Table 8: Choice	e Results: Human	s. Pre-trained AI	and Fine-tuned AI
	. Itesuits. Ituman	s, 1 1 c=ti aintu 1 xi	, and r me-tuneu m

Question 1 Gain	Pre-trained AI temperature 1.0	Fine-tuned AI temperature 1.0	Human
Mean	0.911	0.961	0.884
Std. Dev.	0.285	0.193	0.320
Obs.	4,838	4,838	4,838

Question 2 Loss	Pre-trained AI temperature 1.0	Fine-tuned AI temperature 1.0	Human
Mean	0.110	0.404	0.577
Std. Dev.	0.313	0.491	0.494
Obs.	4,838	4,838	4,838

Table 9: Logit Model Estimation: Humans, Pre-trained AI, and Fine-tuned AI (Gain, 3

Question 1 Gain		Pre-trained AI temperature 1.0	Fine-tuned AI temperature 1.0	Human
	Female dummy	0.033*** (0.007)	0.019*** (0.005)	0.041*** (0.009)
Marginal Effects	Age (10 years)	0.057*** (0.003)	0.030*** (0.002)	0.015*** (0.003)
	Annual income (1 million yen)	-0.027*** (0.001)	-0.014*** (0.001)	-0.005*** (0.001)
obs.		4,838	4,838	4,838
McFadden R²		0.3472	0.4500	0.0140
Predicted Probability		0.911 (0.004)	0.961 (0.003)	0.884 (0.005)

Attributes)

Notes: Rows 2-4 show the marginal effects of each attribute based on the logit model. The numbers in parentheses indicate standard errors. " *** " denotes statistical significance at the 1% level. Row 7 presents the average predicted probabilities based on the estimates from the logit model.

3 Attributes)					
	Question 2 Loss	Pre-trained AI temperature 1.0	Fine-tuned AI temperature 1.0	Human	
	Female dummy	0.028*** (0.009)	0.264*** (0.011)	-0.006 (0.014)	
Marginal	Age	0.023***	0.133***	0.021***	

(0.004)

-0.006***

(0.002)

4,838

0.0204

0.110

(0.004)

(0.004)

-0.003

(0.002)

4,838

0.1804

0.404

(0.007)

(0.005)-0.008***

(0.002)

4,838

0.0038

0.577

(0.007)

Effects

(10 years)

Annual income (1 million yen)

obs.

McFadden \mathbf{R}^2

Predicted Probability

Table 10: Logit Model Estimation: Humans, Pre-trained AI, and Fine-tuned AI (Loss,

Notes: Rows 2 to 4 show the marginal effects of each attribute based on the logit model. The numbers in parentheses indicate standard errors. " *** " denotes statistical significance at the 1% level. Row 7 presents the average predicted probabilities based on the estimates from the logit model.



Figure 6: Choice Probabilities: Humans, Pre-trained AI, and Fine-tuned AI (Gain, 3 Attributes)

Notes: AI_10 and AI_FT_10 refer to the responses from the pre-trained AI and the fine-tuned AI, respectively, when asked with a temperature setting of 1.0. The vertical bars at each point in the figure represent 95% confidence intervals.



Figure 7: Choice Probabilities: Humans, Pre-trained AI, and Fine-tuned AI (Loss, 3 Attributes)

Notes: AI_10 and AI_FT_10 refer to the responses from the pre-trained AI and the fine-tuned AI, respectively, when asked with a temperature setting of 1.0. The vertical bars at each point in the figure represent 95% confidence intervals.

Table 11: Logit Model Prediction Accuracy: Humans, Pre-trained AI, and Fine-tuned

AI

Question 1	Pre-trained AI	Fine-tuned AI	Human
Gain	temperature 1.0	temperature 1.0	
Predicted value, Observed value			
(2, 2)	2.34%	1.24%	0%
	(n=113)	(n=60)	(n=0)
(2, 1)	1.34%	0.78%	0%
	(n=65)	(n=38)	(n=0)
(1, 2)	6.57%	2.65%	11.60%
	(n=318)	(n=128)	(n=561)
(1,1)	89.75%	95.33%	88.40%
	(n=4,342)	(n=4,612)	(n=4,277)
Prediction Accuracy	92.08%	96.57%	88.40%

Question 2 Loss	Pre-trained AI temperature 1.0	Fine-tuned AI temperature 1.0	Human
Predicted value,			
Observed value			
(2, 2)	88.98%	47.66%	0.93%
(2, 2)	(n=4,305)	(n=2,306)	(n=45)
(2 1)	11.02%	18.13%	1.03%
(2, 1)	(n=533)	(n=877)	(n=50)
(1, 2)	0%	11.95%	41.38%
	(n=0)	(n=578)	(n=2,002)
(1,1)	0%	22.26%	56.66%
	(n=0)	(n=1,077)	(n=2,741)
Prediction Accuracy	88.98%	69.93%	57.59%

Prediction Accuracy88.98%69.93%57.59%Notes: The percentage values are rounded to the third decimal place; therefore, the sum of the values in cells (2,2) and (1,1) does not necessarily match the prediction accuracy.

	Between Human and Pre-trained AI	Between Human and Fine-tuned AI	Between Pre-trained AI and Fine-tuned AI
Gain (N=4,838)	0.135	0.123	0.081
Loss (N=4,838)	0.467	0.263	0.353
Total (N=9,676)	0.343	0.198	0.253

Table 12: Wasserstein Distance of Order 2 between Three Groups

Notes: All values in the table for AI are based on the responses under the default setting of temperature = 1.0.