# Large Language Models are Near-Optimal Decision-Makers with a Non-Human Learning Behavior

Hao Li[a, †], Gengrui Zhang[b, †], Petter Holme[c,d], Shuyue Hu[e, ‡], Zhen Wang[a, ‡]

[a]School of Cybersecurity, Northwestern Polytechnical University, China.
[b]Department of Psychology, University of Southern California, United States.
[c]Department of Computer Science, Aalto University, Finland.
[d]Center for Computational Social Science, Kobe University, Japan.
[e]Shanghai Artificial Intelligence Laboratory, China.

li.hao@mail.nwpu.edu.cn; gengruiz@usc.edu; petter.holme@aalto.fi
hushuyue@pjlab.org.cn; w-zhen@nwpu.edu.cn;

## Abstract

Human decision-making belongs to the foundation of our society and civilization, but we are on the verge of a future where much of it will be delegated to artificial intelligence. The arrival of Large Language Models (LLMs) has transformed the nature and scope of AI-supported decision-making; however, the process by which they learn to make decisions, compared to humans, remains poorly understood. In this study, we examined the decision-making behavior of five leading LLMs across three core dimensions of real-world decision-making: uncertainty, risk, and set-shifting. Using three well-established experimental psychology tasks designed to probe these dimensions, we benchmarked LLMs against 360 newly recruited human participants. Across all tasks, LLMs often outperformed humans, approaching near-optimal performance. Moreover, the processes underlying their decisions diverged fundamentally from those of humans. On the one hand, our finding demonstrates the ability of LLMs to manage uncertainty, calibrate risk, and adapt to changes. On the other hand, this disparity highlights the risks of relying on them as substitutes for human judgment, calling for further inquiry.

**Keywords:** Decision-making, Artificial intelligence, Large language models, Experimental psychology

Decision-making is a unifying theme of almost all the social and behavioral sciences. Moreover, it has been closely tied to artificial intelligence (AI), ever since the inception of the latter. It is no coincidence that one of AI's founders, Herbert Simon, began his career studying organizational decision-making [1]. Simon once pointed out that "the capacity of the human mind for formulating and solving complex problems is very small compared with the size of the problems whose solution is required for objectively rational behavior in the real world" [1]. And so the first generation of AI built decision systems [2] by implementing logic and circumstantial knowledge. In contrast, the last decade of AI development has had a very different driving force—to imitate human intellectual output in broad generality. By training self-supervised artificial neural networks on immense textual corpora—so-called Large Language Models (LLMs)—it is commonly accepted that AI has now passed the Turing test [3, 4]. LLMs do make decisions if we prompt them to do so, but are those decisions made with a capacity as "very small" as the humans they try to emulate? That is the central question we pursue in this paper.

LLMs are remarkably versatile, capable of generating meaningful output across a wide range of tasks, and are, unsurprisingly, already being deployed in real-world settings to make or influence decisions [5–7]. For a concrete example, in a legal case in Colombia concerning the medical expenses of an autistic boy, the judge not only queried ChatGPT, asking "Is an autistic minor exonerated from paying fees for their therapies?", but also cited both the prompt and the answer in the official ruling [8]. Likewise, LinkedIn, a major professional networking platform, has deployed an LLM-based chatbot to assist recruiters in automating

---

[†]These authors contributed equally; this work was partly done during their internship at Shanghai Artificial Intelligence Laboratory.
[‡]Corresponding authors.

critical hiring decisions, such as the shortlisting of job candidates, thereby directly shaping individuals' career opportunities [9]. Globally, surveys estimate that 71% of major organizations have already regularly used generative AI in at least one business function [10]. For better or worse, generative AI is rapidly becoming integral to decision-making infrastructures, where each decision made, like those made by humans, has the power to shape individual lives and, potentially, ripple outward to influence entire societies [11–14].

However, are generative AI systems truly ready to assist or even replace humans in decision-making? Recent studies have reported strong performance for the decision-making of these systems in domains such as healthcare [15, 16], finance [17, 18], and strategic games [19–22]. Yet, arguably more critical—but far less explored—is how they arrive at these decisions, particularly when operating outside the narrowly defined scenarios. Domain-specific tasks, though valuable, often introduce numerous confounds, such as background knowledge, ethical considerations, and cultural biases, which can obscure the working mechanisms that drive system behaviors. Moreover, although these tasks capture real-world complexity, they often conflate conceptually distinct dimensions of decision-making, such as risk and uncertainty. For instance, a medical diagnosis benchmark may fold epistemic uncertainty about a patient's true condition into the known risks of various treatments [16]. Consequently, even strong performance on these tasks offers little diagnostic insight into whether the system effectively excels at managing uncertainty, calibrating risk, or both, let alone how these dimensions are reasoned about. In contrast, what fundamentally characterizes human decision-making capacity and their agency is not merely the correctness of a task-specific choice, but the ability to navigate different decision dimensions while simultaneously aligning their choices with goals [23–25].

In this study, we compare the decision-making of generative AI systems and humans, particularly under conditions of uncertainty, risk, or set-shifting. Here, decision-making refers to the process by which individuals assess multiple alternatives and choose actions that align with defined goals. Uncertainty involves acting with incomplete information and ambiguous future outcomes, requiring a balance between short-term and long-term consequences [26, 27]. Risk entails evaluating potential gains and losses based on known probabilities, demanding careful judgment of outcome likelihoods [28, 29]. Set-shifting refers to dynamic environments where conditions evolve over time, requiring the ability to adapt strategies as new information emerges [30, 31]. Uncertainty, risk, and set-shifting are three key dimensions that dominate most real-world decisions [32–34]. They each represent a conceptually distinct construct, yield their own behavioral patterns, and engage their own cognitive mechanism [35–37].

Our methodology for investigating if and how generative AI systems navigate these decision dimensions aligns with recent calls for machine psychology [38–43], which advocates the use of experimental psychology paradigms. Unlike domain-specific tasks, psychological tests, originally designed for humans, strip away contextual confounds and isolate distinct decision dimensions. As a result, they can provide simplified yet tight experimental control, enable precise theory testing, and reveal general cognitive mechanisms that may extend across diverse domains. We adapted three well-established psychological tests: the Iowa Gambling Task for uncertainty [44], the Cambridge Gambling Task for risk [45], and the Wisconsin Card Sorting Task for set-shifting [46] (see Methods for task details). We treated LLMs as subjects in these tests. To mitigate the potential memorization effects of these systems [39], we reworded task descriptions and redesigned payoff structures while preserving the essence of the original tests.

We considered five leading LLMs: GPT (*gpt-4o-2024-08-06*), GPTo4m (*o4-mini-2025-04-16*), Claude (*claude-3-5-sonnet-20240620*), Gemini (*gemini-1.5-pro-002*), and DeepSeek (*DeepSeek-R1-2025-01-20*). To provide a yardstick for their decision-making, we benchmarked them with 360 newly recruited human participants (120 per task), presenting both groups identical experimental instructions (see Methods, and *SI Appendix, Supplementary Note 1* for details). Across all three tests, we consistently observed that LLMs were able to perform significantly better than human participants; however, the ways in which they arrived at their decisions were fundamentally different from those of humans. This key finding offers direct evidence of these systems' general decision-making competence—particularly in managing uncertainty, calibrating risk, and adapting to change. On the other hand, it cautions against using these systems as stand-ins for human decision-making in real-world contexts or behavioral research [47–49]. More broadly, it highlights the critical need for policymakers and system designers to carefully consider how much autonomy to delegate to such systems, ensure transparent communication with end-users about these systems' potential cognitive differences, and maintain meaningful human oversight, especially in domains where human-like reasoning and judgment are essential.

## Results

In each test, most evaluated LLMs outperformed human participants and approached near-optimal performance, yet relied on decision-making strategies that were notably different from those of humans. Computational models revealed that the posterior parameter estimates for LLMs diverged significantly from those of human participants, highlighting fundamental differences in underlying cognitive processes.
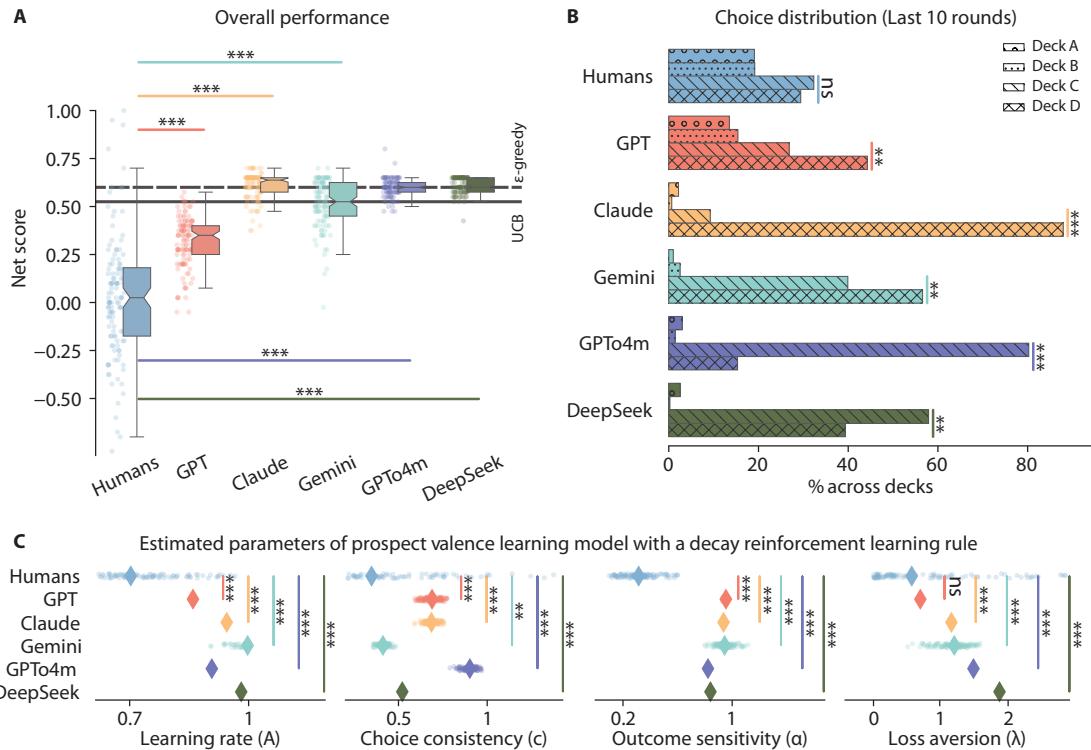
**Fig. 1**: **All the LLMs significantly outperformed humans in the Iowa Gambling Task, yet differed in choice preferences and exhibited distinct parameter estimates in the prospect valence learning model compared to humans. Panel (A)** shows the net scores of LLMs, human participants, and two well-established strategies (Upper Confidence Bound and $\epsilon$−greedy). A Mann-Whitney U test indicated that GPT-4o, GPTo4m, Claude, Gemini, and DeepSeek significantly outperformed human participants. Specifically, Claude achieved the highest median net score and GPTo4m showed the lowest variance. **Panel (B)** shows the distributions of deck selections over the last 10 rounds. Both human participants and LLMs predominantly favored the advantageous decks (C and D). However, a two-proportion Z test revealed that GPT, Claude, and Gemini showed a stronger preference for deck D over deck C, whereas GPTo4m and DeepSeek showed the opposite preference. In contrast, human participants selected decks C and D in similar proportions. **Panel (C)** presents the posterior estimates of the parameters in the prospect valence learning model. Mann-Whitney U tests indicated that compared to humans, all the LLMs demonstrated higher learning rates ($A$), greater sensitivity ($\alpha$) to outcomes, and more tendency ($c$) to make deterministic decisions. Moreover, except for GPT, they also showed stronger reactions ($\lambda$) to penalties than humans. GPT showed no significant difference from humans on the parameter of loss aversion. All inferential (Mann–Whitney U tests) and descriptive statistics (means, medians, and standard deviations) for the Iowa Gambling Task are reported in *SI Appendix*, Table S1, including pairwise comparisons (panels a–c) and parameter estimates (panel d). All parameter estimates demonstrated satisfactory convergence with R-hat ($\hat{R}$) values below 1.01.

## Decision-making under uncertainty

The Iowa Gambling Task tested whether players can prioritize long-term benefits over short-term gains in an uncertain environment. During the task, participants repeatedly selected cards from four decks (labeled A, B, C, and D). Each choice offered an immediate reward, with the possibility of an occasional penalty. Unbeknownst to players, decks A and B were disadvantageous, offering higher immediate rewards but lower long-term expected payoffs, whereas decks C and D were advantageous, providing smaller immediate rewards but greater long-term gains. The outcomes of individual card selections were unpredictable, requiring participants to infer the underlying payoff structure through repeated trials and feedback.

We measured the overall performance of human participants and LLMs using net scores [50], defined as the difference in the proportion of selections between the advantageous and disadvantageous decks. Fig. 1A compares the net scores of human participants, two GPT models, Claude, Gemini, and DeepSeek. The five types of LLMs all significantly outperformed human participants, demonstrating higher task proficiency. Among the models, Claude achieved the highest median net score, indicating better overall task performance, whereas GPTo4m showed the lowest variance, reflecting greater consistency across the task. Additionally, when benchmarked against two well-established strategies for decision-making under uncertainty, namely, upper confidence bound [51] and $\epsilon$−greedy [52], Claude, Gemini, GPTo4m, and DeepSeek approached their

near-optimal performance. Linear regression models further revealed that, over time, these LLMs also learned faster than humans, demonstrating steeper slopes in their proportion of advantageous deck selections (*SI Appendix,* Fig. S1). Thus, compared to humans, LLMs managed to identify advantageous decks earlier and adapted their choices more effectively to the task's reward-penalty structure.

Beyond performance, we analyzed the distribution of choices across decks to assess differences in decision strategies. As illustrated in Fig. 1B, during the final ten rounds, both human participants and LLMs predominantly selected the two advantageous decks (i.e., C and D), which have equivalent long-term expected returns. However, their choice distributions differed substantially. Human participants selected decks C and D in similar proportions, consistent with prior experimental findings in the literature [44]. In contrast, LLMs exhibited systematic but heterogeneous preferences across the advantageous decks, with some models showing a stronger inclination toward deck D, which was associated with less frequent penalties, and others favoring deck C, which had more frequent but smaller penalties. This suggests that LLM behavior is more sensitive to penalty frequency, whereas human choices appear less influenced by the frequency of losses, resulting in a more balanced distribution across the two advantageous options.

We compared the behaviors of humans and LLMs using the Prospect Valence Learning Model with a decay Reinforcement Learning rule [53] (see *SI Appendix, Supplementary Note 2* for details). This model was often used to characterize decision-making processes in this task through four parameters: learning rate ($A$), choice consistency ($c$), outcome sensitivity ($\alpha$), and loss aversion ($\lambda$). Our analysis revealed significant differences in the posterior estimates of these parameters between LLMs and humans (Fig. 1C). Specifically, all the LLMs showed a significantly higher learning rate ($A$), indicating a stronger reliance on cumulative past outcomes, which allowed them to identify and exploit patterns more effectively over time. In contrast, humans weighted recent and past outcomes more evenly, reflecting a more flexible but less pattern-oriented strategy. All the LLMs also exhibited higher choice consistency ($c$), meaning that they more reliably selected options with the highest learned expected value, with less variability or noise in their choices, whereas humans displayed lower consistency, indicating more stochastic or exploratory behavior. Additionally, LLMs generally showed higher outcome sensitivity ($\alpha$) and loss aversion ($\lambda$) compared to humans, with the exception of GPT, which did not significantly differ from humans in loss aversion. This suggests that LLMs typically tended to have a stronger reaction to outcomes, particularly negative ones, compared to humans.

## Decision-making under risk

The Cambridge Gambling Task assessed decision-making under risk by providing explicit information about the gains and losses associated with each choice, while the outcome of each choice remained probabilistic. Players were shown a row of ten boxes, divided into two types (red and blue), with a hidden gold coin randomly assigned to one box in each round. Players predicted which box type contained the coin by placing a proportion of their bets. The ratio of red to blue boxes, which was explicitly shown, varied from a weak asymmetry (e.g., 6:4) to a strong asymmetry (e.g., 9:1), affecting risk levels. A stronger asymmetry increased the probability of the coin being in the majority type, thereby reducing risk.

We used total scores to assess the overall performance in this test. As shown in Fig. 2A, Claude, DeepSeek, and the two GPT models outperformed, or at least matched, human participants in total scores. In particular, both GPTo4m and DeepSeek achieved the highest total scores, with no statistically significant difference between them. Their performance was closer to that of a strategy optimized for expected utility maximization (see Methods for details). In contrast, Gemini performed worse than human participants, yielding the lowest overall score among the LLMs. Following the convention [45], we assessed decision-making quality by calculating the proportion of rounds in which players selected the majority type, which reflected the ability to make probabilistically informed decisions. As shown in *SI Appendix,* Fig. S2, all the LLMs displayed near-ceiling decision-making quality, and consistently chose the majority type with proportions close to 1, regardless of the degree of asymmetry in box distributions. In contrast, humans were unable to consistently choose the majority type, particularly in weak-asymmetry (e.g., a 6:4 red-to-blue ratio, representing high-risk) conditions. These findings highlight the superior accuracy and consistency of LLMs in predicting the most likely outcome. In contrast, humans showed greater variability and a higher tendency toward suboptimal choices, especially when the associated risk is high.

As shown in Fig. 2B, under high-risk conditions with weak asymmetry, humans adopted cautious strategies by placing a low proportion (e.g., 25%) of their bets. However, as asymmetry intensified (reducing risk), humans raised their bets significantly, aligning risk-taking with favorable odds of success to maximize rewards. Thus, humans demonstrated a dynamic and flexible approach to risk adjustment. In contrast, LLMs displayed a more stable and consistent betting strategy across varying degrees of asymmetry, with less variation in risk adjustment. Among the models, both GPTo4m and DeepSeek notably placed the highest bets (around 90%) across all levels of the box distribution, indicating the strongest risk-taking tendency. Claude consistently placed high bets (above 60%), indicating a propensity for risk-taking. GPT maintained moderate bet levels (generally below 50%), suggesting a balanced but slightly risk-averse strategy. Gemini
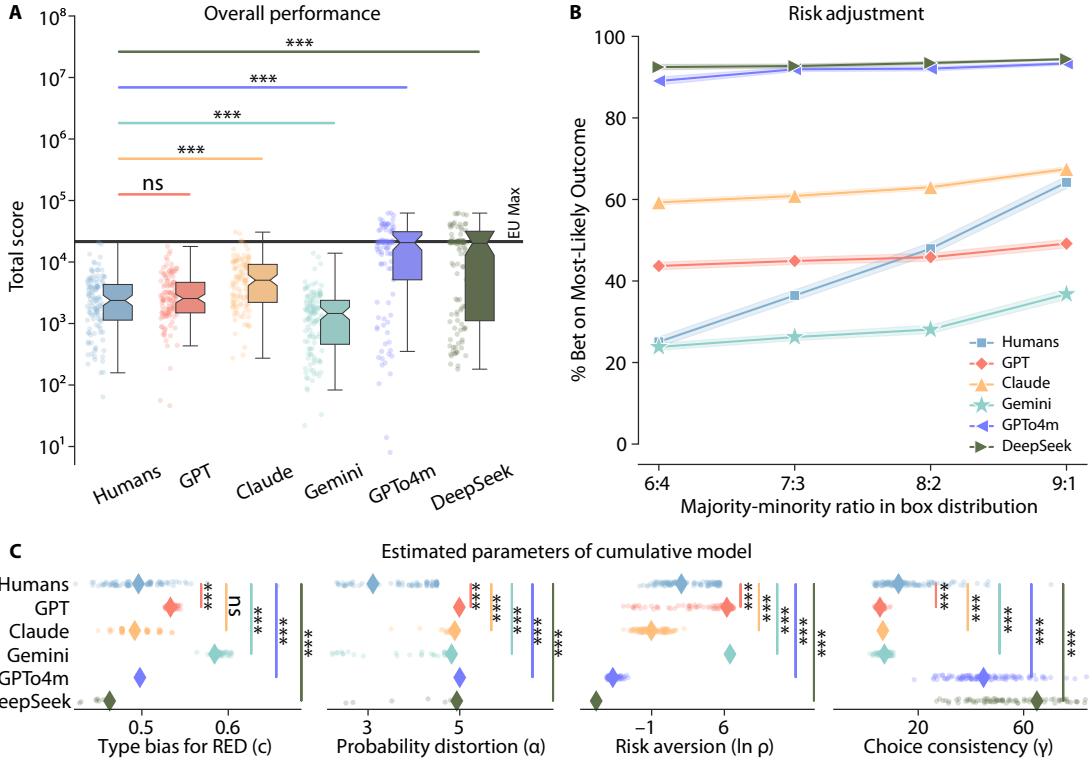
**Fig. 2**: **LLMs, except for GPT and Gemini, outperformed humans on the Cambridge Gambling Task. However, all LLMs consistently exhibited a weaker tendency for risk adjustment and distinct parameter estimates in the cumulative model compared to humans. Panel (A)** shows the total scores of LLMs, human participants, and an Expected Utility Maximization (EU-max) strategy. A Mann-Whitney U test indicated that Claude, GPTo4m, and DeepSeek outperformed human participants, GPT matched human participants, while Gemini underperformed. **Panel (B)** illustrates the risk adjustment in terms of the proportions of bets placed on the most likely outcome. LLMs displayed a more stable and consistent betting behavior across varying levels of asymmetry in box distributions. In contrast, human participants raised their bets significantly as asymmetry intensified. **Panel (C)** presents the posterior estimates of four parameters in the cumulative model: type bias for RED boxes ($c$), probability distortion ($\alpha$), risk aversion ($\rho$), and choice consistency ($\gamma$). A series of Mann–Whitney U tests revealed that, compared to human participants, GPT, GPTo4m, and Gemini exhibited significantly higher type bias for one choice ($c$), while DeepSeek showed a significantly lower value in the opposite direction. Claude did not differ significantly from human participants on this parameter. As for probability distortion ($\alpha$), all five LLMs displayed significantly greater distortion compared to humans. Turning to risk aversion ($\rho$), GPT and Gemini exhibited significantly higher values than human participants, while Claude, GPT4o4m, and DeepSeek showed significantly lower values. Lastly, GPTo4m and DeepSeek exhibited significantly higher choice consistency, while GPT, Claude, Gemini, and humans showed relatively more variability. All inferential (Mann–Whitney U tests) and descriptive statistics (means, medians, and standard deviations) for the Cambridge Gambling Task are reported in *SI Appendix,* Table S2, including pairwise comparisons (panels a–b) and parameter estimates (panel c). All parameter estimates demonstrated satisfactory convergence with R-hat ($\hat{R}$) values below 1.01.

exhibited the most risk-averse behavior, with consistently low bets (typically under 30%), but showed a noticeable increase in low-risk conditions. These results highlight a key distinction: humans flexibly adjusted their risk-taking in response to changing probabilities, while LLMs followed more rigid, consistent strategies regardless of risk level.

The Cumulative Model [54] accounts for players' decisions in this task using a probabilistic choice process (*SI Appendix, Supplementary Note 2*). This model incorporates four parameters: type bias for red ($c$), probability distortion ($\alpha$), risk aversion ($\rho$), and choice consistency ($\gamma$). The posterior estimates of these parameters showed significant differences between humans and LLMs (Fig. 2C). Human participants, Claude, and GPT4o4m showed minimal or no bias between response options, with type bias values close to 0.5. In contrast, Gemini exhibited the strongest bias toward red boxes, while DeepSeek showed a similarly strong bias in the opposite direction, with estimates substantially diverging from 0.5. For probability distortion ($\alpha$), all the LLMs exhibited significantly higher levels of distortion than humans, indicating that LLMs

tend to perceive high-probability events as even more likely and low-probability events as even less likely than they objectively are. The risk aversion ($\rho$), however, differed across LLMs: GPT and Gemini showed greater risk aversion, whereas Claude, GPTo4m, and DeepSeek showed strong risk-seeking tendencies. In contrast, human participants were more balanced—less risk-averse than GPT and Gemini, but also less risk-seeking than other LLMs. Finally, choice consistency ($\gamma$) showed that GPTo4m and DeepSeek followed highly deterministic strategies, closely aligning with model-predicted expected values, while humans and other LLMs exhibited more variability in their choices.

## Decision-making under set-shifting

The Wisconsin Card Sorting Task assessed decision-making with changing conditions. Participants were presented with items that vary in attributes such as color, shape, and the number of symbols. They were asked to match items to one of four cards based on a matching rule (i.e., an attribute), whereas the rule was not explicitly stated. Participants received feedback ('correct' or 'incorrect') for their every match. Upon successfully achieving a predetermined number of correct matches, the rule would change without notice. Thus, this task required participants to identify the rule, detect the change, and continuously adapt to the new rule.

To assess performance, we examined the total number of correct matches. As shown in Fig. 3A, all the LLMs, except for DeepSeek, significantly outperformed human participants, and approached the performance of a strategy optimized for expected utility maximization (see Methods for details). Specifically, Gemini and GPTo4m achieved the highest median number of correct matches. Gemini exhibited the lowest variance, followed closely by GPTo4m, indicating that both models performed accurately and consistently.

We further analyzed their decision-making using four commonly used metrics [55]: the number of rounds to complete the first set (TRSET1), failure to maintain set (FSET), perseverative errors, and non-perseverative errors. As indicated by TRSET1, compared to humans, all the LLMs, except for DeepSeek, used fewer rounds to complete the first set, suggesting that almost all LLMs more efficiently identified and applied the underlying rule (*SI Appendix,* Fig. S3). The FSET measures how often participants changed their matching strategy after achieving five or more consecutive correct matches without negative feedback ('incorrect'). Humans and LLMs displayed comparable values of FSET, suggesting that they were equally capable of maintaining a correct strategy once identified (*SI Appendix,* Fig. S3). The perseverative errors are the number of incorrect matches caused by applying a previously learned rule to a newly changed condition. In contrast, the non-perseverative errors are random errors unrelated to the task. As shown in Fig. 3B, human participants tended to make more non-perseverative errors than perseverative ones. However, LLMs, except for DeepSeek, exhibited the opposite pattern and produced more perseverative errors than non-perseverative ones. DeepSeek, in contrast, matched the pattern of humans by showing more non-perseverative errors.

The Sequential Learning Model [56] describes how participants adjust their decision in the Wisconsin Card Sorting Task through three parameters: reward sensitivity ($r$), punishment sensitivity ($p$), and choice consistency ($d$) (*SI Appendix, Supplementary Note 2*). As shown in Fig. 3C, the posterior estimates for these parameters revealed significant differences between LLMs and humans. For both reward sensitivity ($r$) and punishment sensitivity ($p$), LLMs typically exhibited higher values than humans, suggesting that they adapted more quickly to both positive and negative feedback. Choice consistency ($d$) showed that LLMs generally made more deterministic decisions consistent with model-predicted expected values, in contrast to the greater variability observed in human choices. The only exception was DeepSeek, which did not differ significantly from humans in choice consistency and was less sensitive to rewards compared to humans. However, it was more sensitive to punishment than humans.

# Discussion

By virtue of passing the Turing test, LLMs have to possess decision-making capabilities, but most of the actual studies about them have been restricted to domain-specific settings, where multiple decision dimensions are often intertwined. Focusing on decision-making under uncertainty, risk, and set-shifting, we used three standard psychological tests to investigate five LLMs. Notably, four of them—GPT-4o, Claude, GPTo4m, and DeepSeek—consistently matched or exceeded human performance across all tasks. Moreover, their performance generally approached optimality. Our findings offer evidence for the general decision-making abilities of LLMs, particularly in navigating uncertainty, calibrating risk, and adapting flexibly to changing conditions.

The AI's super-human performance also comes with an apparent non-human cognition. LLMs often employed strategies that diverged substantially from those of human participants. In the Iowa Gambling Task, LLMs were more sensitive to the frequency of losses, leading to divergent preferences between the two advantageous decks. In contrast, human choices were less influenced by loss frequency and showed a more balanced distribution across those decks. In the Cambridge Gambling Task, LLMs made optimal
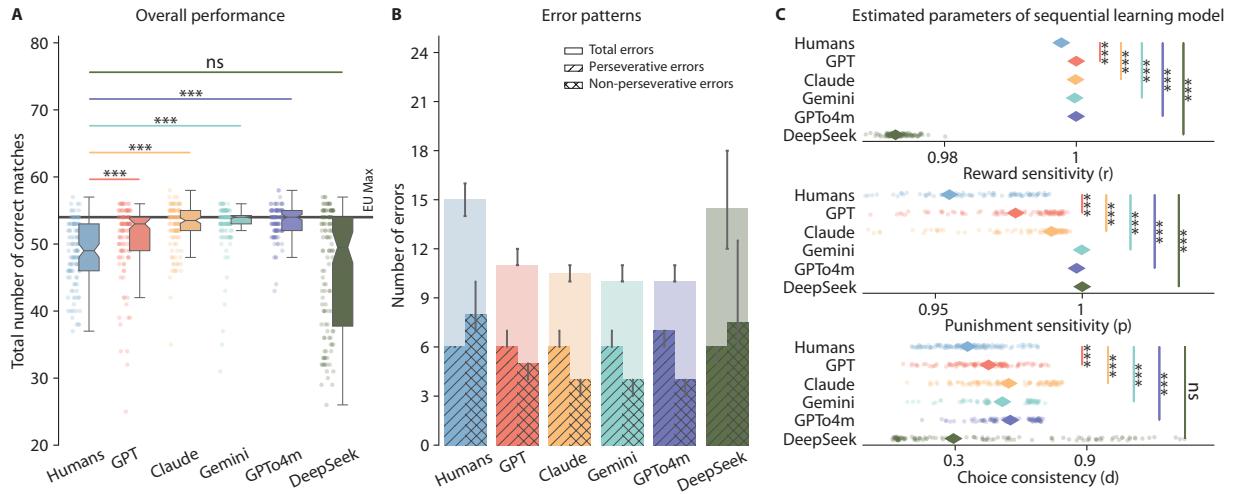
**Fig. 3**: **All the LLMs outperformed, or at least matched, humans in the Wisconsin Card Sorting Task, while exhibiting generally distinct error patterns and parameter estimates in the sequential learning model compared to humans.** **Panel (A)** shows the total number of correct matches for human participants, LLMs, and an Expected Utility Maximization (EU-max) strategy. A two-tailed Mann-Whitney U test revealed that GPT, Claude, Gemini and GPTo4m all significantly outperformed human participants; however, DeepSeek matched human participants. **Panel (B)** illustrates comparisons of median across three error types made by human participants and LLMs. Bars represent the median number of total, perseverative, and non-perseverative errors, with 95% confidence intervals of the median shown as error bars. Human participants made significantly more non-perseverative errors unrelated to the task than perseverative errors. In contrast, except for DeepSeek, LLMs produced more perseverative errors than non-perseverative ones. The error pattern of DeepSeek matched that of human participants, with more non-perseverative errors than perseverative ones. **Panel (C)** presents the posterior estimates of the parameters in the sequential learning model. Mann-Whitney U tests revealed that compared to human participants, LLMs except for DeepSeek demonstrated higher reward sensitivity ($r$), greater sensitivity to punishment ($p$), and higher choice consistency ($d$). As an exception, compared to humans, DeepSeek showed lower reward sensitivity, higher punishment sensitivity, and a comparable level of choice consistency. All inferential (Mann–Whitney U tests) and descriptive statistics (means, medians, and standard deviations) for the Wisconsin Card Sorting Task are reported in *SI Appendix*, Table S3, including pairwise comparisons (panels a–c) and parameter estimates (panel d). All parameter estimates demonstrated satisfactory convergence with R-hat ($\hat{R}$) values below 1.01.

probabilistic choices by consistently selecting the option with the highest expected return, yet showed minimal adjustment in their betting strategies in response to changing levels of risk. In contrast, human participants exhibited greater variability in their betting behavior despite making less optimal probabilistic choices. In the Wisconsin Card Sorting Task, LLMs adapted to rule changes faster and made fewer random errors, while humans were more prone to make mistakes. Taken together, these findings suggest that LLMs act as if relying on cognitive processes that diverge substantially from those of humans. Therefore, replacing humans with machines in decision-making at a large scale will have profound consequences for our society.

Our results underscore the risks associated with using LLMs as substitutes for human participants in behavioral research and practical decision-making. The models not only fail to reproduce typical human decision-making but also can generate misleading signals. A notable example is their consistent lack of risk adjustment, a trait that most humans exhibit. This deficiency is not primarily a deviation from normative behavior; rather, it is a sensitive trans-diagnostic marker of cognitive dysfunction, widely seen across a range of psychiatric conditions such as addiction and depression [57, 58]. As such, relying on LLMs in tasks involving probabilistic reasoning can obscure important psychological insights. Moreover, LLMs also do not capture the breadth of reasoning strategies typical of human cognition. Their generated responses are less diverse than those of human participants and remain largely unaffected by demographic cues or contextual framing, according to our robustness checks (see Methods, and SI Appendix, Supplementary Note 3). In contrast, human decisions are shaped by individual differences and context, such as age, gender, and cultural background [59–62]. Thus, using LLMs to replace human participants in behavioral research risks oversimplifying the phenomena under study and misrepresenting how people respond in actual decision-making contexts.

If LLMs do not faithfully replicate human cognition, what exactly are they emulating? Our parameter estimation from computational models suggests that LLMs excel at recognizing and exploiting historical patterns, enabling them to efficiently gather information under uncertainty. They also exhibit systematic probability distortions, allowing their choices under risk to align closely with the options most likely to succeed. Across tasks, LLMs tend to show heightened sensitivity to outcomes, particularly negative ones, and generally make decisions consistent with the principle of expected utility, reflecting a form of rationality. These behavioral tendencies appear to underpin their strong performance in decision-making tasks. Moreover, we found that LLMs' decision-making patterns remained largely consistent across different prompt formulations and temperature settings (*SI Appendix*, Supplementary Note 3). These findings suggest that while LLMs do not behave as cognitive replicas of human decision-making, they operate as consistent, rational, and outcome-driven agents.

Despite differences in risk attitudes (e.g., Claude being more risk-seeking, while GPT and Gemini are more risk-averse), their behavior is primarily shaped by the structural features of the tasks they encounter.

While a consistent, rational, outcome-driven agent may be appealing from a utilitarian perspective, it does not necessarily translate into human acceptance of AI. In our post-experiment survey (*SI Appendix*, Fig. S4), participants expressed little interest in seeking assistance from AI, even when they were informed of the LLMs' high performance on the same tasks and despite their own lower scores. Furthermore, participants did not believe that AI could improve their decision-making, either in terms of performance or efficiency. These findings point to a form of algorithm aversion [63–65], which may arise from concerns about trust, autonomy, or the perceived appropriateness of algorithmic input, particularly in decisions involving ambiguous outcomes, complex judgment, or personal values. Thus, even when LLMs demonstrate high technical performance, their effective integration into human decision-making processes will require attention not only to their capabilities but also to users' perceptions, interpretations, and willingness to accept algorithmic guidance.

Even if people were willing to adopt LLMs to assist in their decisions, it is important to ask: when and for what kinds of decisions can these systems provide meaningful assistance? LLMs tend to perform well in settings with clearly defined goals and outcomes, thanks to their rationality and outcome sensitivity. Yet humans, famously sub-optimal decision makers, often succeed by departing from strict rationality and discounting cost-benefit logic. Consider the Wright Brothers: their crash-prone test flights, many ending in failure, defied utility maximization principles but drastically accelerated the path to powered flight. In such cases, heuristics, imagination, and a willingness to tolerate error or loss can outperform conservative or algorithmically optimized strategies. LLMs, bound to rational evaluation functions, may struggle to replicate this productive 'irrationality,' which remains a distinctly human advantage.

Overall, decision-making is never solely about performance. Many human choices cannot be objectively evaluated, and what seems like a poor decision now may later prove wise. In such cases, what truly matters is how the decision is made—how we reason, weigh uncertainty, balance risks, and adapt to change. While our findings show that LLMs can perform well on stylized psychological decision tasks, they also reveal a crucial disconnection between human-level performance and human-like cognition. This suggests that LLM intelligence should not be viewed as a mirror of human thought, but rather as a novel form of reasoning. For system designers, the issue is not just performance, but whether people are willing to rely on LLMs for decision-making, especially when LLMs do not think or act like them. For policymakers, the challenge goes beyond evaluating performance: it requires assessing whether, and to what extent, these systems reason in ways that align with human values and norms, and uphold accountability. For everyday users, the question is how much autonomy we are willing to share with these systems, especially in decisions with real consequences. After all, a decision is also about what kind of reasoning we are willing to trust, and the kind of minds we are willing to build. These decisions not only reflect who we are but also shape who we may become, and at times, they may alter the course of history.

## Methods

### *Iowa Gambling Task*

This task tests if participants can prioritize long-term benefits over short-term gains in an uncertain setting [44]. In our experiments, this task lasted for 80 rounds and comprised four decks (labeled as A, B, C, and D), from which participants must choose one deck per round. Each choice offered an immediate reward, with the possibility of an occasional penalty. Decks A and B were disadvantageous decks, while C and D were advantageous decks. Specifically, each choice of A or B resulted in a high immediate reward, but the penalty, when it occurred, was substantial, offsetting the immediate rewards in the long run. As a result, the expected payoffs from these decks were lower. Conversely, each choice of C or D resulted in a low immediate reward, though the penalty, when it occurred, was minor, yielding a higher expected payoff in the long run.

Additionally, compared to Deck A (or C), the penalty from Deck B (or D) occurred less frequently. Participants were not told the total number of rounds or the reward and penalty associated with each deck. Their decisions were made under uncertainty and relied on learning over time. This task can be mathematically framed as a multi-armed bandit problem, for which the upper confidence bound [51] and $\epsilon$-greedy [52] are two well-established solution strategies.

### Cambridge Gambling Task

This task tests decision-making under risk [45]. In our experiments, this task consisted of 64 rounds in total, with every 8 rounds forming a sub-session. Participants saw a row of 10 boxes categorized into two types: Red and Blue; they must choose under which box type (Red or Blue) a gold coin was hidden and must place a bet on the type chosen. The ratio of red to blue boxes was explicitly shown and indicated the risk associated with each choice. There were 8 possible red-to-blue ratios: 1:9, 2:8, 3:7, 4:6, 6:4, 7:3, 8:2, and 9:1. A stronger asymmetry (e.g., 1:9) in the red-to-blue ratio indicated a higher probability of the coin being in the majority type, suggesting a lower risk in choosing the majority type. Bet levels were fixed at 5%, 25%, 50%, 75%, or 95% of their current score. A correct guess resulted in gaining scores equal to the bet, while an incorrect guess resulted in losing scores equal to the bet. Different from the original design, we presented all five bet options simultaneously, rather than sequentially, as the sequential presentation aimed to test human impulsive behaviors, which were not applicable to LLMs. Participants had complete information about the number of rounds, the risk (the red-to-blue ratios), and the possible rewards and penalties associated with each choice. The expected utility maximization strategy in this context is to always choose the majority box type and place the highest possible (95%) bet. However, this strategy does not reflect typical decision-making behavior in healthy individuals, as it lacks risk adjustment [57].

### Wisconsin Card Sorting Task

This task tests participants' ability to make decisions under set-shifting conditions [46]. In our experiments, this task lasted 64 rounds, and comprised four cards (labeled as A, B, C, and D). In each round, participants were presented with an item and required to choose one card that matched the item's pattern. Both the cards and items were marked with a set of symbols, and the matching pattern was based on one of three attributes of the symbols: color, shape, or number. A correct match resulted in a reward, while an incorrect match resulted in no points. The matching rule remained consistent for a sequence of eight consecutive correct responses, after which it changed without explicit notification. Participants were informed that the matching rule could change but were not told when or how. This design required participants to adapt their decision-making strategy in response to shifting conditions over time. The expected utility maximization strategy for this task is to switch to a different attribute after receiving negative feedback and to continue using the current attribute following positive feedback.

### Large Language Models

We provided LLMs with a system prompt and a decision-making prompt (see *SI Appendix*, Supplementary Note 1 for prompt details). Through the system prompt, LLMs received exactly the same experimental instructions as human participants. The decision-making prompt presented LLMs with their previous choices and the corresponding outcomes from past rounds, and asked LLMs to make a choice for the current round. However, just like human participants, LLMs were not given explicit instructions on how to utilize the provided information. We maintained the default parameters for all LLMs unless specified otherwise, as this better simulates typical user conditions, which served as the baseline. To prevent LLMs from exploiting biases to solve the tasks [66], we performed a cyclic permutation of the order of options for all tasks.

### Robustness Checks

To assess if our findings on LLMs are robust to prompt variations, we conducted robustness checks by varying the experimental settings as follows (see *SI Appendix*, Supplementary Note 3 for details): i) adjusting the LLMs' temperature from the default value (i.e. 1) to 0 and 0.5, ii) applying arithmetic transformation to the scores of each round (e.g., multiplying them by a constant factor or adding a constant value), iii) restructuring the prompts to reflect various decision-making contexts (including both economic and medical scenarios), and iv) introducing role-play prompts to test the impact of demographic information (e.g., different age ranges, genders, and ethnicity) and risk preferences (e.g., risk-taking and risk-averse). In total, there were 19 variants in the Iowa Gambling Task as well as the Cambridge Gambling Task, and 15 variants in the Wisconsin Card Sorting Task. We repeated each variant 10 times using GPT-4o. Overall, we observed that the decision-making patterns of LLMs remained qualitatively unchanged across these variants ( *SI Appendix*, Fig. S5, S6, S7, S8).

### Human Subjects

We recruited a total of 360 participants, including 37.8% women, with a mean age of 18.95 years (Table S4). The experiments were pre-registered (AsPredicted #182473, #186129, and #203115) and were conducted at Northwest A&F University and Northwestern Polytechnical University in Xi'an, from July to December 2024. Participants came from multiple departments to minimize interaction among them. Each task was administered using oTree [67] and completed by 120 participants. Details of these tasks were maintained under strict confidentiality until participants arrived at the lab. Upon arrival, each participant was allocated a computer, isolated from others by partitions to guarantee independent completion. Participants began by a tutorial that described the task (*SI Appendix,* Fig. S9, S10, S11). Then, the experiment, which lasted for multiple rounds, started. Each round comprised a choice page, where participants made a decision in a task, and a result page, where the score of each round and the cumulated score were shown (*SI Appendix,* Fig. S12, S13, S14, S15). After completing all rounds, participants saw their final scores (*SI Appendix,* Fig. S16), and then their demographic data (*SI Appendix,* Fig. S17), feedback on the task, and attitudes toward artificial intelligence (*SI Appendix,* Fig. S18) were collected. Throughout the experiment, communication between participants was strictly prohibited. Each session lasted approximately 30 minutes. We provided an average payment of 30 CNY, which included a 15 CNY show-up fee, with the remaining amount being calculated based on experiment scores.

### Ethics Statement

This study was approved by the Northwestern Polytechnical University Ethics Committee on the use of human participants in research and carried out in accordance with all relevant guidelines. Informed consent was obtained from all participants.

### Data, Materials, and Software Availability.

Data, Materials, and Software Availability. Data and code for the current study are available through the GitHub repository.

# References

[1] Simon, H.A.: Administrative Behavior: A Study of Decision-Making Processes in Administrative Organization. The Free Press, Glencoe (1947)

[2] Turban, E., Watkins, P.R.: Integrating expert systems and decision support systems. MIS Q. **10**(2), 121–136 (1986)

[3] Jones, C.R., Bergen, B.K.: Large Language Models Pass the Turing Test. Preprint arXiv:2503.23674 (2025)

[4] Mei, Q., Xie, Y., Yuan, W., Jackson, M.O.: A turing test of whether ai chatbots are behaviorally similar to humans. Proceedings of the National Academy of Sciences **121**(9), 2313925121 (2024)

[5] Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., Arx, S., Bernstein, M.S., Bohg, J., Bosselut, A., Brunskill, E., et al.: On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258 (2021)

[6] Handler, A., Larsen, K.R., Hackathorn, R.: Large language models present new questions for decision support. Int. J. Info. Manag. **79**, 102811 (2024)

[7] Chen, Z., Ma, J., Zhang, X., Hao, N., Yan, A., Nourbakhsh, A., Yang, X., McAuley, J., Petzold, L.R., Wang, W.Y.: A Survey on Large Language Models for Critical Societal Domains: Finance, Healthcare, and Law. Preprint arXiv:2405.01769 (2024)

[8] Taylor, L.: Colombian judge says he used ChatGPT in ruling. Technology section (2023). https://www.theguardian.com/technology/2023/feb/03/colombia-judge-chatgpt-ruling Accessed 2025-05-08

[9] Alba, D., Yin, L.: Uncovering What OpenAI's GPT Sees When Used For Ranking Resumes. Bloomberg Green Daily newsletter article (2024). https://www.bloomberg.com/news/newsletters/2024-03-08/companies-should-think-twice-before-using-generative-ai-in-hiring Accessed 2025-05-08

[10] Singla, A., Sukharevsky, A., Yee, L., Chui, M., Hall, B.: The state of AI: How organizations are rewiring to capture value. Accessed 2025-05-11 (2025). https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-state-of-ai

[11] Noy, S., Zhang, W.: Experimental evidence on the productivity effects of generative artificial intelligence. Science **381**(6654), 187–192 (2023)

[12] Extance, A.: Chatgpt has entered the classroom: how llms could transform education. Nature **623**(7987), 474–477 (2023)

[13] Williams, C.Y., Miao, B.Y., Kornblith, A.E., Butte, A.J.: Evaluating the use of large language models to provide clinical recommendations in the emergency department. Nature Communications **15**(1), 8236 (2024)

[14] Goh, E., Gallo, R.J., Strong, E., Weng, Y., Kerman, H., Freed, J.A., Cool, J.A., Kanjee, Z., Lane, K.P., Parsons, A.S., et al.: Gpt-4 assistance for improvement of physician performance on patient care tasks: a randomized controlled trial. Nat. Med. **31**, 1233–1238 (2025)

[15] Singhal, K., Azizi, S., Tu, T., Mahdavi, S.S., Wei, J., Chung, H.W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S., et al.: Large language models encode clinical knowledge. Nature **620**(7972), 172–180 (2023)

[16] Hager, P., Jungmann, F., Holland, R., Bhagat, K., Hubrecht, I., Knauer, M., Vielhauer, J., Makowski, M., Braren, R., Kaissis, G., et al.: Evaluation and mitigation of the limitations of large language models in clinical decision-making. Nat. Med. **30**(9), 2613–2622 (2024)

[17] Chen, Z., Chen, W., Smiley, C., Shah, S., Borova, I., Langdon, D., Moussa, R., Beane, M., Huang, T.-H., Routledge, B., Wang, W.Y.: FinQA: A dataset of numerical reasoning over financial data. In: Moens, M.-F., Huang, X., Specia, L., Yih, S.W.-t. (eds.) Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (2021)

[18] Xie, Q., Han, W., Chen, Z., Xiang, R., Zhang, X., He, Y., Xiao, M., Li, D., Dai, Y., Feng, D., *et al.*: Finben: A holistic financial benchmark for large language models. Adv. Neural Inf. Process. Syst. **37**, 95716–95743 (2024)

[19] Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., Anandkumar, A.: Voyager: An open-ended embodied agent with large language models. Transact. Mach. Learn. Res. **3** (2024)

[20] Tan, W., Zhang, W., Xu, X., Xia, H., Ding, Z., Li, B., Zhou, B., Yue, J., Jiang, J., Li, Y., et al.: Cradle: Empowering foundation agents towards general computer control. arXiv preprint arXiv:2403.03186 (2024)

[21] Akata, E., Schulz, L., Coda-Forno, J., Oh, S.J., Bethge, M., Schulz, E.: Playing repeated games with large language models. online ahead of print in Nat. Hum. Behav. (2025)

[22] Wang, Z., Song, R., Shen, C., Yin, S., Song, Z., Battu, B., Shi, L., Jia, D., Rahwan, T., Hu, S.: Large language models overcome the machine penalty when acting fairly but not when acting selfishly or altruistically. arXiv preprint arXiv:2410.03724 (2024)

[23] Johnson, J.G., Busemeyer, J.R.: Decision making under risk and uncertainty. Wiley Interdiscip. Rev. Cogn. Sci. **1**(5), 736–749 (2010)

[24] Hastie, R., Dawes, R.M.: Rational Choice in an Uncertain World: The Psychology of Judgment and Decision Making. Sage, ??? (2010)

[25] Einhorn, H.J., Hogarth, R.M.: Behavioral decision theory: Processes of judgement and choice. Annu. Rev. Psychol. **32**(1981), 53–88 (1981)

[26] Knight, F.H.: Risk, Uncertainty and Profit vol. 31. Houghton Mifflin, ??? (1921)

[27] Camerer, C., Weber, M.: Recent developments in modeling preferences: Uncertainty and ambiguity. J. Risk Uncertain. **5**, 325–370 (1992)

[28] Machina, M.J.: Decision-making in the presence of risk. Science **236**(4801), 537–543 (1987)

[29] Kahneman, D., Tversky, A.: Prospect theory: An analysis of decision under risk. In: Handbook of the Fundamentals of Financial Decision Making: Part I, pp. 99–127. World Scientific, ??? (2013)

[30] Brehmer, B.: Dynamic decision making: Human control of complex systems. Acta Psychol. **81**(3), 211–241 (1992)

[31] Payne, J.W., Bettman, J.R., Johnson, E.J.: The Adaptive Decision Maker. Cambridge University Press, ??? (1993)

[32] Kochenderfer, M.J. (ed.): Decision Making Under Uncertainty: Theory and Application. MIT Press, Cambridge MA (2015)

[33] Ruggeri, K., Alí, S., Berge, M.L., Bertoldo, G., Bjørndal, L.D., Cortijos-Bernabeu, A., Davison, C., Demić, E., Esteban-Serna, C., Friedemann, M., *et al.*: Replicating patterns of prospect theory for decision under risk. Nat. Hum. Behav. **4**(6), 622–633 (2020)

[34] Uddin, L.Q.: Cognitive and behavioural flexibility: neural mechanisms and clinical considerations. Nat. Rev. Neurosci. **22**(3), 167–179 (2021)

[35] Strategy, A.: Deciding advantageously before knowing the. Science **275**, 1293–1293 (1997)

[36] Sacré, P., Kerr, M.S., Subramanian, S., Fitzgerald, Z., Kahn, K., Johnson, M.A., Niebur, E., Eden, U.T., González-Martínez, J.A., Gale, J.T., *et al.*: Risk-taking bias in human decision-making is encoded via a right–left brain push–pull system. Proc. Natl. Acad. Sci. USA **116**(4), 1404–1413 (2019)

[37] Konishi, S., Nakajima, K., Uchida, I., Kameyama, M., Nakahara, K., Sekihara, K., Miyashita, Y.: Transient activation of inferior prefrontal cortex during cognitive set shifting. Nat. Neurosci. **1**(1), 80–84 (1998)

[38] Binz, M., Schulz, E.: Using cognitive psychology to understand gpt-3. Proceedings of the National

Academy of Sciences **120**(6), 2218523120 (2023)

[39] Hagendorff, T.: Machine psychology: Investigating emergent capabilities and behavior in large language models using psychological methods. Preprint arXiv:2303.13988 (2023)

[40] Kosinski, M.: Evaluating large language models in theory of mind tasks. Proceedings of the National Academy of Sciences **121**(45), 2405460121 (2024)

[41] Hagendorff, T.: Deception abilities emerged in large language models. Proceedings of the National Academy of Sciences **121**(24), 2317967121 (2024)

[42] Chen, Y., Liu, T.X., Shan, Y., Zhong, S.: The emergence of economic rationality of gpt. Proc. Natl. Acad. Sci. USA **120**(51), 2316205120 (2023)

[43] Lehr, S.A., Saichandran, K.S., Harmon-Jones, E., Vitali, N., Banaji, M.R.: Kernels of selfhood: Gpt-4o shows humanlike patterns of cognitive dissonance moderated by free choice. Proceedings of the National Academy of Sciences **122**(20), 2501823122 (2025)

[44] Bechara, A., Damasio, A.R., Damasio, H., Anderson, S.W.: Insensitivity to future consequences following damage to human prefrontal cortex. Cognition **50**(1-3), 7–15 (1994)

[45] Rogers, R.D., Everitt, B., Baldacchino, A., Blackshaw, A.J., Swainson, R., Wynne, K., Baker, N., Hunter, J., Carthy, T., Booker, E., *et al.*: Dissociable deficits in the decision-making cognition of chronic amphetamine abusers, opiate abusers, patients with focal damage to prefrontal cortex, and tryptophan-depleted normal volunteers: evidence for monoaminergic mechanisms. Neuropsychopharmacology **20**(4), 322–339 (1999)

[46] Berg, E.A.: A simple objective technique for measuring flexibility in thinking. J. Gen. Psychol. **39**(1), 15–22 (1948)

[47] Abdurahman, S., Atari, M., Karimi-Malekabadi, F., Xue, M.J., Trager, J., Park, P.S., Golazizian, P., Omrani, A., Dehghani, M.: Perils and opportunities in using large language models in psychological research. PNAS Nexus **3**(7), 245 (2024)

[48] Wang, A., Morgenstern, J., Dickerson, J.P.: Large language models that replace human participants can harmfully misportray and flatten identity groups. Nat. Mach. Intell. **7**, 400–411 (2025)

[49] Grossmann, I., Feinberg, M., Parker, D.C., Christakis, N.A., Tetlock, P.E., Cunningham, W.A.: Ai and the transformation of social science research. Science **380**(6650), 1108–1109 (2023)

[50] Bull, P.N., Tippett, L.J., Addis, D.R.: Decision making in healthy participants on the iowa gambling task: new insights from an operant approach. Front. Psychol. **6**, 391 (2015)

[51] Auer, P., Cesa-Bianchi, N., Fischer, P.: Finite-time analysis of the multiarmed bandit problem. Mach. Learn. **47**, 235–256 (2002)

[52] Watkins, C.J.C.H., et al.: Learning from delayed rewards. PhD thesis (1989)

[53] Ahn, W.-Y., Busemeyer, J.R., Wagenmakers, E.-J., Stout, J.C.: Comparison of decision learning models using the generalization criterion method. Cogn. Sci. **32**(8), 1376–1402 (2008)

[54] Romeu, R.J., Haines, N., Ahn, W.-Y., Busemeyer, J.R., Vassileva, J.: A computational model of the cambridge gambling task with applications to substance use disorders. Drug Alcohol Depend. **206**, 107711 (2020)

[55] Gläscher, J., Adolphs, R., Tranel, D.: Model-based lesion mapping of cognitive control using the wisconsin card sorting test. Nat. Commun. **10**(1), 20 (2019)

[56] Bishara, A.J., Kruschke, J.K., Stout, J.C., Bechara, A., McCabe, D.P., Busemeyer, J.R.: Sequential learning models for the wisconsin card sort task: Assessing processes in substance dependent individuals. J. Math. Psychol. **54**(1), 5–13 (2010)

[57] Clark, L., Bechara, A., Damasio, H., Aitken, M., Sahakian, B., Robbins, T.: Differential effects of insular

and ventromedial prefrontal cortex lesions on risky decision-making. Brain **131**(5), 1311–1322 (2008)

[58] Effah, R., Ioannidis, K., Grant, J.E., Chamberlain, S.: Exploring decision-making performance in young adults with mental health disorders: a comparative study using the cambridge gambling task. Psychol. Med. **54**(9), 1890–1896 (2024)

[59] Byrnes, J.P., Miller, D.C., Schafer, W.D.: Gender differences in risk taking: A meta-analysis. Psychol. Bull. **125**(3), 367 (1999)

[60] Weber, E.U., Hsee, C.: Cross-cultural differences in risk perception, but cross-cultural similarities in attitudes towards perceived risk. Manag. Sci. **44**(9), 1205–1217 (1998)

[61] Tversky, A., Kahneman, D.: The framing of decisions and the psychology of choice. Science **211**(4481), 453–458 (1981)

[62] Tymula, A., Rosenberg Belmaker, L.A., Ruderman, L., Glimcher, P.W., Levy, I.: Like cognitive function, decision making across the life span shows profound age-related changes. Proc. Natl. Acad. Sci. USA **110**(42), 17143–17148 (2013)

[63] Dietvorst, B.J., Simmons, J.P., Massey, C.: Algorithm aversion: people erroneously avoid algorithms after seeing them err. J. Exp. Psychol. Gen. **144**(1), 114 (2015)

[64] Castelo, N., Bos, M.W., Lehmann, D.R.: Task-dependent algorithm aversion. J. Market. Res. **56**(5), 809–825 (2019)

[65] Karataş, M., Cutright, K.M.: Thinking about god increases acceptance of artificial intelligence in decision-making. Proc. Natl. Acad. Sci. USA **120**(33), 2218961120 (2023)

[66] Zhao, Z., Wallace, E., Feng, S., Klein, D., Singh, S.: Calibrate before use: Improving few-shot performance of language models. In: International Conference on Machine Learning, pp. 12697–12706 (2021). PMLR

[67] Chen, D.L., Schonger, M., Wickens, C.: otree—an open-source platform for laboratory, online, and field experiments. Journal of Behavioral and Experimental Finance **9**, 88–97 (2016)

[68] Carpenter, B., Gelman, A., Hoffman, M.D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M.A., Guo, J., Li, P., Riddell, A.: Stan: A probabilistic programming language. J. Stat. Softw. **76**, 1–32 (2017)

[69] Ahn, W.-Y., Haines, N., Zhang, L.: Revealing neurocomputational mechanisms of reinforcement learning and decision-making with the hBayesDM package. Comp. Psychiatr. **1**, 24–57 (2017)

# Supplementary note 1: Prompt for LLMs

The prompts provided to the LLMs consist of two parts: the system prompt and the decision-making prompt. Through the system prompt, LLMs received exactly the same experimental instructions as human subjects, with no additional information about tasks. The decision-making prompt presents LLMs with their previous choices and the corresponding outcomes from past rounds, and asks LLMs to make a choice for the current round. However, just like participants, LLMs are not given explicit instructions on how to utilize the provided information.

## Iowa Gambling Task

Fig. S20 shows the outcomes for each choice of these decks. `<choice_{i}>`, `<reward_{i}>`, and `<penalty_{i}>` represent the choice made, the reward received, and the penalty incurred (if any) in the $i$-th round, respectively. `<points_so_far>` indicates the cumulative payoff accumulated by the LLM up to the current round.

---

**System prompt**

```
In this game, you find yourself in a mysterious room with four ancient treasure
↪  chests. Opening each chest will yield a reward but may also simultaneously
↪  result in a penalty, depending on the chosen chest. With each turn, you will
↪  choose one chest to open. Please consider carefully, as your choice may
↪  significantly impact your points. Specifically, the rewards will increase
↪  your points, while penalties will deduct your points. At the start of the
↪  game, you will receive a loan of 2000 points. The game has several rounds in
↪  which your points will accumulate, and your goal is to maximize your points
↪  by the end of the game.

The only hint I can give you, and the most important thing to note is this: Out
↪  of these chests, there are some that are worse than others, and to win you
↪  should try to stay away from bad chests. No matter how much you find
↪  yourself losing, you can still win the game if you avoid the worst chests.
↪  Also note that the computer does not change the order of the chests once the
↪  game begins. It does not make you lose at random, or make you lose money
↪  based on the last chest you picked.

Your response must always present in the following format:

<reasoning>Reasons for your choice this round</reasoning>

<choice>Any number between 1-4 indicates your choice of chest</choice>
```

---

**Decision-making prompt**

```
Here is the historical information from the past round(s), and you may use it as
↪  a reference for your following choice.

In round 1, you chose chest number <choice_{1}>. You earned <reward_{1}>  points
↪  in rewards (and received a penalty of <penalty_{1}> points).
In round 2, you chose chest number <choice_{2}>. You earned <reward_{2}>  points
↪  in rewards (and received a penalty of <penalty_{2}> points).

                                    ⋮

In round i-1, you chose chest number <choice_{i-1}>. You earned <reward_{i-1}>
↪  points in rewards (and received a penalty of <penalty_{i-1}> points).

Your total points so far: <points_so_far> points. Now this is the i-th round of
↪  the game. Please make your choice.
```

# Cambridge Gambling Task

`<round_interactions>` represents the number of rounds in each subsession. `<choice_type_{i}>` represents the type of the box chosen in the $i$-th round, and `<choice_percent_{i}>` represents the bet proportion in the $i$-th round. `<coin_type_{i}>` represents the type of the box where the coin is placed. `<payoff_{i}>` represents the reward (or penalty) received in the $i$-th round. `<remain_by_subsession>` represents the player's total earnings in the current subsession. `<num_type_RED_boxes>` and `<num_type_BLUE_boxes>` represent the number of Type RED and Type BLUE boxes in the current round.

---

**System prompt**

```
In front of you are 10 treasure chests, some of which are type F and others
↪   are type J. A gold coin is hidden in one of these chests. Your task is to
↪   guess in which type of chest the gold coin is hidden.
Every <round_interactions> rounds constitute a phase. At the start of each
↪   phase, your points will be set to 100. Your goal is to maximize the points
↪   in each phase, and the sum of the points from each phase is your total
↪   score in the game.

Each round, you need to bet by deducting a certain percentage from your
↪   existing points to choose the type of chest. The available percentages
↪   are: 5%, 25%, 50%, 75%, 95%. If the gold coin is hidden in any one of the
↪   chests of the type you guessed, you win; otherwise, you lose. If you win,
↪   you will receive double the points of your bet; if you lose, you gain
↪   nothing. Points are rounded to the nearest whole number.
An important game setting is: the position of the gold coin is completely
↪   random in each round and has no connection to previous locations. That is,
↪   each round's coin position is an independent event.

Here is a summary of the accessible choices, strictly limited to the following
↪   options. Each choice maps to a specific type guess and bet percentage,
↪   with the list of choices and their corresponding as follows:

Choice 0 maps to F, 5% bet.
Choice 1 maps to F, 25% bet.
Choice 2 maps to F, 50% bet.
Choice 3 maps to F, 75% bet.
Choice 4 maps to F, 95% bet.
Choice 5 maps to J, 5% bet.
Choice 6 maps to J, 25% bet.
Choice 7 maps to J, 50% bet.
Choice 8 maps to J, 75% bet.
Choice 9 maps to J, 95% bet.

Your response must always present in the following XML format:
<reasoning>Reasons for your choice this round should cover why you chose the
↪   type and bet percentage</reasoning>
<choice>Any number between 0-9 indicates your choice</choice>
The available range for <choice> is 0-9.
```

---

**Decision-making prompt**

```
Here is the historical information from the past round(s), and you may use it
↪   as a reference for your following choice.
```

```
    In round 1, you chose the <choice_type_{1}> chest and bet
 ↪   <choice_percent_{1}>%. (Un)Fortunately, the coin was hidden under the
 ↪   <coin_type_{1}> chest, and you received <payoff_{1}> points as a reward
 ↪   (or penalty).
    In round 2, you chose the <choice_type_{2}> chest and bet
 ↪   <choice_percent_{2}>%. (Un)Fortunately, the coin was hidden under the
 ↪   <coin_type_{2}> chest, and you received <payoff_{2}> points as a reward
 ↪   (or penalty).


                                      ⋮


    In round i-1, you chose the <choice_type_{i-1}> chest and bet
 ↪   <choice_percent_{i-1}>%. (Un)Fortunately, the coin was hidden under the
 ↪   <coin_type_{i-1}> chest, and you received <payoff_{i-1}> points as a
 ↪   reward (or penalty).

    Your total points in this phase so far: <remain_by_subsession> points. Now
 ↪   this is the i-th round of the game. In front of you are
 ↪   <num_type_RED_boxes> Type F chest(s) and <num_type_BLUE_boxes> Type J
 ↪   chest(s). Please make your choice.
```

## Wisconsin Card Sorting Task

<chest_attribute_{A}> describes the symbols on Card A (e.g., "Card A has 2 green flowers."). <item_attribute_{i}> describes the characteristics of the item in the $i$-th round (e.g., "The item has 2 blue hearts."). <choice_id_{i}> indicates the card selected in the $i$-th round (A, B, C, or D). <reason_{i}> explains the reasoning process for the $i$-th round, which aims to improve the LLM's performance on this task.

---

**System prompt**

```
    In the game, you have 4 chests in front of you.
    In each round, you will be presented with one item, and your task is to choose
 ↪   one of the 4 chests to match the presented item based on its pattern.
    The pattern will be one of the following three: color, shape, or number. There
 ↪   will be no combination of these patterns to define the match.
    If the match is correct, you will receive a ``Match Correct''; if incorrect,
 ↪   you will get a ``Match Failed.''
    Note: You must determine whether to match based on color, number, or shape.
 ↪   Once you figure out the rule, you can follow it for a while, but stay
 ↪   alert|the rule changes periodically! Pay close attention to feedback; if
 ↪   you receive error messages, it's time to adjust your rule. That's all!

    <chest_attribute_{A}>.

    <chest_attribute_{B}>.

    <chest_attribute_{C}>.

    <chest_attribute_{D}>.

    Your response must always present in the following format:
    <reasoning>A brief reason for your choice this round</reasoning>
    <choice>Any number between 1-4 indicates your choice of chest A, B, C,
 ↪   D</choice>
```

---

```
    Decision-making prompt

        Here is the historical information from the past round(s), and you may use it
        ↪  as a reference for your following choice.

        In round 1, <item_attribute_{1}>, You chose chest <choice_id_{1}>. Your
        ↪  reasoning process is <reason_{1}>. Match Correct (or Failed).

        In round 2, <item_attribute_{2}>, You chose chest <choice_id_{2}>. Your
        ↪  reasoning process is <reason_{2}>. Match Correct (or Failed).




                                           .
                                           .
                                           .



        In round {i-1}, <item_attribute_{i-1}>, You chose chest <choice_id_{i-1}>.
        ↪  Your reasoning process is <reason_{i-1}>. Match Correct (or Failed).

        Now this is the i-th round of the game. <item_attribute_{i}>. Please make your
        ↪  choice.
```

## Supplementary note 2: Computational models

We used computational models to quantitatively understand the decision-making processes of humans and LLMs. We utilized hierarchical Bayesian modeling to estimate the parameters of these models, employing the Stan software package [68] for posterior inference. The implementation refers to the hBayesDM package [69]. The descriptions of each model are as follows.

### Iowa Gambling Task

We used the Prospect Valence Learning Model with decay Reinforcement Learning rule [53] in the Iowa Gambling Task, which includes four parameters: learning rate ($A$), choice consistency ($c$), outcome sensitivity ($\alpha$), and loss aversion ($\lambda$). This model employs the prospect theory utility function to evaluate outcomes, which is characterized by a diminishing sensitivity to both gains and losses (parameterized by $\alpha$) and an asymmetric weighting of losses relative to gains (parameterized by $\lambda$). The utility function $u(t)$ for outcome $x(t)$ at the round $t$ is defined as

$$u(t) = \begin{cases} x(t)^{\alpha} & \text{if } x(t) \geq 0, \\ -\lambda|x(t)|^{\alpha} & \text{if } x(t) < 0. \end{cases}$$

Let $E_j(t)$ represent the expected value of the $j$-th deck at the round $t$, and $A$ represent the learning rate governing the forgetting of historical information. The indicator $\delta_j(t)$ denotes whether the $j$-th deck was selected at round $t$, thereby influencing the update of its expectancy

$$E_j(t) = A \cdot E_j(t-1) + \delta_j(t) \cdot u(t).$$

The model utilizes the softmax rule to select a deck for the round $t+1$, i.e.

$$\Pr(j, t) = \frac{e^{\theta \cdot E_j(t)}}{\sum_{k=1}^{4} e^{\theta \cdot E_k(t)}},$$

where $\theta$ is conventionally set to $3^c - 1$ [53], with $c$ being the choice consistency parameter.

### Cambridge Gambling Task

We used the Cumulative Model [54] in the Cambridge Gambling Task, which involves four parameters: type bias for RED ($c$), probability distortion ($\alpha$), risk aversion ($\rho$), and choice consistency ($\gamma$). The probabilities

of selecting red or blue boxes in this model are defined as follows:

$$\Pr(Red) = \frac{cr^\alpha}{cr^\alpha + (1-c)(1-r)^\alpha}, \quad \Pr(Blue) = 1 - \Pr(Red),$$

where $r$ represents the proportion of red boxes, $c$ represents the type bias for selecting red boxes, and $\alpha$ is the probability distortion parameter, indicating whether the participant distorts the asymmetry in probabilities. After selecting a specific box type, for each betting proportion $b_i \in \{5\%, 25\%, 50\%, 75\%, 95\%\}$, the expected value is given as follows:

$$E(b_i \mid \text{Chosen type}) = \Pr(\text{Chosen type}) \cdot u(\text{Wins})$$
$$+ \big(1 - \Pr(\text{Chosen type})\big) \cdot u(\text{Loss}),$$

where $u(\text{Wins})$ and $u(\text{Loss})$ correspond to the utility functions under the winning and losing outcomes, respectively. These two utility functions are defined as follows:

$$u(\text{Wins}) = \log\big(1 + \text{CurrScore} \cdot (1 + b_i)\big)$$

$$u(\text{Loss}) = \log\big(1 + \rho \cdot \text{CurrScore} \cdot (1 - b_i)\big),$$

where CurrScore represents the current score in the current sub-session, and $\rho$ is the risk aversion parameter. The softmax rule is applied to determine the betting proportion, i.e.

$$\Pr\big(b_i \mid \text{Chosen Type}\big) = \frac{\exp\big\{\gamma \cdot E\big(b_i \mid \text{Chosen Type}\big)\big\}}{\sum_j \exp\big\{\gamma \cdot E\big(b_j \mid \text{Chosen Type}\big)\big\}}.$$

Note that we reworded the experiments to mitigate the memorization effect. After rewording, types F and J correspond to red and blue boxes, respectively.

### Wisconsin Card Sorting Task

We applied the Sequential Learning Model [56] in the Wisconsin Card Sorting Task, which incorporates three parameters: reward sensitivity ($r$), punishment sensitivity ($p$), and choice consistency ($d$). For each presented card $k$, let $\mathbf{m}_k(t)$ be a $3 \times 1$ vector, in which each generic element $m_{k,i}(t) = 0$ or 1 denotes whether the card $k$ matches the item at round $t$ according to a specific rule $i \in \{\text{color}, \text{shape}, \text{number}\}$. Let $\mathbf{a}(t)$ be a $3 \times 1$ attention vector representing the attentional weights assigned to each rule at the round $t$. The probability of selecting card $k$ at the round $t$ is given by $\Pr(k, t)$:

$$\Pr(k, t) = \frac{\mathbf{m}_k(t)^\top \mathbf{a}(t)^d}{\sum_{j=1}^4 \mathbf{m}_j(t)^\top \mathbf{a}(t)^d},$$

where $d$ represents the choice consistency parameter, and $\mathbf{a}(t)^d$ denotes the element-wise exponentiation of the vector $\mathbf{a}(t)$. After selecting card $k$, the attention vector $\mathbf{a}(t+1)$ for the round $t+1$ is updated based on the feedback signal, such that

$$\mathbf{a}(t+1) = \begin{cases} (1-r)\mathbf{a}(t) + r\mathbf{s}(t) & \text{if 'correct'} \\ (1-p)\mathbf{a}(t) + p\mathbf{s}(t) & \text{if 'incorrect'} \end{cases}$$

where $r$ and $p$ are the reward sensitivity and punishment sensitivity parameters, respectively, and $\mathbf{s}(t)$ is a $3 \times 1$ vector denoting the signal amplitude. For each generic element $s_i(t)$, its value is given by

$$s_i(t) = \begin{cases} \frac{m_{k,i}(t)a_i(t)}{\mathbf{m}_k(t)^\top \mathbf{a}(t)} & \text{if 'correct'}, \\ \frac{(1-m_{k,i}(t))a_i(t)}{(\mathbf{1}^\top - \mathbf{m}_k(t)^\top)\mathbf{a}(t)} & \text{if 'incorrect'}. \end{cases}$$

## Supplementary note 3: Robustness checks

In total, we examined 19 variants in both the Iowa Gambling Task and Cambridge Gambling Task, and 15 variants in the Wisconsin Card Sorting Task. The variants included: (i) adjusting the LLMs' temperature from the default values to 0 and 0.5; (ii) applying arithmetic transformations to the scores of each round, such as multiplying by a constant factor or adding a constant value (specifically, 0.5x, 2x, +100, +200 for

Iowa Gambling Task; 5x, 10x, +100, +200 for Cambridge Gambling Task; and not applicable to Wisconsin Card Sorting Task, as it does not involve scores); (iii) restructuring the prompts to reflect different decision-making contexts, including both economic and medical scenarios; and (iv) introducing role-play prompts to assess the impact of demographic information (e.g., different age ranges, ethnicity, gender, and other categories such as elder, kid, American, Asian, Black, Hispanic, and White) and risk preferences (e.g., risk-taking and risk-averse behavior).

Each variant was repeated 10 times using GPT-4o. Overall, we observed that the decision-making patterns of the LLMs remained qualitatively the same across these variants, confirming the robustness of our findings and highlighting the fundamental differences in decision-making patterns between LLMs and humans. The following summarizes the key observations in these variants.

## Iowa Gambling Task

As shown in *SI Appendix*, Fig. S5, GPT-4o consistently reproduced the card choice pattern observed in the baseline condition reported in the main text, favoring deck D over deck C across all 19 prompt variants. This pattern demonstrates that the model's sensitivity to penalty frequency remained robust across varied prompt designs. These results highlight GPT-4o's stable decision-making strategy under uncertain conditions.

## Cambridge Gambling Task

As shown in *SI Appendix*, Fig. S6, the GPT-4o consistently selected the majority box type across all prompt variants, accurately identifying the most probable outcome regardless of asymmetry in box distributions. This pattern held across all variants except for economics and medicine, indicating a broadly robust capacity for probabilistic reasoning.

In contrast, *SI Appendix*, Fig. S7 reveals that the GPT-4o exhibited minimal adjustment in its betting behavior across varying levels of risk. Despite increasing asymmetry (e.g., from 6:4 to 9:1 ratios), bet proportions remained largely stable across conditions. Such a lack of sensitivity to shifting odds was consistent across all experimental conditions, which suggests a rigid decision strategy that does not adapt to changing risk levels.

## Wisconsin Card Sorting Task

As shown in *SI Appendix*, Fig. S8, the LLM consistently produced fewer Non-perseverative errors than Perseverative errors across all the variants. The relative reduction in Non-perseverative errors, typically regarded as random or inattentive mistakes, suggests that the LLM maintained strong task focus and demonstrated enhanced rule-learning efficiency and overall performance.

**Table S1**: Statistical comparisons and parameter estimates for the Iowa Gambling Task

(a) Comparisons of overall performance between LLMs and humans

| Model | $W$ | $p$-value | Cohen's $d$ |
|---|---|---|---|
| GPT | 2 556 | $< .001^{***}$ | $-1.135$ |
| Claude | 773.5 | $< .001^{***}$ | $-2.456$ |
| Gemini | 1 209 | $< .001^{***}$ | $-1.973$ |
| GPTo4m | 815.5 | $< .001^{***}$ | $-2.438$ |
| DeepSeek | 773 | $< .001^{***}$ | $-2.497$ |

*Note.* Cohen's $d$ represents the standardized difference in net scores between human participants and each LLM, with LLMs serving as the reference group. Negative values indicate that LLMs outperformed humans on the Iowa Gambling Task.

(b) Within-group comparisons of deck preferences (Deck C vs. Deck D)

| Model | # of Deck C | # of Deck D | % of Deck C$^*$ | % of Deck D$^*$ | $Z$ | $p$-value | Cohen's $h$ |
|---|---|---|---|---|---|---|---|
| Humans | 388 | 353 | 32.333 | 29.417 | 0.5 | 0.625 | 0.063 |
| GPT | 322 | 531 | 26.833 | 44.250 | $-2.8$ | $0.005^{**}$ | $-0.366$ |
| Claude | 111 | 1055 | 9.250 | 87.917 | $-12.2$ | $< .001^{***}$ | $-1.813$ |
| Gemini | 479 | 679 | 39.917 | 56.583 | $-2.6$ | $0.010^{**}$ | $-0.335$ |
| GPTo4m | 963 | 184 | 80.250 | 15.333 | 10.1 | $< .001^{***}$ | 1.416 |
| DeepSeek | 695 | 472 | 57.917 | 39.333 | 2.9 | $0.004^{**}$ | 0.374 |

*Note.* Cohen's $h$ reflects the standardized difference in proportions between deck C and deck D selections, computed as a two-proportion $Z$ test with deck D as the reference $(C - D)$. Positive values indicate a greater preference for deck C, while negative values indicate a greater preference for deck D. *The percentage of deck C and D is calculated from the total number of deck A, B, C and D (e.g., % of Deck C = $\frac{\text{\# of Deck C}}{\text{\# of Deck A} + \text{\# of Deck B} + \text{\# of Deck C} + \text{\# of Deck D}}$).

(c) Prospect Valence Learning model with a decay reinforcement learning rule: Pairwise parameter comparisons with human participants

| Param. | Model | $W$ | $p$-value | Cohen's $d$ |
|---|---|---|---|---|
| $A$ | GPT | 1 680.0 | $< .001^{***}$ | $-1.696$ |
| | Claude | 360.0 | $< .001^{***}$ | $-2.642$ |
| | Gemini | 38.0 | $< .001^{***}$ | $-3.113$ |
| | GPTo4m | 963.0 | $< .001^{***}$ | $-2.226$ |
| | DeepSeek | 0.0 | $< .001^{***}$ | $-3.045$ |
| $c$ | GPT | 2 893.0 | $< .001^{***}$ | $-1.283$ |
| | Claude | 2 937.0 | $< .001^{***}$ | $-1.284$ |
| | Gemini | 5 545.0 | $0.002^{**}$ | 0.004 |
| | GPTo4m | 1 141.0 | $< .001^{***}$ | $-2.288$ |
| | DeepSeek | 4 107.0 | $< .001^{***}$ | $-0.511$ |
| $\alpha$ | GPT | 0.0 | $< .001^{***}$ | $-7.331$ |
| | Claude | 0.0 | $< .001^{***}$ | $-7.167$ |
| | Gemini | 0.0 | $< .001^{***}$ | $-6.603$ |
| | GPTo4m | 0.0 | $< .001^{***}$ | $-5.849$ |
| | DeepSeek | 0.0 | $< .001^{***}$ | $-6.058$ |
| $\lambda$ | GPT | 6 720.0 | 0.373 | 0.474 |
| | Claude | 4 674.0 | $< .001^{***}$ | $-0.086$ |
| | Gemini | 4 854.0 | $< .001^{***}$ | $-0.104$ |
| | GPTo4m | 4 063.0 | $< .001^{***}$ | $-0.483$ |
| | DeepSeek | 3 480.0 | $< .001^{***}$ | $-0.952$ |

*Note.* Cohen's $d$ represents the standardized difference in parameter estimates between human participants and each LLM, with LLMs serving as the reference group. Negative values indicate that LLMs had higher parameter values than humans.

(d) Prospect Valence Learning model with a decay reinforcement learning rule: Descriptive statistics for parameters

| Param. | Model | Mean | Median | SD |
|---|---|---|---|---|
| $A$ | Humans | 0.7059 | 0.7028 | 0.1263 |
| | GPT | 0.8575 | 0.8579 | 0.0015 |
| | Claude | 0.942 | 0.9421 | 0.0008 |
| | Gemini | 0.9869 | 0.9942 | 0.0179 |
| | GPTo4m | 0.9048 | 0.9049 | 0.0014 |
| | DeepSeek | 0.978 | 0.9782 | 0.0012 |
| $c$ | Humans | 0.414 | 0.3484 | 0.2977 |
| | GPT | 0.6871 | 0.6893 | 0.0455 |
| | Claude | 0.6861 | 0.686 | 0.0343 |
| | Gemini | 0.4132 | 0.4125 | 0.0255 |
| | GPTo4m | 0.8991 | 0.9022 | 0.0355 |
| | DeepSeek | 0.5216 | 0.5214 | 0.0038 |
| $\alpha$ | Humans | 0.3185 | 0.3149 | 0.1219 |
| | GPT | 0.9522 | 0.9543 | 0.0089 |
| | Claude | 0.9366 | 0.937 | 0.0036 |
| | Gemini | 0.9479 | 0.9448 | 0.0575 |
| | GPTo4m | 0.8227 | 0.8227 | 0.0013 |
| | DeepSeek | 0.8407 | 0.8408 | 0.0007 |
| $\lambda$ | Humans | 1.1005 | 0.5903 | 1.1493 |
| | GPT | 0.7151 | 0.7161 | 0.0092 |
| | Claude | 1.1705 | 1.1708 | 0.0035 |
| | Gemini | 1.1865 | 1.212 | 0.2196 |
| | GPTo4m | 1.4933 | 1.4932 | 0.0037 |
| | DeepSeek | 1.874 | 1.8741 | 0.0015 |

**Table S2**: Statistical comparisons and parameter estimates for the Cambridge Gambling Task

(a) Comparisons of overall performance (Total Scores) between LLMs and humans

| Model | $W$ | $p$-value | Cohen's $d$ |
|---|---|---|---|
| GPT | 6 354 | 0.116 | −0.100 |
| Claude | 4 175 | < .001*** | −0.691 |
| Gemini | 9 282.5 | < .001*** | 0.484 |
| GPTo4m | 4 774 | < .001*** | −1.155 |
| DeepSeek | 5 029 | < .001*** | −1.043 |

*Note.* Cohen's $d$ represents the standardized difference in total scores between human participants and each LLM, with LLMs serving as the reference group. Negative values indicate that LLMs outperformed humans on the Cambridge Gambling Task.

(b) Cumulative model: Pairwise parameter comparisons with humans

| Param. | Model | $W$ | $p$-value | Cohen's $d$ |
|---|---|---|---|---|
| | GPT | 1 071.0 | < .001*** | −1.786 |
| | Claude | 8 251.0 | 0.051 | 0.174 |
| $c$ | Gemini | 121.0 | < .001*** | −4.207 |
| | GPTo4m | 4 811.0 | < .001*** | −0.108 |
| | DeepSeek | 12 900.0 | < .001*** | 1.564 |
| | GPT | 0.0 | < .001*** | −2.068 |
| | Claude | 389.0 | < .001*** | −1.710 |
| $\alpha$ | Gemini | 1 864.0 | < .001*** | −1.156 |
| | GPTo4m | 0.0 | < .001*** | −2.078 |
| | DeepSeek | 410.0 | < .001*** | −1.771 |
| | GPT | 2 940.0 | < .001*** | −0.885 |
| | Claude | 12 992.0 | < .001*** | 1.826 |
| $\rho$ | Gemini | 0.0 | < .001*** | −2.610 |
| | GPTo4m | 14 400.0 | < .001*** | 4.346 |
| | DeepSeek | 14 400.0 | < .001*** | 5.466 |
| | GPT | 12 697.0 | < .001*** | 1.299 |
| | Claude | 12 236.0 | < .001*** | 1.174 |
| $\gamma$ | Gemini | 11 560.0 | < .001*** | 1.088 |
| | GPTo4m | 564.0 | < .001*** | −2.201 |
| | DeepSeek | 218.0 | < .001*** | −1.068 |

*Note.* Cohen's $d$ represents the standardized difference in parameter estimates between human participants and each LLM, with LLMs serving as the reference group. Negative values indicate that LLMs had higher parameter values than humans.

(c) Cumulative model: Descriptive statistics for parameters

| Param. | Model | Mean | Median | SD |
|---|---|---|---|---|
| | Humans | 0.4955 | 0.4961 | 0.0301 |
| | GPT | 0.5338 | 0.5331 | 0.0035 |
| $c$ | Claude | 0.4908 | 0.4918 | 0.0225 |
| | Gemini | 0.5866 | 0.5842 | 0.0057 |
| | GPTo4m | 0.4978 | 0.4978 | 0.0002 |
| | DeepSeek | 0.4617 | 0.4627 | 0.0051 |
| | Humans | 2.9341 | 3.1058 | 1.4003 |
| | GPT | 4.9822 | 4.9822 | 0.0002 |
| $\alpha$ | Claude | 4.7595 | 4.8839 | 0.5637 |
| | Gemini | 4.3296 | 4.8154 | 0.9767 |
| | GPTo4m | 4.9915 | 4.9915 | 0.0001 |
| | DeepSeek | 4.8124 | 4.9277 | 0.5382 |
| | Humans | 63.8021 | 6.4272 | 88.4852 |
| | GPT | 418.6155 | 512.7302 | 306.6256 |
| $\rho$ | Claude | 0.7864 | 0.3560 | 1.0657 |
| | Gemini | 691.8412 | 692.5117 | 5.4954 |
| | GPTo4m | 0.0104 | 0.0088 | 0.0056 |
| | DeepSeek | 0.0018 | 0.0018 | 0.0000 |
| | Humans | 15.8826 | 12.5300 | 11.2113 |
| | GPT | 5.5439 | 5.4863 | 1.0043 |
| $\gamma$ | Claude | 6.5740 | 6.5812 | 0.0470 |
| | Gemini | 7.1626 | 7.1715 | 1.6764 |
| | GPTo4m | 48.7207 | 44.7819 | 17.8695 |
| | DeepSeek | 125.3505 | 65.0248 | 144.5221 |

**Table S3**: Statistical comparisons and parameter estimates for the Wisconsin Card Sorting Task

(a) Comparisons of overall performance between LLMs and humans

| Model | $W$ | $p$-value | Cohen's $d$ |
|---|---|---|---|
| GPT | 4 910.5 | $< .001^{***}$ | $-0.399$ |
| Claude | 3 579.0 | $< .001^{***}$ | $-0.925$ |
| Gemini | 3 109.5 | $< .001^{***}$ | $-0.934$ |
| GPTo4m | 2 854.0 | $< .001^{***}$ | $-1.218$ |
| DeepSeek | 7 391.5 | 0.722 | 0.332 |

*Note.* Cohen's $d$ represents the standardized difference in overall performance between human participants and each LLM, with LLMs serving as the reference group. Negative values indicate that LLMs outperformed humans on the Wisconsin Card Sorting Task.

(b) Within-group comparisons of error patterns (perseverative vs. non-perseverative)

| Model | $W$ | $p$-value | Cohen's $d$ |
|---|---|---|---|
| Humans | 4 821.0 | $< .001^{***}$ | $-0.868$ |
| GPT | 9 266.5 | $< .001^{***}$ | $-0.030$ |
| Claude | 11 079.5 | $< .001^{***}$ | 0.611 |
| Gemini | 12 124.5 | $< .001^{***}$ | 0.480 |
| GPTo4m | 12 636.5 | $< .001^{***}$ | 1.195 |
| DeepSeek | 5 040.5 | $< .001^{***}$ | $-1.013$ |

*Note.* Cohen's $d$ reflects the standardized difference between perseverative errors and non-perseverative errors. Negative values indicate more non-perseverative errors relative to perseverative errors.

(c) Sequential learning model: Pairwise parameter comparisons with humans

| Param. | Model | $W$ | $p$-value | Cohen's $d$ |
|---|---|---|---|---|
| | GPT | 0.0 | $< .001^{***}$ | $-18.875$ |
| | Claude | 0.0 | $< .001^{***}$ | $-17.957$ |
| $r$ | Gemini | 0.0 | $< .001^{***}$ | $-16.902$ |
| | GPTo4m | 0.0 | $< .001^{***}$ | $-18.696$ |
| | DeepSeek | 14 400.0 | $< .001^{***}$ | 16.079 |
| | GPT | 4 333.0 | $< .001^{***}$ | $-0.263$ |
| | Claude | 2 859.0 | $< .001^{***}$ | $-0.294$ |
| $p$ | Gemini | 0.0 | $< .001^{***}$ | $-1.238$ |
| | GPTo4m | 0.0 | $< .001^{***}$ | $-1.209$ |
| | DeepSeek | 0.0 | $< .001^{***}$ | $-1.241$ |
| | GPT | 5 132.0 | $< .001^{***}$ | $-0.470$ |
| | Claude | 2 973.0 | $< .001^{***}$ | $-1.171$ |
| $d$ | Gemini | 2 963.0 | $< .001^{***}$ | $-1.080$ |
| | GPTo4m | 2 219.0 | $< .001^{***}$ | $-1.433$ |
| | DeepSeek | 7 602.0 | 0.455 | $-0.253$ |

*Note.* Cohen's $d$ represents the standardized difference in parameter estimates between human participants and each LLM, with LLMs serving as the reference group. Negative values indicate that LLMs had higher parameter values than humans.

(d) Sequential learning model: Descriptive statistics for parameters

| Param. | Model | Mean | Median | SD |
|---|---|---|---|---|
| | Humans | 0.9974 | 0.9975 | 0.0002 |
| | GPT | 0.9997 | 0.9997 | 0.0000 |
| $r$ | Claude | 0.9996 | 0.9996 | 0.0000 |
| | Gemini | 0.9995 | 0.9995 | 0.0000 |
| | GPTo4m | 0.9997 | 0.9997 | 0.0000 |
| | DeepSeek | 0.9729 | 0.9727 | 0.0022 |
| | Humans | 0.9238 | 0.9548 | 0.0867 |
| | GPT | 0.9486 | 0.9772 | 0.1010 |
| $p$ | Claude | 0.9531 | 0.9895 | 0.1108 |
| | Gemini | 0.9997 | 0.9997 | 0.0000 |
| | GPTo4m | 0.9979 | 0.9980 | 0.0002 |
| | DeepSeek | 0.9999 | 0.9999 | 0.0000 |
| | Humans | 0.3626 | 0.3563 | 0.1454 |
| | GPT | 0.4308 | 0.4519 | 0.1449 |
| $d$ | Claude | 0.5496 | 0.5436 | 0.1728 |
| | Gemini | 0.5098 | 0.5149 | 0.1266 |
| | GPTo4m | 0.5507 | 0.5513 | 0.1153 |
| | DeepSeek | 0.4392 | 0.2924 | 0.4037 |

**Table S4**: Summary of participant demographics by task and session. IGT denotes Iowa Gambling Task, CGT denotes Cambridge Gambling Task, and WCST denotes Wisconsin Card Sorting Task.

| Date | Task | Location | Rounds | Participants | Mean age | SD age | %Women |
|------|------|----------|--------|--------------|----------|--------|--------|
| 14 Jul 2024 | IGT | Xi'an | 80 | 40 | 19.6 | 1.6 | 42.5 |
| 20 Jul 2024 | IGT | Xi'an | 80 | 40 | 20.7 | 1.5 | 50.0 |
| 20 Jul 2024 | IGT | Xi'an | 80 | 40 | 20.1 | 1.5 | 55.0 |
| 7 Sept 2024 | CGT | Xi'an | 64 | 40 | 18.4 | 1.3 | 35.0 |
| 7 Sept 2024 | CGT | Xi'an | 64 | 40 | 18.7 | 2.0 | 32.5 |
| 7 Sept 2024 | CGT | Xi'an | 64 | 40 | 18.2 | 1.9 | 27.5 |
| 7 Dec 2024 | WCST | Xi'an | 64 | 30 | 18.3 | 0.5 | 36.7 |
| 7 Dec 2024 | WCST | Xi'an | 64 | 30 | 18.5 | 0.8 | 23.3 |
| 8 Dec 2024 | WCST | Xi'an | 64 | 30 | 18.1 | 0.7 | 46.7 |
| 8 Dec 2024 | WCST | Xi'an | 64 | 30 | 18.3 | 0.7 | 23.3 |

**Table S5**: The parameters in computational models.

| Model | Parameter | Range of Values | Interpretation |
|-------|-----------|-----------------|----------------|
| PVL-DecayRI model [53] | $A$ | $0 \leq A \leq 1$ | Learning rate |
| | $c$ | $0 \leq c \leq 5$ | Choice consistency |
| | $\alpha$ | $0 \leq \alpha \leq 2$ | Outcome sensitivity |
| | $\lambda$ | $0 \leq \lambda \leq 10$ | Loss aversion |
| Cumulative model [54] | $c$ | $0 \leq c \leq 1$ | Type bias for RED |
| | $\alpha$ | $0 \leq \alpha \leq 5$ | Probability distortion |
| | $\rho$ | $0 \leq \rho < +\infty$ | Risk aversion |
| | $\gamma$ | $0 \leq \gamma < +\infty$ | Choice consistency |
| Sequential learning model [56] | $r$ | $0 \leq r \leq 1$ | Reward sensitivity |
| | $p$ | $0 \leq p \leq 1$ | Punishment sensitivity |
| | $d$ | $0 \leq d \leq 5$ | Choice consistency |

**Fig. S1**: **LLMs learned faster than humans in the Iowa Gambling Task, showing steeper increases in advantageous deck selections over time.** The plot shows the proportion of advantageous choices (decks C or D) across 80 rounds. Human participants (slope = 0.374, $r = 0.83$, $p < .001$), GPT (slope = 0.466, $r = 0.648$, $< .001$), Claude (slope = 0.925, $r = 0.69$, $p < .001$), GPTo4m (slope = 0.768, $r = 0.813$, $p < .001$), Gemini (slope = 1.027, $r = 0.812$, $p < .001$), and DeepSeek (slope = 0.877, $r = 0.822$, $p < .001$) all exhibited a significant increasing trend of choosing advantageous decks, with LLMs demonstrating faster learning speeds compared to humans.



**Fig. S2**: **LLMs demonstrated consistently higher decision-making quality than humans across all levels of risk conditions.** This panel shows the decision-making quality in terms of the percentage of selecting the most likely outcome. Human performance improved with increasing asymmetry, but remained below that of all LLMs, which approached ceiling-level accuracy across conditions.

**Fig. S3**: **LLMs identified the correct matching rule more quickly than humans but showed similar stability in maintaining the correct strategy.** TRSET1 in the left panel represents the number of rounds to complete the first set. LLMs required significantly fewer rounds compared to human participants ( GPT: $W = 10{,}612$, $p < .001$, Cohen's $d = 0.976$; Claude: $W = 9{,}828.5$, $p < .001$, Cohen's $d = 0.644$; Gemini: $W = 10{,}809$, $p < .001$, Cohen's $d = 0.982$; GPTo4m: $W = 10{,}756$, $p < .001$, Cohen's $d = 1.018$; nevertheless, DeepSeek ($W = 7{,}627$, $p = 0.423$, Cohen's $d = 0.128$) matched human participants.). In the right panel, FSET measures how often participants changed their matching strategy after achieving five or more consecutive correct matches without negative feedback ('incorrect'). Human participants and LLMs exhibited comparable FSET values ( GPT: $W = 6{,}553.5$, $p = 0.208$, Cohen's $d = -0.189$; Claude: $W = 6{,}906.5$, $p = 0.569$, Cohen's $d = -0.142$; Gemini: $W = 6{,}443$, $p = 0.1414$, Cohen's $d = -0.236$; GPTo4m: $W = 6{,}698$, $p = 0.3276$, Cohen's $d = -0.16$; DeepSeek: $W = 7{,}528$, $p = 0.5241$, Cohen's $d = 0.033$).

**A**

**Attitudes towards AI usage**



**B**

**Let AI assist me in playing**



Fig. S4: **Participants generally exhibit an overall negative attitude toward AI assistance across all tasks. Panel (A)** shows the participants significantly disagreed with statements reflecting positive perceptions of AI, including believing AI can assist their decision-making ( IGT: $V = 546.5$, $p < .001$; CGT: $V = 829.5$, $p < .001$; WCST: $V = 590.0$, $p < .001$), believing AI makes it easier to make decisions (IGT: $V = 575.0$, $p < .001$; CGT: $V = 1,113.0$, $p < .001$; WCST: $V = 780.0$, $p < .001$), believing AI can help them finish faster ( IGT: $V = 624.5$, $p < .001$; CGT: $V = 1,275.0$,$p < .001$; WCST: $V = 801.0$, $p < .001$), believing using AI can improve their game scores ( IGT: $V = 774.0$, $p < .001$; CGT: $V = 1,598.5$, $p = 0.0181$; WCST: $V = 932.0$, $p = 0.0034$), and letting AI assist them in playing ( IGT: $V = 335.5$, $p < .001$; CGT: $V = 808.5$, $p < .001$; WCST: $V = 1,028.0$, $p < .001$). Conversely, participants displayed relatively neutral or slightly positive attitudes toward letting AI play the game for them ( IGT: $V = 1,557.0$, $p = 0.0032$; CGT: $V = 1,563.0$, $p = 0.0012$; WCST: $V = 1,470.0$, $p = 0.0126$). **Panel (B)** shows that participants expressed reluctance to accept AI assistance, regardless of whether GPT outperformed them. After each task, participants rated their agreement with the statement "Let AI assist me in playing." Responses were split based on whether the participant's performance was equal to or better than GPT (blue), or worse than GPT (red). We found that participants maintained negative attitudes toward AI assistance even when informed of the LLMs' high performance on the same tasks and despite their own lower scores ( IGT: $V = 115.0$, $p < .001$; CGT: $V = 317.0$, $p < .001$; WCST: $V = 361.5$, $p < .001$). Statistical significance from neutrality was assessed using the Wilcoxon signed-rank test. Bars represent the mean agreement level, with 95% confidence intervals of the mean shown as error bars.
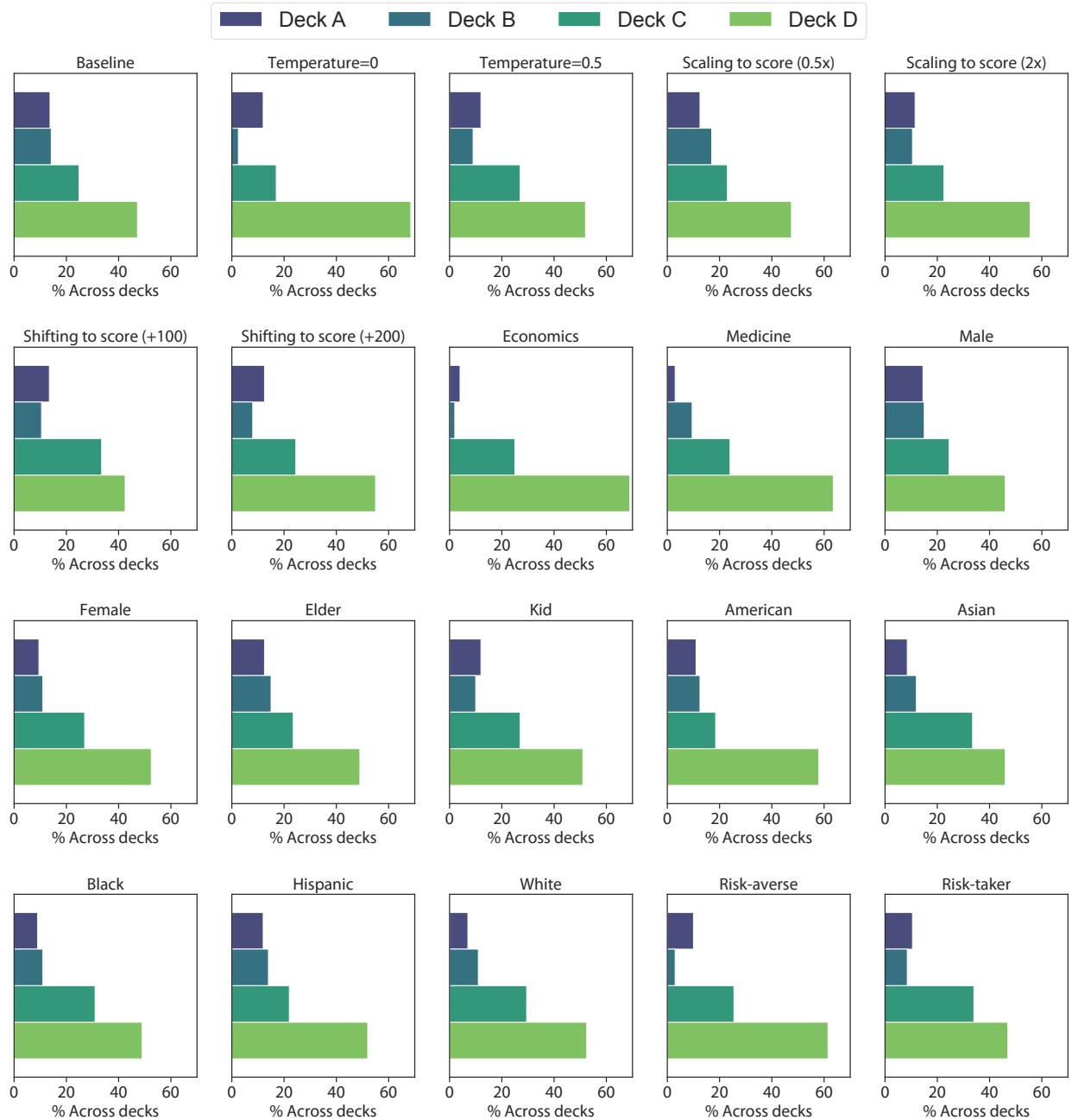
**Fig. S5**: **Robustness checks on Iowa Gambling Task.** Choice Distribution in the Last 20 Rounds for different prompt variations, using GPT-4o over 10 sessions. The LLM predominantly selected the two advantageous decks across all experimental variants, exhibiting a stronger preference for deck D, which imposed penalties less frequently.
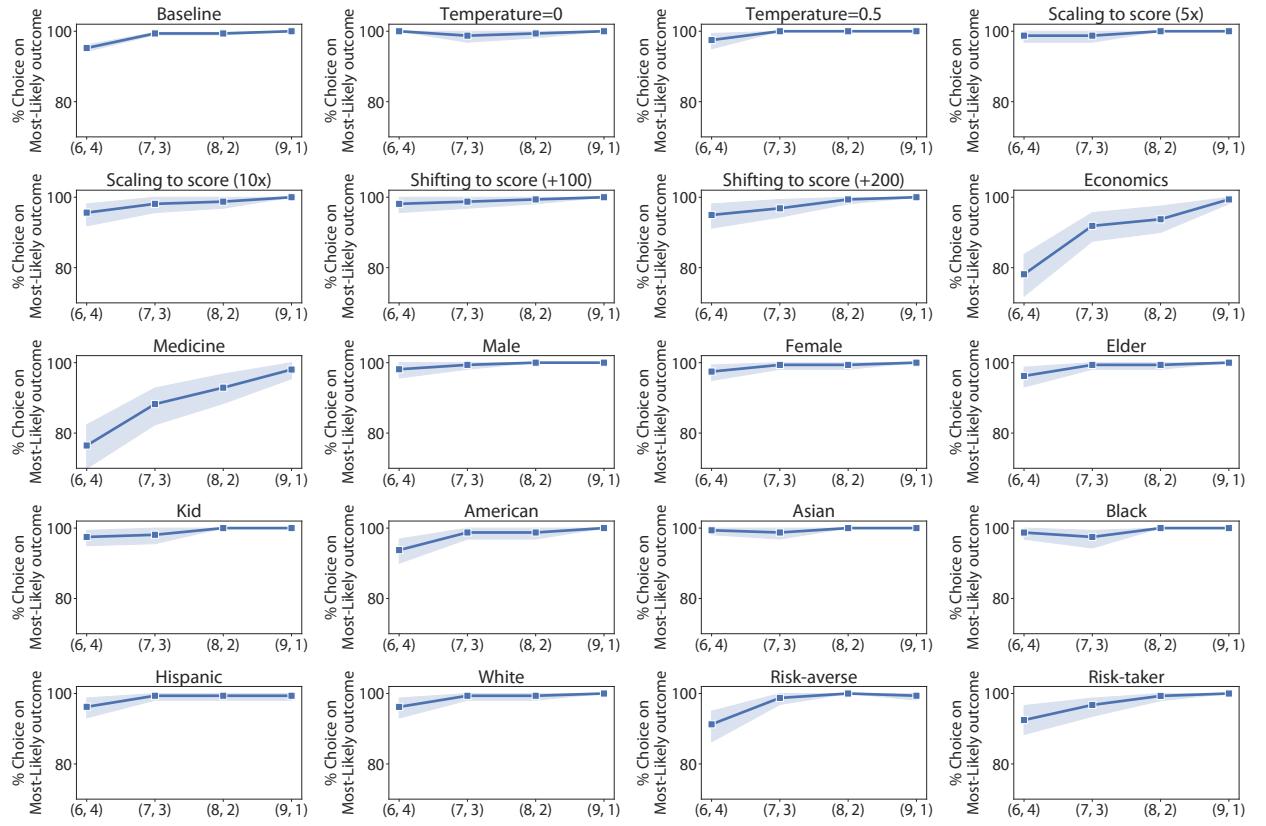
**Fig. S6**: **Robustness checks for the Cambridge Gambling Task.** Decision-making quality for different prompt variations, using GPT-4o over 10 sessions. The LLM consistently opted for the majority box type irrespective of the degree of asymmetry in box distributions, thereby reliably predicting the most probable outcome across the various variants.

**Fig. S7**: **Robustness checks for the Cambridge Gambling Task.** Risk adjustment for different prompt variations, using GPT-4o over 10 sessions. The LLM maintained a consistent betting pattern across different levels of asymmetry, indicating a lack of effective risk adjustment.

**Fig. S8**: **Robustness checks for the Wisconsin Card Sorting Task.** Two WCST metrics: Perseverative errors and Non-perseverative errors for different prompt variations, using GPT-4o over 10 sessions. Median values are represented by diamond markers. The LLM demonstrated fewer Non-perseverative errors than Perseverative errors across the variants, thereby evidencing enhanced task proficiency.

## Introduction

### Welcome to the behavioral experiment

Welcome to the behavioral experiment! All collected data will be used for research purposes only.

You will receive a certain amount of payoff after finishing the experiment, which consists of a show-up fee of 15 RMB and an experiment bonus, typically ranging from 5-25 RMB.

The bonus depends on your final score. A Higher score means a higher payoff. The system will show your score cumulated over time and will show your final score at the end of the experiment.

**Note that** you will receive no payoff if you withdraw midway. Please do not communicate with other participants during the experiment.

You can raise your hand for assistance if needed. Ensure all the electronic devices are on silent or flight mode during the entire experiment. Thank you for your cooperation.

☐ I acknowledge and agree to the provided terms. I voluntarily participate in the Behavioral Experiment.

Next

## Introduction

### Welcome to the behavioral experiment

In this game, you find yourself in a mysterious room with four ancient treasure chests. Each chest holds an unknown thing that can be either a reward or a penalty. With each turn, you will choose one chest to open. Please consider carefully as your choice may significantly impact your points. Specifically, the rewards will increase your points, while penalties will deduct your points. The game has several rounds in which you points will accumulate, and your goal is to maximize your points by the end of the game.

**The only hint I can give you**, and the most important thing to note is this: Out of these chests, there are some that are worse than others, and to win you should try to stay away from bad chests. No matter how much you find yourself losing, you can still win the game if you avoid the worst chests. Also note that the computer does not change the order of the chests once the game begins. It does not make you lose at random, or make you lose money based on the last chest you picked.

☐ I understand the rules and know how to maximize my score.

Next

**Fig. S9**: Pages of introduction on Iowa Gambling Task. Once participants have understood the task, they can begin by clicking the "Next" button.

## Introduction

### Welcome to the behavioral experiment

Welcome to the behavioral experiment! All collected data will be used for research purposes only.

You will receive a certain amount of payoff after finishing the experiment, which consists of a show-up fee of 15 RMB and an experiment bonus, typically ranging from 5-25 RMB.

The bonus depends on your final score. A Higher score means a higher payoff. The system will show your score cumulated over time and will show your final score at the end of the experiment.

**Note that** you will receive no payoff if you withdraw midway. Please do not communicate with other participants during the experiment.

You can raise your hand for assistance if needed. Ensure all the electronic devices are on silent or flight mode during the entire experiment. Thank you for your cooperation.

☐ I acknowledge and agree to the provided terms. I voluntarily participate in the Behavioral Experiment.

Next

## Introduction

### Welcome to the behavioral experiment

In front of you are 10 treasure chests, some of which are type F and others are type J. A gold coin is hidden in one of these chests. Your task is to guess in which **type** of chest the gold coin is hidden. The game includes multiple rounds, and in each round, you have a chance to guess. Guessing will deduct from your existing points as a cost.

Every **8** rounds constitute a phase. At the start of each phase, your points will be set to **100.0**. Your goal is to maximize the points in each phase, and the sum of the points from each phase is your total score in the game.

Each round, you need to bet by deducting a certain percentage from your existing points to choose the type of chest. The available percentages are: 5%, 25%, 50%, 75%, 95%.

If the gold coin is hidden in **any one** of the chests of the type you guessed, you win; otherwise, you lose.

If you win, you will receive double the points of your bet; if you lose, you gain nothing. Points are rounded to the nearest whole number.

**An important game setting is:** the position of the gold coin is completely random in each round and has no connection to previous locations. That is, each round's coin position is an independent event.

☐ I understand the rules and know how to maximize my score.

Next

**Fig. S10**: Pages of introduction on Cambridge Gambling Task. Once participants have understood the task, they can begin by clicking the "Next" button.

## Introduction

### Welcome to the behavioral experiment

Welcome to the behavioral experiment! All collected data will be used for research purposes only.

You will receive a certain amount of payoff after finishing the experiment, which consists of a show-up fee of 15 RMB and an experiment bonus, typically ranging from 5-25 RMB.

The bonus depends on your final score. A Higher score means a higher payoff. The system will show your score cumulated over time and will show your final score at the end of the experiment.

**Note that** you will receive no payoff if you withdraw midway. Please do not communicate with other participants during the experiment.

You can raise your hand for assistance if needed. Ensure all the electronic devices are on silent or flight mode during the entire experiment. Thank you for your cooperation.

☐ I acknowledge and agree to the provided terms. I voluntarily participate in the Behavioral Experiment.

Next

## Introduction

### Welcome to the behavioral experiment

In the game, you have 4 chests in front of you.

In each round, you will be presented with one item, and your task is to choose one of the 4 chests to match the presented item based on its pattern.

The pattern will be one of the following three: color, shape, or number. There will be no combination of these patterns to define the match.

If the match is correct, you will receive a "Match Correct"; if incorrect, you will get a "Match Failed."

**Note:** You must determine whether to match based on color, number, or shape. Once you figure out the rule, you can follow it for a while, but stay alert—the rule changes periodically! Pay close attention to feedback; if you receive error messages, it's time to adjust your rule. That's all!

☐ I understand the rules and know how to maximize my score.

Next

**Fig. S11**: Pages of introduction on Wisconsin Card Sorting Task. Once participants have understood the task, they can begin by clicking the "Next" button.

# Choose Your Chest

Round 1

Please choose one of the following treasure chests.

| Chest1 | Chest2 | Chest3 | Chest4 |

Chest 1 ○ Chest 2 ○ Chest 3 ○ Chest 4 ○

Next

# Results in this round

Round 1

| Chest1 | Chest2 | Chest3 | Chest4 |

This round you earned： **100 points**

Your total points so far： **2100 points**

Next

**Fig. S12**: Choice and Result Pages on Iowa Gambling Task.

# Choose Your Choice

Round 1

Please choose one of the following choices.

In front of you, there are 8 type F chests and 2 type J chests.

| Type F | Type F | Type F | Type F | Type F | Type F | Type F | Type F | Type J | Type J |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|

Your current phase points are 100. Guess under which type of chest the gold coin is and decide on your betting ratio:

**Type F Chests**

| F 5%: ● | F 25%: ○ | F 50%: ○ | F 75%: ○ | F 95%: ○ |
|---------|----------|----------|----------|----------|

**Type J Chests**

| J 5%: ○ | J 25%: ○ | J 50%: ○ | J 75%: ○ | J 95%: ○ |
|---------|----------|----------|----------|----------|

[Next]

# Results in this round

Round 1

You chose the type F chest, and your bet ratio was 5%, with the gold coin under chest number 2:

| Type F | Type F | Type F | Type F | Type F | Type F | Type F | Type F | Type J | Type J |
|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|

Fortunately, you guessed right, and in this round you received: **5 points**

Total points earned in this phase: **105 points**

[Next]

**Fig. S13**: Choice and Result Pages on Cambridge Gambling Task.

36

**Fig. S14**: Choice Page on Wisconsin Card Sorting Task.

**Fig. S15**: Result Page on Wisconsin Card Sorting Task.



**Fig. S16**: Final Result Page will display the participant's score for the current task.

# Information Collection

Your accumulated score in this experiment is **380**.

We will record your information for the purpose of distributing the experiment compensation.
(Your information will only be used for this experiment and will not be disclosed to any third parties.)

Please ensure that the Name, Phone Number, and Student Number you provide are correct; otherwise, you will not receive the experiment payoff.

Experiment ID:

Name:

Phone Number:

Student Number:

University：

Major:

Ethnicity:

Gender:

.......    ⌄

Age:

Country:

Religion:

Alipay Account:

Submit

**Fig. S17**: Demographic Information Collection Page at the end of all tasks.

## Survey Items

**During the game:**

I need to frequently make complex decisions

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

I need to evaluate, compare, and weigh the available information to make decisions

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

I often need to engage in extensive thought before making decisions

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

I am capable of discerning better options and making good decisions

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

**This game:**

Involving ambiguous or uncertain information

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

Placing significant time pressure on me, making it hard for me to decide in time

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

**I think:**

Using AI can improve my game score

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

Using AI can help me finish the game faster

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

AI can assist my decision-making during the game

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

AI can make it easier for me to make decisions in the game

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

**The current performance of the AI per round is (X±Y), and your average score per round is Z; if there were another game, I would be willing to:**

Let AI play the game for me

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

Let AI assist me in playing the game

○ Strongly Agree    ○ Agree    ○ Not Sure    ○ Disagree    ○ Strongly Disagree

下一页

**Fig. S18**: Post-Task Survey on task feedback and attitudes toward AI assistance.
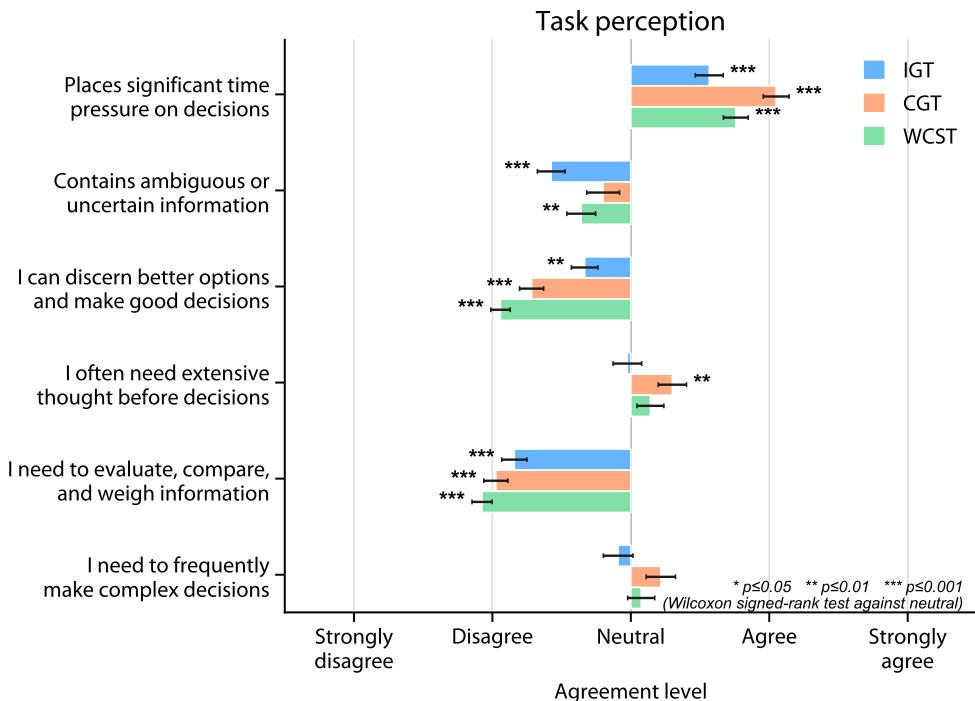
**Fig. S19**: Participants' subjective perceptions of the three decision-making tasks. Bars represent the mean agreement level, with 95% confidence intervals of the mean shown as error bars.
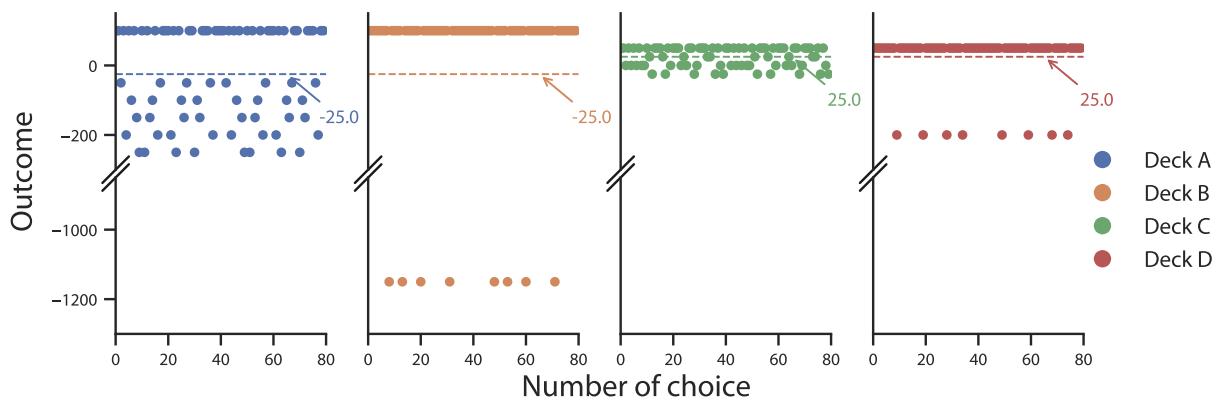


**Fig. S20**: The outcomes for four decks (A, B, C, D) based on the number of choices, with dashed lines representing the average outcome for each deck. Decks A and B provide higher immediate rewards but result in a net loss for every 10 choices due to penalties. Deck A offers five small penalties, while Deck B offers one large penalty. In contrast, Decks C and D provide lower immediate rewards but yield a net gain for every 10 choices. Deck C offers five small penalties and Deck D offers one large penalty.