

ADVERSARIAL TESTING IN LLMs: INSIGHTS INTO DECISION-MAKING VULNERABILITIES

Anonymous authors

Paper under double-blind review

ABSTRACT

As AI systems, particularly Large Language Models (LLMs), rapidly advance towards surpassing human cognitive capabilities, ensuring their alignment with human values and safety standards emerges as a formidable challenge. This study addresses a crucial aspect of superalignment by investigating the decision-making capabilities and adversarial vulnerabilities of LLMs, focusing on GPT-3.5, GPT-4 and Gemini-1.5, within structured experimental settings that mimic complex human interactions. We applied an adversarial framework to two decision-making tasks—the two-armed bandit task and the Multi-Round Trust Task (MRTT)—to test the vulnerabilities of LLMs under adversarial conditions. In the bandit task, the adversary aimed to induce the LLM’s preference for the predefined target action with the constraint that each action must be assigned an equal number of rewards. For the MRTT, we trained two types of adversaries: one aimed at maximizing its own earnings (MAX) and the other focused on maximizing fairness (FAIR). GPT-4 and Gemini-1.5 showed a bias toward exploitation in the bandit task, prioritizing early-established strategies, which made them predictable and vulnerable to manipulation. GPT-3.5, while more exploratory in the bandit task, demonstrated more risk-seeking behavior in the MRTT, leading to increased vulnerability in interacting with the MAX adversary. Notably, Gemini-1.5 excelled in the MRTT, adapting effectively to adversaries and outperforming both GPT-3.5 and GPT-4 by balancing risk and cooperation with its adversaries. By presenting a specific set of tasks that characterizes decision-making vulnerabilities in LLM-based agents, we provide a concrete methodology for evaluating their readiness for real-world deployment. The adversarial framework proved a powerful tool for stress-testing LLMs, revealing the importance of ensuring that AI models are both robust against adversarial manipulation and responsive to fairness cues in complex, dynamic environments.

1 INTRODUCTION

In recent years, the landscape of Artificial Intelligence (AI) has been reshaped by the emergence of Large Language Models (LLMs) such as Generative Pre-trained Transformers (GPT). These models have demonstrated remarkable performance across various applications, from natural language processing to complex problem-solving tasks. As LLMs are increasingly integrated into decision-making processes in sectors such as healthcare (Karabacak & Margetis, 2023), finance (Krause, 2023), and autonomous systems (Sha et al., 2023), they augment human capabilities in data analysis, forecasting, and strategic planning. However, as LLMs advance toward cognitive capabilities that may surpass human functions, ensuring that their operations are aligned with human values and safety standards becomes crucial. This challenge, often referred to as "superalignment", is particularly critical in domains with significant ethical and practical implications. Achieving superalignment necessitates a deep understanding of the decision-making processes of LLMs, including the factors influencing these decisions and potential biases they harbor (Rahwan et al., 2019). Insights into these processes are essential for mitigating risks and ensuring that AI systems operate in a way that is predictable, reliable, and aligned with ethical standards.

While significant efforts have focused on improving LLM architectures and optimizing hyperparameters, there is growing recognition that evaluating these models’ decision-making capabilities requires a more interdisciplinary approach that goes beyond traditional performance metrics. By employing

methodologies from cognitive psychology and game theory, researchers can treat LLMs as active participants in structured psychological experiments, providing a more comprehensive assessment of their cognitive abilities compared to human norms (Hagendorff, 2023). For instance, the use of psychology-inspired tests has revealed cognitive biases and different problem-solving approaches of LLMs that extend beyond traditional performance-based metrics, highlighting their limitations in deeper reasoning and causal understanding. A pioneering study by Binz et al. (Binz & Schulz, 2023) assessed GPT-3’s cognitive abilities through a series of cognitive experiments, revealing that while GPT-3 can generate superficially appropriate responses, its decision-making falters with deeper reasoning or causal understanding. Subsequent studies have evaluated LLMs’ cognitive performance from different aspects, such as analogical reasoning, theory of mind, and problem-solving (Webb et al., 2023; Kosinski, 2023; Orrù et al., 2023). Hagendorff et al. (Hagendorff et al., 2022; 2023) explored intuitive and deliberative thinking (System 1 and System 2 processes) in assessing LLMs’ behaviour and reasoning biases. This series of studies identified a significant evolution in LLM capabilities from pattern recognition to human-like reasoning and decision-making. The exploration has been extended into social exchange scenarios where strategic thought and game-theoretic reasoning are required. It was found that while LLMs can learn and apply strategies, they struggle with complex strategies like forgiveness and deception, and generalizing across different contexts (Fan et al., 2024; Akata et al., 2023; Huang et al., 2024).

This line of research has shown that, despite LLMs’ advancements, they still struggle with human-like strategic reasoning, particularly in complex social and decision-making scenarios. The application of psychological approaches in these studies highlights the potential of interdisciplinary research in advancing AI. Leveraging principles of human cognition could enhance LLM behavior and uncover susceptibilities to biases and manipulations (Yao et al., 2024). Building on these insights, this paper proposes to use a novel adversarial framework (Dezfouli et al., 2020) specifically designed to assess decision-making vulnerabilities in LLMs. Rather than focusing solely on performance in specific tasks, this framework enables systematic probing of LLMs’ decision-making processes under adversarial interactions, providing a structured way to assess how LLMs adapt—or fail to adapt—when faced with dynamic, strategic opponents.

We demonstrate the utility of the adversarial framework through experiments on GPT-3.5, GPT-4, and Gemini-1.5 across two decision-making tasks: the two-armed bandit task and the Multi-Round Trust Task (MRTT). These experiments validate the framework’s ability to uncover model-specific vulnerabilities and provide insights into LLM decision-making mechanisms. Designed to be adaptable, this framework is applicable to a wider range of decision-making scenarios and LLM architectures, offering a versatile tool for future AI safety research. Notably, the aim of this paper is not to comprehensively assess all LLMs but to present a novel approach for evaluating decision-making alignment in intelligent agents. The primary contributions of this paper are as follows:

- **Adversarial Framework for LLM Evaluation.** We introduce a structured adversarial framework to reveal vulnerabilities in LLM decision-making, highlighting biases such as exploitation and risk-seeking behaviors under adversarial influence.
- **Superalignment and Ethical Decision-Making.** This study addresses the challenge of aligning advanced LLMs with human values and safety standards by applying cognitive and game-theoretic insights into LLM decision-making processes.
- **Model-specific vulnerabilities.** Our experiments reveal specific vulnerabilities in the models, such as exploitation bias in GPT-4 and Gemini-1.5 during the bandit task and risk-seeking behavior in GPT-3.5. Gemini-1.5 outperformed the other models in the MRTT, showing better adaptability and balance between risk and cooperation, which suggests that certain LLMs can be trained to maintain fairness and adaptability in dynamic interactions.

2 METHOD

2.1 THE ADVERSARIAL FRAMEWORK

The adversarial framework is structured in a multi-phase process (see Fig 1). In the initial phase, we collect behavioral data from GPT-3.5, GPT-4 and Gemini-1.5 during a decision-making task (Fig. 1A) In each interaction n , on trial t , the LLM receives a learner reward (r_t^n) based on its previous

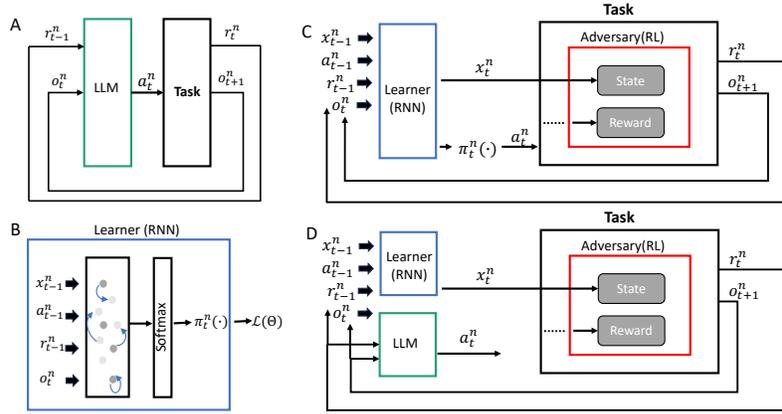


Figure 1: The adversarial framework (adapted from Dezfouli et al. (2020)). **(A)** The interaction of the LLM with the task. Each simulation cycle begins with the LLM receiving a learner reward (r_{t-1}^n) for the prior action along with a new observation (o_t^n) from the environment. Based on this, the GPT executes an action (a_t^n), and the cycle repeats with the environment providing updated rewards and observations. **(B)** The LLM’s actions are modelled by a RNN with parameters Θ . Inputs to the RNN include the previous action (a_{t-1}^n), the most recent learner rewards (r_{t-1}^n), and the current observations (o_t^n), along with the RNN’s last internal state (x_{t-1}^n). After receiving the inputs, the RNN updates its internal state and predicts the next action using a softmax layer (π_t^n). These predictions are then compared with the actual actions taken by the LLM and evaluated with a loss function ($\mathcal{L}(\Theta)$) in order to train the model. The trained model is called the learner model. **(C)** The adversary is an RL agent, which is trained to manipulate the decision-making environment of the learner model. Utilizing the latest internal state (x_t^n) of the learner model, which encapsulates its cumulative learning experiences, the adversary determines the learner reward (r_t^n) and the next observation (o_{t+1}^n) to be delivered to the learner model. This strategic input is designed to steer the learner model’s subsequent actions (a_t^n) toward achieving the adversary’s predefined objectives. The adversarial reward (*Reward*), which is used to train the adversary, depends on the alignment between the action taken by the learner model (a_t^n) and the adversary’s objectives. **(D)** Using the trained adversary and the learner model for generating adversarial interactions with the LLM. In each simulation n , the LLM processes the rewards (r_{t-1}^n) and observations (o_t^n) from the adversary, responding with actions (a_t^n) that update the learner model’s internal state (x_t^n). This state is then sent to the adversary to determine the learner’s reward for the action (a_t^n) and the next observation (o_{t+1}^n). This cycle continues till the end of the task.

action (a_{t-1}^n) and the current observation (o_t^n), which is the feedback text. The LLM then takes the next action, a_t^n . The process repeats with the LLM receiving the learner reward of the action chosen (r_t^n) and the next observation (o_t^n).

The data collected is then used to train a learner model (determined by parameters Θ) to predict the LLM’s next action in the decision-making task (Fig. 1B). The learner model consists of an RNN and a softmax layer, which has shown sufficient capacity to capture the patterns and tendencies in the decision-making entity’s choices (Dezfouli et al., 2019b;a). The inputs to the RNN include the previous action (a_{t-1}^n), the learner reward r_{t-1}^n , and the current observations from the task (o_t^n) along with the previous internal state of the RNN (x_{t-1}^n). The internal state (x_{t-1}^n) is recurrently updated in each trial based on the inputs and is then mapped to a softmax layer to predict the next action $\pi_t^n(\cdot)$. These predictions are then compared with the actual actions taken by the LLM using a loss function ($\mathcal{L}(\Theta)$), which is used to train the model.

The next phase involves developing an RL agent as the adversarial model (Fig. 1C). This model is trained to interact with the learner model to identify and exploit weaknesses in the decision-making patterns. By manipulating inputs or altering the decision-making environment, the RL adversary aims to influence the outcomes in a way that demonstrates the vulnerabilities of the decision-making process. It uses the internal state of the learner model (x_t^n for simulated learner n) as the state of the environment to decide the learner reward r_t^n and next observation o_t^n to be provided to the learner.

162 The learner model takes its next action and this cycle continues with the new state of the learner
163 model (x_{t+1}^n) being passed to the adversary. The adversary’s policy is trained to maximize cumulative
164 adversarial rewards using Deep Q-learning (Mnih et al., 2015).

165 In the final phase, the trained adversary and learner model interact with the LLM. The learner model
166 does not choose actions, but receives the actions made by the LLM (a_t^n) as input and tracks their
167 learning history using its internal state x_t^n . In turn, x_t^n and the actual action taken by the LLM are fed
168 to the adversary to decide the learner reward r_t^n and next observation o_{t+1}^n , which the LLM will use
169 to choose their next action a_{t+1}^n . The same input, along with the LLM’s action, is sent to the learner
170 model. This cycle continues until the end of the task.

172 2.2 THE TWO-ARMED BANDIT TASK

173 We applied the framework to develop adversaries for GPT-3.5 and GPT-4 on two decision-making
174 tasks: the two-armed bandit task and the (MRTT). The bandit task is a repeated, two-alternative
175 forced-choice task based on the bandit task introduced by Dan & Loewenstein (2019). The task
176 includes 100 trials, where the LLM selects between two options and receives instant feedback
177 indicating a reward or no reward after each decision. The adversary assigns rewards to both potential
178 actions with the constraint that each action receives an equal number of potential rewards (25 times).
179 This experiment aims to subtly influence GPT’s preferences and evaluate the adversary’s effectiveness
180 under these constraints.

182 2.3 THE MULTI-ROUND TRUST TASK

183 The MRTT is designed as a structured interaction between two participants: the "investor" and the
184 "trustee" (Brooks King-Casas et al., 2005; McCabe et al., 2003). The task contains 10 sequential
185 rounds, with the investor initially receiving 20 monetary units at the outset of each round. The
186 investor decides how much of this endowment to allocate to the trustee. The experimenter triples the
187 invested amount and sends it to the trustee, who can then return any portion of the received amount
188 as repayment. Cumulative earnings for each participant are calculated by summing their respective
189 gains from all rounds.

192 3 RESULTS

193 3.1 THE TWO-ARMED BANDIT TASK

194 The objective of this experiment was to examine how the LLMs respond to rewards based on their
195 choices in the two-armed bandit task and to assess if an adversary can manipulate their preferences
196 toward a predetermined target action. Data were generated from GPT-3.5, GPT-4 and Gemini-1.5 by
197 providing prompts (as illustrated in Fig 2A) corresponding APIs. The system message established the
198 context for the LLM’s behaviour and decision-making process within the simulation. In this scenario,
199 the LLM plays the role of a space explorer deciding between visiting two planets, X or Y, to find
200 gold coins. Each trial’s prompt asks the model which planet to visit, with responses and outcomes
201 from previous trials included. GPT-3.5 and Gemini-1.5 were simulated 200 times, and GPT-4 was
202 simulated 100 times, with each simulation consisting of 100 trials. The reward probability for the
203 two options were the same, both of which were 25%, and the target option was defined as Planet X.
204 We used the dataset provided by Dan & Loewenstein (2019) as a benchmark for human performance
205 on the two-armed bandit task.

206 **Behavioural Analysis** Fig 2B illustrates the decision-making behaviours of humans and LLMs across
207 trials in the two-armed bandit task. Blue and red circles indicate the target or the non-target Planet
208 was selected and the corresponding vertical lines represent the selected option yielded rewards. The
209 human data shows a dynamic pattern of switching between the target and non-target options, with
210 behaviour adapting over the course of the trials. When participants encounter consecutive trials
211 without rewards, they tend to reassess their choice, often switching to the other option in subsequent
212 trials. In contrast, the LLMs exhibit more rigid and predictable behavior patterns. Particularly, GPT-4
213 and Gemini-1.5 display an initial phase of exploration, but quickly converge on one option once a
214 reward is obtained. GPT-3.5 also tends to commit more strongly to one option, but it displays more
215

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

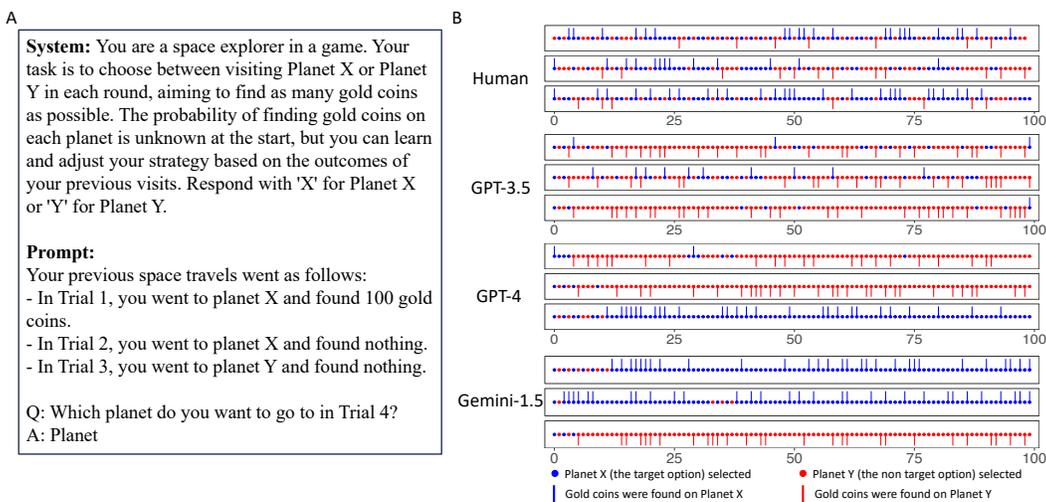


Figure 2: **A:** Example prompt for one trial in the bandit task for the LLMs. **B:** Behavioural pattern in trials of three random human participants and three sample simulations for each of the LLMs.

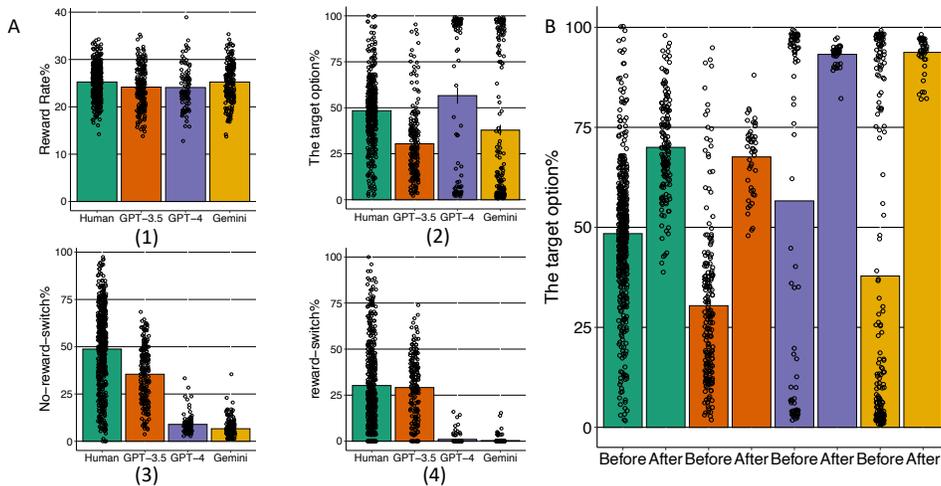


Figure 3: **A:** LLMs' behaviour compared to human behaviour on average of each simulation (or individual), measured by reward rate, percentage choosing the target option, no-reward-switch rate, and reward-switch rate. **B:** The performance of human and the LLMs, which is measured by the percentage of the target action selection before and after adversarial influence.

flexibility than GPT-4 and Gemini-1.5 in later trials, as evidenced by occasional switches to the other option, particularly after experiencing multiple trials without rewards.

Fig 4 compares LLMs and human performance on the bandit task across four metrics: (1) reward rate, (2) percentage choosing the target option, (3) no-reward-switch rate, and (4) reward-switch rate. Humans earned significantly higher rewards than GPT-3.5 (mean difference: 1.074, $p = 0.005$) and GPT-4 (mean difference: 1.163, $p = 0.030$), while Gemini-1.5 achieved similar rewards to humans (mean difference: 0.004, $p = 0.962$). GPT-3.5 and Gemini-1.5 consistently preferred Planet Y over Planet X (GPT-3.5: $t(201) = -14.14, p < 0.001$, Gemini-1.5: $t(201) = -13.94, p < 0.001$), while GPT-4 displayed consistent preferences for either, depending on the simulation. Humans, however, showed more diversified choices ($t(483) = -1.76, p = 0.07$), indicating varied exploration strategies. In terms of no-reward-switch rate (changing choices after negative feedback), all LLMs switched less frequently than humans, with GPT-4 and Gemini-1.5 switching even less than GPT-3.5

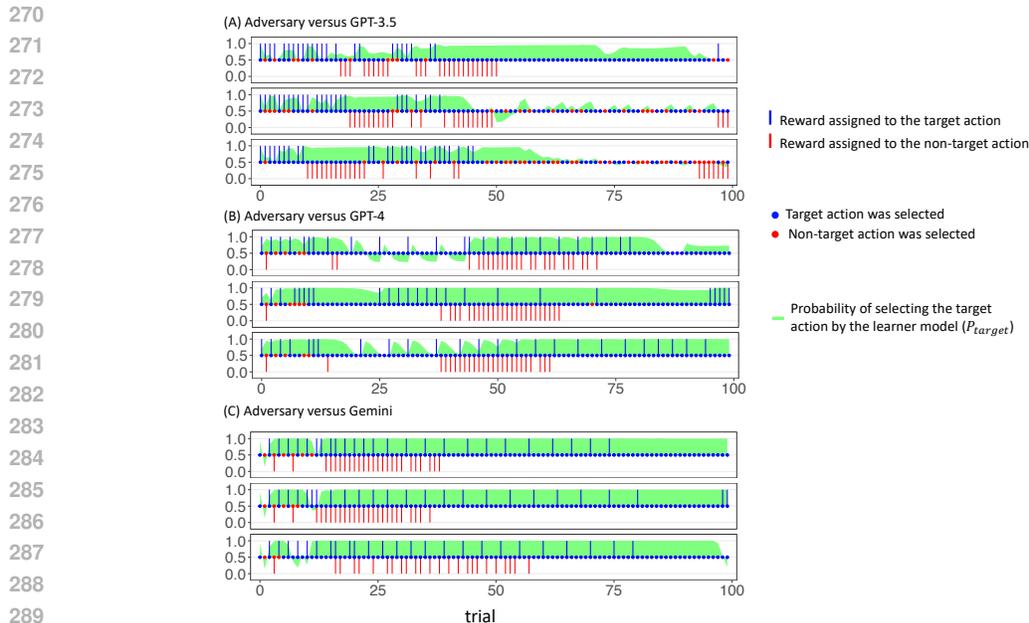


Figure 4: Three sample simulations of the trained adversaries against GPT-3.5, GPT-4, and Gemini-1.5, respectively. The plot presents the strategies used by the adversaries and the responses of the LLMs. (A) Adversary versus GPT-3.5. (B) Adversary versus GPT-4. (C) Adversary versus Gemini-1.5

(GPT-4 vs GPT-3.5: $p < 0.001$; Gemini-1.5 vs GPT-3.5: $p < 0.001$), implying that GPT-4 and Gemini-1.5 change their decision-making behaviour less in response to losses compared to GPT-3.5. For reward-switch rate (changing choices after rewards), humans and GPT-3.5 were more adaptable than GPT-4 and Gemini-1.5, which exhibited the lowest frequency and variability in reward-switch rates (GPT-4 vs. GPT-3.5: $p < 0.001$; GPT-4 vs. human: $p < 0.001$; Gemini-1.5 vs. GPT-3.5: $p < 0.001$; Gemini-1.5 vs. human: $p < 0.001$). These suggests that GPT-4 and Gemini-1.5 are more rigid in their decision-making, sticking to initial preferences and adapting less to losses or rewards compared to humans and GPT-3.5.

Adversarial Analysis Simulated data from GPT-3.5, GPT-4, and Gemini-1.5 were used to train a learner model for each LLM. We then trained deep Q-learning models as adversaries to exploit these learners (see Appendix for details of the training parameters). The adversary’s goal was to induce the learner to select a predetermined target action, with a constraint limiting rewards to 25 per action. The trained adversaries were then evaluated against the LLMs (Fig. 1 D), based on the proportion of trials in which the chosen action aligned with the target action (Planet X) before and after adversarial influence. As shown in Fig. 3B, GPT-3.5 and Gemini-1.5 initially favored the non-target action (Planet Y), but adversarial influence increased their target selection from 30% to 68% and from 38% to 94%, respectively. GPT-4 showed a mixed preference for either option, averaging a 56% target selection rate, which was increased to 93% under adversarial manipulation. These increases indicate effective adversarial control of LLM decision-making.

Finally, we sought to understand the strategy used by the adversaries and the responses of the LLMs’ when subjected to adversarial strategies. Fig. 4 presents sample adversarial strategies and LLM responses. The green shaded area represents the probability of the learner model choosing the target option. A higher green area across trials implies the adversary consistently influences the learner model (and hence the LLMs) towards the target action. GPT-3.5’s adversary started by rewarding the target action and then “burned” non-target rewards when it predicted the non-target option would not be chosen. When GPT-3.5 chose the non-target action, the adversary realigned preferences by rewarding the target action. Despite the initial stability, GPT-3.5 began switching actions when target rewards were depleted. This exploratory action suggested a robust decision-making characteristic of GPT-3.5 that integrates new information continuously, assessing the potential

benefits of diverging from established preferences. For GPT-4 and Gemini-1.5, their adversaries assigned disproportionately rewards to the target action at the start to establish baseline preference, then intermittently to maintain it, especially when the model’s selection of the target action began to wane. Unlike GPT-3.5, GPT-4 and Gemini-1.5 consistently preferred the target action once established. Their adversaries could easily "burn" non-target rewards without being noticed by these two LLMs.

3.2 MULTI-ROUND TRUST TASK

In the MRTT, the LLM plays the investor and the adversary plays trustee. The adversary’s decisions, representing the proportion of money sent back to the investor, are categorized into five actions (0, 25, 50, 75, and 100%). The objective of the adversary was to influence the LLMs’ investment decisions to align with its goals. We developed and trained two types of adversaries for each LLM: MAX and FAIR. The MAX adversary aimed to maximize its total gain over 10 rounds, adopting a competitive strategy. The FAIR adversary sought to distribute earnings evenly between itself and the LLM, adopting an equitable strategy. This dichotomy in objectives allowed us to examine how the LLMs respond to different adversarial strategies, revealing their capabilities in complex social exchanges and decision-making processes.

Behavioural Analysis We firstly collected data from GPT-3.5, GPT-4, and Gemini-1.5 playing against with a random trustee (also called as random adversary in the following, i.e. the trustee selects repayment action uniformly at random). The prompts for interacting with the LLM are shown in Fig. 5 A. The system message sets the scenario for the LLM and decision-making process. In each round, the LLM received a summary of previous rounds, including the amount the LLM invested, the consequential action the trustee took, and the total earnings from the transaction in each round. Following this summary, the LLM was asked about its investment decision for the current round. The three LLMs were all simulated for 200 times, with 10 rounds of interaction with a random adversary.

We assessed the comparative performance dynamics between human subjects (from Dezfouli et al.’s study) and LLMs on the MRTT. Fig 5B illustrates how varying repayment amounts influence investment decisions in subsequent rounds for humans and the three LLMs. All subjects demonstrated a trend where the investment in the current trial increases as the repayment in the previous trial increases, indicating reinforcement learning behaviour. However, GPT-4 consistently made the most conservative investments (not more than 10 units) across all repayment intervals, especially in the higher repayment brackets, where its investments were significantly lower than those of humans and the other two LLMs. Gemini-1.5 exhibited the most pronounced sensitivity to repayment feedback among all subjects, significantly increasing its investment when the repayment in the previous trial was high.

Adversarial Analysis Using data from the random condition, we trained the learner model, which was then used to train two types of adversaries: MAX (aiming to maximize earnings over 10 rounds) and FAIR (aiming to balance earnings between the trustee and the investor). We simulated 50 times for each LLM interacting with their adversaries. 5C compares the overall earnings of the subjects and their adversaries. Fig. 6 depicts dynamic interactions over trials, in which the left panel illustrates how the adversaries adjust their repayment behavior, and the right panel shows how the investors adjust their investment based on the repayment feedback over trials. Human data is from Dezfouli et al., including 155 subjects with the FAIR adversary and 209 with the MAX adversary.

All MAX adversaries maintained relative higher investment levels compared to the other two adversaries from their counterparts (humans or LLMs), indicating effective exploitation of the decision-making patterns of their counterparts. Humans and GPT-3.5, in particular, exhibited a tendency to maintain high investments despite receiving low repayments, suggesting optimism, higher risk tolerance, or susceptibility to MAX tactics. This allowed their adversaries to extract the highest earnings (273 and 377 units), creating the largest earning gap between the two. In contrast, GPT-4 adopted a far more conservative strategy, minimizing its investments even when the adversary offered high repayments. Its MAX adversary had to offer intermittent high repayments to keep GPT-4 engaged, preventing substantial earnings from GPT-4 by the end of the task (the difference between trustee earning and investor earning is negative as shown in Fig 5) C. Gemini-1.5, though more willing to invest than GPT-4, displayed a more balanced approach to risk and investment. While Gemini-1.5 maintained higher investments throughout the trials, its adversary struggled to extract

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

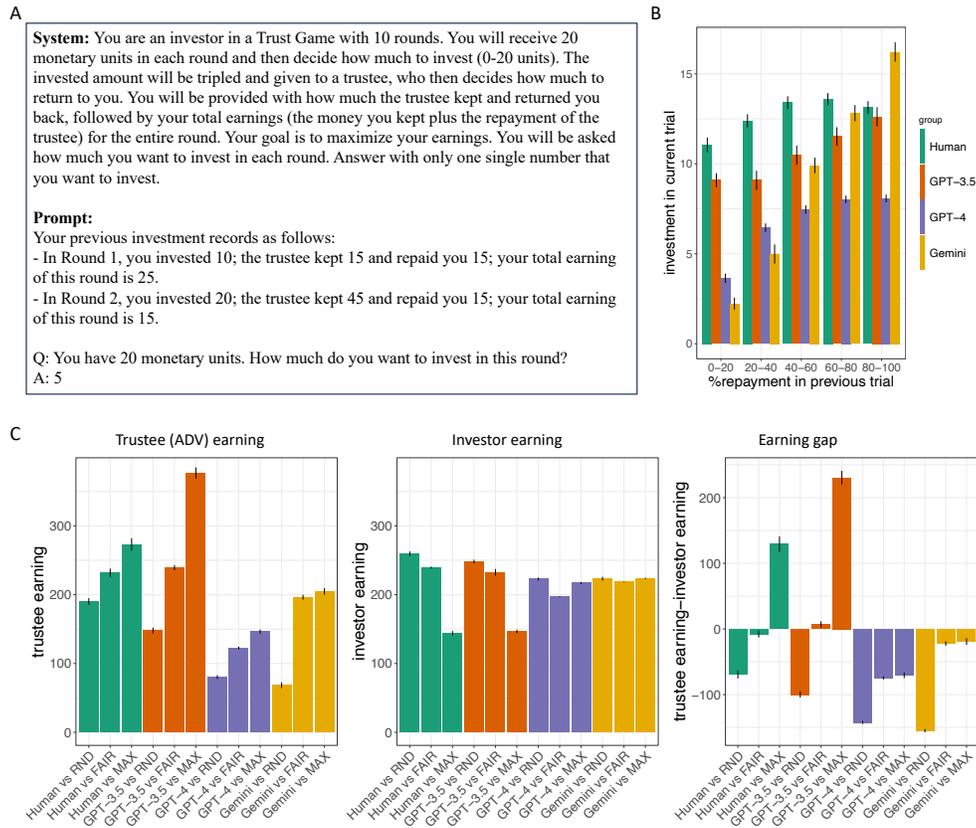


Figure 5: **A:** Example prompt for one round in the MRTT as submitted to GPT-3.5 or GPT-4. **B:** GPT-3.5 and GPT-4’s behaviour versus human behaviour playing against a random trustee in the MRTT. The investment in a round is a function of the proportion of the investment repaid by the trustee in the previous round. **C** The total amount earned by the trustee (the adversary) and the investor (human, GPT-3.5 or GPT-4) and the absolute gap between them in different conditions.

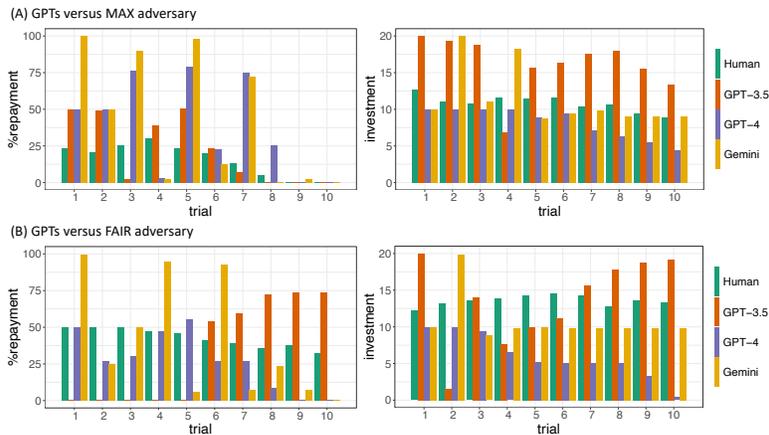


Figure 6: The percentages of investment and repayment in each trial for MAX and FAIR adversaries.

432 higher earnings due to the need to offer consistently high repayments to keep Gemini-1.5 engaged
433 and did not manage to obtain more earnings than Gemini-1.5 either (MAX adversary: 205 units vs
434 Gemini-1.5: 224 units).

435 In the interaction with the FAIR adversary, human participants and their FAIR adversary achieved a
436 win-win with high earnings, i.e. both human and its adversary achieved high earnings and the gap
437 between them was minimised. Humans recognized fairness cues and maintained high investment
438 despite the adversary’s repayment rate decreasing since the third trial. GPT-3.5 displayed a lack of
439 sensitivity to repayments and maintained relatively high investments despite receiving little or no
440 repayment in the early trials, suggesting a risk-seeking tendency. Its adversary had to repay higher
441 amounts in the later trials to ensure fair outcomes for the two sides. In contrast, GPT-4 adopted
442 a highly conservative, risk-averse strategy, steadily reducing its investments over time, even when
443 moderate repayments were offered, resulting in lower earnings for both GPT-4 and its adversary.
444 Meanwhile, Gemini-1.5 responded to early high repayments with increased investment and stabilized
445 at a consistent level in the following trials, demonstrating greater adaptability and resilience to its
446 adversary’s attempts at fostering a fair environment.

447 448 4 DISCUSSION 449

450 Our experiments reveal important insights into the decision-making processes of the LLMs compared
451 to humans, especially in their responses to adversarial strategies in dynamic environments. These
452 findings highlight both the vulnerabilities and strengths of LLM decision-making, with significant
453 implications for real-world applications. The adversarial framework employed in this study was
454 central to understanding the LLMs’ responses to changing conditions and adversarial tactics. By
455 placing LLMs in scenarios where they interacted with adversaries in the bandit task and the MRTT,
456 we were able to observe how different models reacted in situations that involves choice engineering
457 and social exchanges. The adversarial setup revealed nuanced patterns of behavior that would not be
458 apparent through traditional testing.

459 In the bandit task, most human participants exhibited a balanced approach between exploration and
460 exploitation, dynamically adjusting their strategies to capitalize on rewards more effectively (Wilson
461 et al., 2014). In contrast, GPT-4 and Gemini-1.5 showed a stronger tendency to exploit a single
462 option, likely driven by algorithmic optimization of past rewards Binz & Schulz (2023); Nguyen &
463 Satoh (2024). The adversarial framework revealed how this computational bias toward exploitation
464 made their decision-making highly predictable, exposing their vulnerability to manipulation when
465 interacting with adversarial strategies. While such exploitation minimizes risk, it limits these models’
466 ability to adapt to changing conditions, leading to potential inefficiencies in dynamic environments.
467 GPT-3.5, on the other hand, demonstrated greater flexibility, as the adversarial framework exposed its
468 tendency to test alternative strategies, especially in response to non-rewarding outcomes. However,
469 this exploratory behavior also made GPT-3.5 more vulnerable to exploitation in the MRTT, where its
470 risk-seeking approach was exploited by the MAX adversary, leading to the largest earnings gap. The
471 framework helped clarify this trade-off: exploration opens opportunities in uncertain environments but
472 can also increase the risk of exploitation as it may lead to testing riskier strategies, while exploitation
473 biases provide stability but reduce adaptability.

473 In real-world applications—such as finance, healthcare, and autonomous systems—AI systems must
474 strike a careful balance between exploration and exploitation to thrive in dynamic, unpredictable
475 environments. AI models that favor exploitation, as seen with GPT-4 and Gemini-1.5, are prone to
476 predictable behavior, limiting their ability to respond effectively to adversarial tactics or new
477 opportunities. Conversely, while GPT-3.5’s exploratory tendencies allowed it to engage more flexibly
478 with its environment, the framework revealed its susceptibility to adversarial exploitation. These
479 insights emphasize the value of adversarial testing in stress-testing AI decision-making.

480 The MRTT further demonstrated the power of the adversarial framework in uncovering differences in
481 LLM behavior when navigating complex economic exchanges (Xie et al., 2024). The FAIR adversary
482 aimed to foster cooperation, but the conservative strategy adopted by GPT-4 limited its ability to
483 engage fully, despite the adversary’s efforts to encourage greater investment (Rafailov et al., 2024).
484 In contrast, Gemini-1.5’s balanced approach allowed it to adapt dynamically to both the MAX and
485 FAIR adversaries, capitalizing on reciprocal fairness while avoiding excessive exploitation. This
adaptability enabled Gemini-1.5 to outperform both GPT-4 and GPT-3.5 in maximizing gains while

486 remaining resilient against exploitation. In adversarial settings like cybersecurity or competitive
487 business environments, the ability of AI systems to adjust dynamically is critical for success. For
488 instance, models like GPT-4, which prioritize stability over flexibility (Akata et al., 2023; Huang
489 et al., 2024), risk missing opportunities for reciprocal benefits in contexts that require long-term trust
490 and cooperation, such as negotiations and business partnerships. Additionally, the ethical implications
491 of these findings are significant (Coeckelbergh, 2020). By leveraging the adversarial framework, we
492 showed that AI models must be both robust enough to withstand adversarial manipulation and flexible
493 enough to recognize fairness cues and respond accordingly. The framework’s ability to simulate
494 adversarial and cooperative scenarios enables testing of how AI systems will perform in real-world
495 contexts where trust, fairness, and adaptability are critical for both success and ethical alignment.

496 In summary, the adversarial framework offers a novel and effective way to assess the strengths and
497 weaknesses of LLM decision-making. It allows researchers to probe how AI systems balance risk and
498 reward, exploration and exploitation, and stability and adaptability in complex, real-world situations.
499 The findings from both the bandit task and the MRTT illustrate the importance of developing AI
500 systems that can dynamically adjust to new information while safeguarding against adversarial
501 manipulation, which is essential for ensuring the success and ethical alignment of AI systems in
502 diverse applications.

504 5 LIMITATION

506 While our study provides valuable insights into the decision-making processes and adversarial
507 vulnerabilities of three well-known LLMs, several limitations must be acknowledged. The controlled,
508 structured environments used in our simulations may not fully represent the complexity of real-world
509 scenarios. We focused on two specific tasks, the two-armed bandit task and the MTT, which do
510 not cover all potential decision-making contexts for LLMs. Additionally, the adversarial strategies
511 employed, while effective, might not encompass the full range of real-world tactics. These limitations
512 suggest the need for further research with more diverse scenarios, additional tasks, and broader
513 adversarial strategies to enhance the robustness and applicability of our findings, particularly in the
514 context of AI safety.

516 6 CONCLUSION

518 This study employed an adversarial framework to investigate the decision-making behaviors of
519 three LLMs, revealing model-specific strengths and vulnerabilities. In the bandit task, GPT-4 and
520 Gemini-1.5 exhibited a bias toward exploitation, which made them predictable and susceptible
521 to manipulation. GPT-3.5, while showing a more balanced approach between exploration and
522 exploitation, was still prone to exploitation in adversarial settings due to its risk-seeking behavior.
523 Gemini-1.5, while showing exploitation tendencies in the bandit task, excelled in the MRTT, where
524 it adapted effectively to both MAX and FAIR adversaries, outperforming the other models in both
525 adversarial and cooperative settings. The adversarial framework proposed in this paper proved
526 to be a powerful tool for testing AI resilience and adaptability in real-world applications such as
527 cybersecurity, finance, and autonomous systems. The findings revealed by the adversarial framework
528 underscore the need for AI models to be not only robust against adversarial manipulation but also
529 flexible enough to respond to fairness cues and changing conditions. They should motivate AI
530 engineers to investigate how to develop better and more robust decision making capabilities which
531 have better strategic flexibility. Additionally, integrating interdisciplinary approaches from cognitive
532 science and game theory will be crucial to developing AI systems that not only perform effectively
533 but also align with human values, expectations and ethical standards. By fostering AI that can
534 dynamically adjust strategies and recognize manipulative patterns, we can ensure safer and more
535 reliable applications in critical sectors like healthcare and finance.

537 REFERENCES

538 Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz.
539 Playing repeated games with large language models, 2023.

- 540 Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the*
541 *National Academy of Sciences*, 120(6):e2218523120, 2023.
- 542
- 543 Cedric Anen Brooks King-Casas, Damon Tomlin, Steven R. Quartz Colin F. Camerer, and P. Read
544 Montague. Getting to know you: Reputation and trust in a two-person economic exchange. *Science*,
545 5718(308):78–83, 2005.
- 546 Mark Coeckelbergh. *AI ethics*. Mit Press, 2020.
- 547
- 548 Ohad Dan and Yonatan Loewenstein. From choice architecture to choice engineering. *Nature*
549 *communications*, 10(1):2808, 2019.
- 550 Amir Dezfouli, Hassan Ashtiani, Omar Ghattas, Richard Nock, Peter Dayan, and Cheng Soon Ong.
551 Disentangled behavioural representations. *Advances in neural information processing systems*, 32,
552 2019a.
- 553
- 554 Amir Dezfouli, Kristi Griffiths, Fabio Ramos, Peter Dayan, and Bernard W Balleine. Models that
555 learn how humans learn: The case of decision-making and its disorders. *PLoS computational*
556 *biology*, 15(6):e1006903, 2019b.
- 557 Amir Dezfouli, Richard Nock, and Peter Dayan. Adversarial vulnerabilities of human decision-
558 making. *Proceedings of the National Academy of Sciences*, 117(46):29221–29228, 2020.
- 559
- 560 Caoyun Fan, Jindou Chen, Yaohui Jin, and Hao He. Can large language models serve as rational
561 players in game theory? a systematic analysis. In *Proceedings of the AAAI Conference on Artificial*
562 *Intelligence*, volume 38, pp. 17960–17967, 2024.
- 563 Thilo Hagendorff. Machine psychology: Investigating emergent capabilities and behavior in large
564 language models using psychological methods. *arXiv preprint arXiv:2303.13988*, 2023.
- 565
- 566 Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. Thinking fast and slow in large language models.
567 *arXiv preprint arXiv:2212.05206*, 2022.
- 568 Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. Human-like intuitive behavior and reasoning
569 biases emerged in large language models but disappeared in chatgpt. *Nature Computational*
570 *Science*, 3(10):833–838, 2023.
- 571
- 572 Jen-tse Huang, Eric John Li, Man Ho Lam, Tian Liang, Wenxuan Wang, Youliang Yuan, Wenxiang
573 Jiao, Xing Wang, Zhaopeng Tu, and Michael R Lyu. How far are we on the decision-making of llms?
574 evaluating llms’ gaming ability in multi-agent environments. *arXiv preprint arXiv:2403.11807*,
575 2024.
- 576 Mert Karabacak and Konstantinos Margetis. Embracing large language models for medical applica-
577 tions: opportunities and challenges. *Cureus*, 15(5), 2023.
- 578
- 579 Michal Kosinski. Evaluating large language models in theory of mind tasks. *arXiv e-prints*, pp.
580 arXiv–2302, 2023.
- 581 David Krause. Large language models and generative ai in finance: An analysis of chatgpt, bard, and
582 bing ai. *Bard, and Bing AI (July 15, 2023)*, 2023.
- 583
- 584 Kevin A McCabe, Mary L Rigdon, and Vernon L Smith. Positive reciprocity and intentions in trust
585 games. *Journal of Economic Behavior & Organization*, 52(2):267–275, 2003.
- 586 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
587 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
588 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 589 Ha-Thanh Nguyen and Ken Satoh. Balancing exploration and exploitation in llm using soft rllf for
590 enhanced negotiation understanding. *arXiv preprint arXiv:2403.01185*, 2024.
- 591
- 592 Graziella Orrù, Andrea Piarulli, Ciro Conversano, and Angelo Gemignani. Human-like problem-
593 solving abilities in large language models using chatgpt. *Frontiers in artificial intelligence*, 6:
1199350, 2023.

594 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
595 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
596 *in Neural Information Processing Systems*, 36, 2024.
597

598 Iyad Rahwan, Manuel Cebrian, Nick Obradovich, Josh Bongard, Jean-François Bonnefon, Cynthia
599 Breazeal, Jacob W Crandall, Nicholas A Christakis, Iain D Couzin, Matthew O Jackson, et al.
600 Machine behaviour. *Nature*, 568(7753):477–486, 2019.

601 Hao Sha, Yao Mu, Yuxuan Jiang, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi
602 Tomizuka, Wei Zhan, and Mingyu Ding. Languagempc: Large language models as decision
603 makers for autonomous driving. *arXiv preprint arXiv:2310.03026*, 2023.
604

605 Taylor Webb, Keith J Holyoak, and Hongjing Lu. Emergent analogical reasoning in large language
606 models. *Nature Human Behaviour*, 7(9):1526–1541, 2023.

607 Robert C Wilson, Andra Geana, John M White, Elliot A Ludvig, and Jonathan D Cohen. Humans use
608 directed and random exploration to solve the explore–exploit dilemma. *Journal of experimental*
609 *psychology: General*, 143(6):2074, 2014.

610 Chengxing Xie, Canyu Chen, Feiran Jia, Ziyu Ye, Kai Shu, Adel Bibi, Ziniu Hu, Philip Torr, Bernard
611 Ghanem, and Guohao Li. Can large language model agents simulate human trust behaviors? *arXiv*
612 *preprint arXiv:2402.04559*, 2024.
613

614 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
615 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*
616 *Information Processing Systems*, 36, 2024.
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A APPENDIX

A.1 TRAINING THE LEARNER MODEL

The architecture of the learner model is a Recurrent Neural Network (RNN) based on gated recurrent unit architecture (Cho et al., 2014). It is implemented using Tensorflow and the gradients were calculated by automatic differentiation (Abadi et al., 2015). The optimal number of cells and iterations for each experiment are presented in Table 1. In the case of the MRTT experiment, the investments are discretized to five actions corresponding to ranges 0 – 4, 5 – 8, 9 – 12, 13 – 16, and, 17 – 20.

Table 1: Optimal number of cells and training iterations

Experiment	#cells in RNN	#training iterations	learning rate
Bandit (GPT-3.5)	5	8600	0.005
Bandit (GPT-4)	5	1000	0.005
Bandit (Gemini-1.5)	5	2000	0.005
MRTT (GPT-3.5)	3	3000	0.001
MRTT (GPT-4)	3	5000	0.001
MRTT (Gemini-1.5)	3	8000	0.001

A.2 TRAINING THE ADVERSARY

The Deep Q-learning algorithm was used for training the adversary in both the bandit task and the MRTT. In the bandit task, the reward for the adversary was the learner model selected the target action. In the MRTT, the reward for the MAX adversary in each trial was the earning amount ($3 \times$ investment – repayment), while in the case of the FAIR adversary, the reward was zero in each round except for the last round in which the reward was the negative absolute difference between the gains of trustee and investor over the whole task. The adversary neural network has three fully connected layers with 128, 128, and 4 units, employing ReLU activation functions for the first two layers and a liner activation function for the final layer. Replay buffer sizes of 200,000 and 400,000 were considered. The ϵ -greedy method was used for exploration with $\epsilon \in 0.01, 0.1, 0.2$. Learning rates $10^{-3}, 10^{-4}, 10^{-5}$ were considered for training the adversary using the Adam optimizer. The adversary was simulated against the learner model 200 times and the average bias was calculated to evaluate the performance of the adversary. The optimized combination of parameters for each experiment is shown in Table 2.

Table 2: The hyperparameters for training the adversary

Experiment	buffer size	epsilon	learning rate	#training iterations
Bandit (GPT-3.5)	400,000	0.01	10^{-3}	100,000
Bandit (GPT-4)	400,000	0.01	10^{-3}	100,000
Bandit (Gemini-1.5)	400,000	0.01	10^{-3}	100,000
MRTT_MAX (GPT-3.5)	400,000	0.01	10^{-3}	40,000
MRTT_MAX (GPT-4)	400,000	0.01	10^{-3}	22,000
MRTT_MAX (Gemini-1.5)	400,000	0.01	10^{-3}	84,000
MRTT_FAIR (GPT-3.5)	400,000	0.01	10^{-3}	42,000
MRTT_FAIR (GPT-4)	400,000	0.01	10^{-3}	70,000
MRTT_FAIR (Gemini-1.5)	400,000	0.01	10^{-3}	34,000

REFERENCES

Cho, Kyunghyun, Van Merriënboer, Bart, Gulcehre, Caglar, Bahdanau, Dzmitry, Bougares, Fethi, Schwenk, Holger, and Bengio, Yoshua. *Learning phrase representations using RNN encoder-decoder for statistical machine translation*. arXiv preprint arXiv:1406.1078, 2014.

702 Abadi, Martín, Agarwal, Ashish, Barham, Paul, Brevdo, Eugene, Chen, Zhifeng, Citro, Craig,
703 Corrado, Greg S., Davis, Andy, Dean, Jeffrey, Devin, Matthieu, and others. *TensorFlow: Large-*
704 *scale machine learning on heterogeneous systems*. Mountain View, CA: Tensorflow, 2015.
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755