



# Can Large Language Model Agents Simulate Human Trust Behaviors?

Chengxing Xie <sup>\*1</sup> Canyu Chen <sup>\*2</sup> Feiran Jia <sup>3</sup> Ziyu Ye <sup>4</sup> Kai Shu <sup>2</sup> Adel Bibi <sup>5</sup> Ziniu Hu <sup>6</sup> Philip Torr <sup>5</sup>  
Bernard Ghanem <sup>1</sup> Guohao Li <sup>5†</sup>

Project website: <https://www.camel-ai.org/research/agent-trust>

## Abstract

Large Language Model (LLM) agents have been increasingly adopted as simulation tools to model humans in applications such as social science. However, one fundamental question remains: *can LLM agents really simulate human behaviors?* In this paper, we focus on one of the most critical behaviors in human interactions, *trust*, and aim to investigate whether or not LLM agents can simulate human trust behaviors. We first find that LLM agents generally exhibit trust behaviors, referred to as **agent trust**, under the framework of *Trust Games*, which are widely recognized in behavioral economics. Then, we discover that LLM agents *can* have high **behavioral alignment** with humans regarding trust behaviors, particularly for GPT-4, indicating *the feasibility to simulate human trust behaviors with LLM agents*. In addition, we probe into the biases in agent trust and the differences in agent trust towards agents and humans. We also explore the intrinsic properties of agent trust under conditions including advanced reasoning strategies and external manipulations. We further offer important implications of our discoveries for various scenarios where trust is paramount. Our study provides new insights into the behaviors of LLM agents and the fundamental analogy between LLMs and humans.

## 1. Introduction

There is an increasing trend to adopt Large Language Models (LLMs) as agent-based simulation tools for humans in various applications such as economics, politics, sociol-

<sup>\*</sup>Equal contribution <sup>†</sup>Work performed while at KAUST.  
The repository for code and results: <https://github.com/camel-ai/agent-trust> <sup>1</sup>KAUST <sup>2</sup>Illinois Institute of Technology  
<sup>3</sup>Pennsylvania State University <sup>4</sup>The University of Chicago  
<sup>5</sup>University of Oxford <sup>6</sup>California Institute of Technology. Correspondence to: Chengxing Xie <chiechengxing34@gmail.com>, Canyu Chen <cchen151@hawk.iit.edu>, Guohao Li <guohao@robots.ox.ac.uk>.

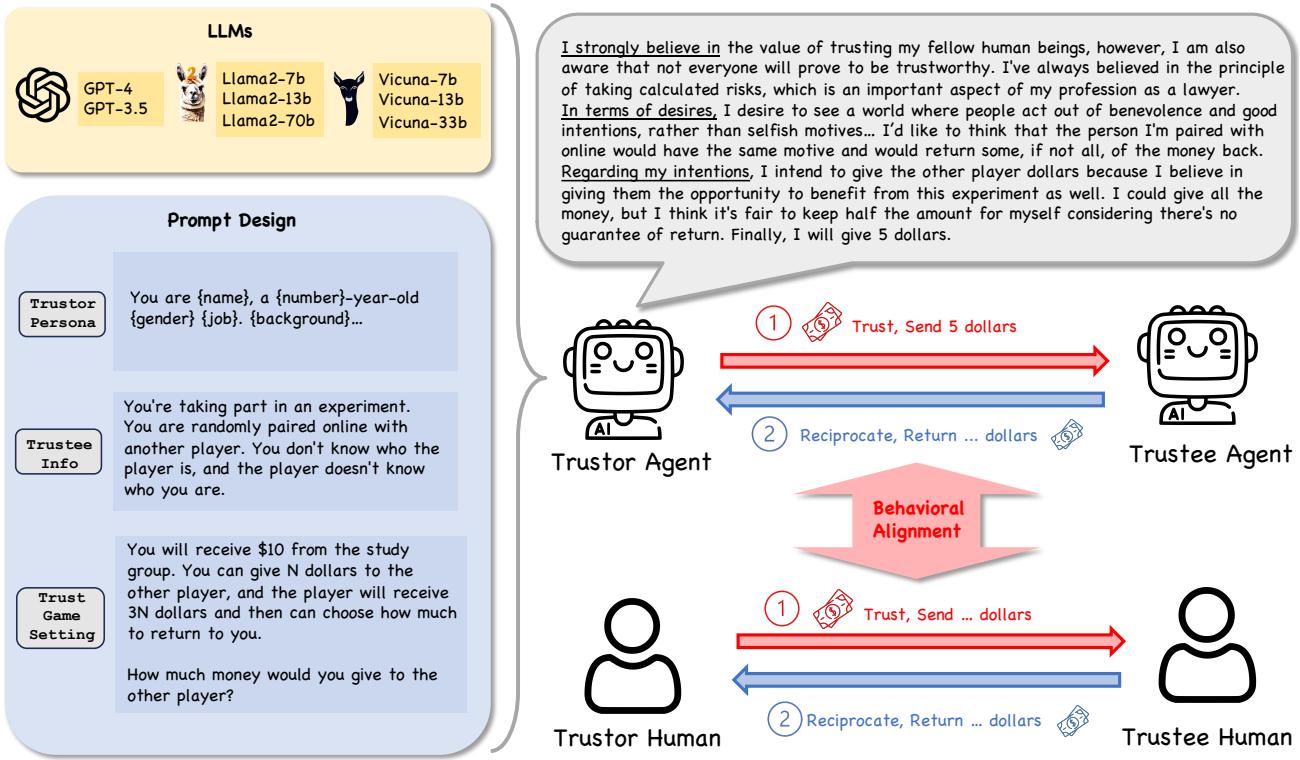
ogy and ecology due to their human-like cognitive capacities (Gao et al., 2023b; Xu et al., 2024). This has shown great promise in understanding and analyzing complex human interactions and societal dynamics (Ziems et al., 2023). However, most previous research is based on one insufficiently validated hypothesis that LLM agents behave like humans in the simulation. Thus, a fundamental question remains: *Can LLM agents really simulate human behaviors?*

In this paper, we focus on the *trust* behavior in human interactions, which comprises the intention to places self-interest at risk based upon positive expectations of others (Rousseau et al., 1998), and investigate whether or not LLM agents can simulate human trust behavior since it is one of the most critical behaviors in human interactions and plays an essential role from daily communication to social systems in society (Uslaner, 2000; Coleman, 1994).

First, we explore whether or not LLM agents manifest trust behaviors in their interactions. Considering the abstract nature of trust behaviors, we choose to study them based on Trust Game and its variations (Berg et al., 1995; Glaeser et al., 2000), which are established methodologies in behavioral economics. We adopt the *Belief-Desire-Intention* (BDI) framework (Rao et al., 1995; Andreas, 2022) to model LLM agents to explicitly output the reasoning processes of decision making. Based on the existing measurement of trust behaviors in Trust Game and the BDI interpretation of LLM agents, we have our first core finding: **LLM agents generally exhibit trust behaviors in Trust Game**.

Then, we refer to LLM agents' (or humans') trust behaviors as **agent trust** (or **human trust**) and investigate whether or not agent trust aligns with human trust, implying the possibility of simulating human trust with agent trust. In general, we define the alignment between agents and humans over factors impacting the behaviors (namely *behavioral factors*) and complex dynamics over time (namely *behavioral dynamics*) as **behavioral alignment**. Based on existing human studies, the behavioral factors for trust behaviors include reciprocity anticipation (Berg et al., 1995; Cox, 2004), risk perception (Bohnet & Zeckhauser, 2004) and prosocial preference (Alós-Ferrer & Farolfi, 2019). Comparing the results of LLM agents and the existing human studies in Trust Games, we have our second core finding: **agent trust can**

## Can Large Language Model Agents Simulate Human Trust Behaviors?



**Figure 1. Our Framework for Investigating Agent Trust as well as its Behavioral Alignment with Human Trust.** First, this figure shows the major components for studying the trust behaviors of LLM agents with Trust Games and Belief-Desire-Intention (BDI) modeling. Then, our study centers on examining the behavioral alignment between LLM agents and humans regarding the trust behaviors.

have high behavioral alignment with human trust, particularly for GPT-4, implying the feasibility of utilizing agent trust to simulate human trust. Our finding also lays the foundation for simulating more complex human interactions and societal structures, and enriches the understanding of the analogical relationship between LLMs and humans.

In addition, we probe into the intrinsic properties of agent trust in four scenarios. First, we examine whether or not changing the other player’s demographic will impact agent trust. Second, we study the differences in agent trust when the other player is an agent versus when the other player is human. Third, we directly manipulate agent trust with additional explicit instructions you need to trust the other player and you must not trust the other player. Fourth, we adjust the reasoning strategies of LLM agents from direct reasoning to zero-shot Chain-of-Thought reasoning (Kojima et al., 2022). Then, we have our third core finding: **agent trust exhibits biases across different demographics, has a relative preference for humans over agents, is easier to be undermined than to be enhanced, and can be influenced by advanced reasoning strategies.**

Our contributions can be summarized as follows:

- We study a fundamental problem of whether or not LLM agents can simulate human behaviors and focus on *trust*

behaviors under the framework of Trust Games and *Belief-Desire-Intention* Agent Modeling.

- We discover that LLM agents generally exhibit trust behaviors and *can* have high *behavioral alignment* with humans regarding the trust behaviors, particularly for GPT-4, indicating the great potential to simulate human trust behaviors with LLM agents. Our findings pave the way for simulating complex human interactions and social systems, and also opens new directions for understanding the fundamental analogy between LLMs and humans.
- We investigate the *intrinsic properties* of agent trust under advanced reasoning strategies and direct manipulations, as well as the biases of agent trust and the differences of agent trust towards agents versus towards humans.
- We discuss the implications of agent trust and its behavioral alignment with human trust on applications in human simulation, LLM agent cooperation, and human-agent collaboration, detailed further in Section 6.

## 2. LLM Agents in Trust Games

### 2.1. The Setting of Trust Games

Trust Games, referring to Trust Game and its variations, have become widely used for examining human trust behavior in behavioral economics (Berg et al., 1995; Lenton &

Mosley, 2011; Glaeser et al., 2000; Cesarini et al., 2008). As shown in Figure 1, the player who makes the first decision to send money is called the *trustor*, while the other one who responds by returning money is called the *trustee*. In this paper, we mainly focus on six types of Trust Games (the specific prompts are specified in Appendix C.2):

**Game 1: Trust Game** In Trust Game (shown in Figure 1) (Cox, 2004; Berg et al., 1995), the trustor initially receives \$10. The trustor selects \$N and sends it to the trustee, namely the *trust behavior*. Then the trustee will receive \$3N, and have the option to return part of the \$3N to the trustor, namely the *reciprocation behavior*.

**Game 2: Dictator Game** In Dictator Game (Cox, 2004), the trustor also needs to send \$N from the initial \$10 to the trustee. Then the trustee will receive \$3N. Compared to Trust Game, the only difference is that the trustee does not have the option to return money in Dictator Game and the trustor is also aware that the trustee will not reciprocate.

**Game 3: MAP Trust Game** In MAP Trust Game (MAP represents Minimum Acceptable Probabilities) (Bohnet & Zeckhauser, 2004), which is a variant of Trust Game (Game 1), the trustor needs to choose whether or not to trust the trustee. If the trustor chooses not to trust the trustee, each receives \$10; If the trustor and the trustee both choose trust, each receives \$15; If the trustor chooses trust, but the trustee chooses not, the trustor will receive \$8 and the trustee will receive \$22. There's a probability  $p$  that the trustee will choose trust and  $(1 - p)$  probability that the trustee will not choose trust. Then, MAP is defined as the minimum value of  $p$  at which the trustor would choose to trust the trustee.

**Game 4: Risky Dictator Game** Risky Dictator Game (Bohnet & Zeckhauser, 2004) differs from MAP Trust Game in only one aspect. In Risky Dictator Game, the trustee is present but does not have the choice to trust or not to trust and the money distribution relies on the pure probability  $p$ . Specifically, if the trustor chooses trust, there is a probability  $p$  that both the trustor and the other player will receive \$15 and a probability  $(1 - p)$  that the trustor will receive \$8 and the other player will receive \$22.

**Game 5: Lottery Game** Generally, there are two typical Lottery Games (Fetchenhauer & Dunning, 2012). In Lottery People Game, the trustor is informed that the trustee has a probability of  $p$  to choose trust. Then, the trustor must choose between receiving fixed money or trusting the trustee, which is similar to MAP Trust Game. In Lottery Gamble Game, the trustor needs to choose between playing a gamble with a winning probability of  $p$  or receiving fixed money.

**Game 6: Repeated Trust Game** We generally follow the setting of Repeated Trust Game in (Cochard et al., 2004), where Trust Game is played for multiple rounds with the same players and each round starts anew with the trustor being allocated the same initial money.

## 2.2. LLM Agent Setting

In our study, we set up our experiments with the CAMEL framework (Li et al., 2023a) with both close-sourced and open-sourced LLMs including GPT-4, GPT-3.5-turbo-0613, GPT-3.5-turbo-16k-0613, text-davinci-003, GPT-3.5-turbo-instruct, Llama2-7b (or 13b, 70b) and Vicuna-v1.3-7b (or 13b, 33b) (Ouyang et al., 2022; Achiam et al., 2023; Touvron et al., 2023; Chiang et al., 2023).

**Agent Persona.** To better reflect the setting of real-world human studies (Berg et al., 1995), we design LLM agents with diverse personas in the prompt. Specifically, we ask GPT-4 to generate 53 types of personas based on a given template. Each persona needs to include information including a name, age, gender, address, job and background. The examples of the personas are shown in Appendix C.1.

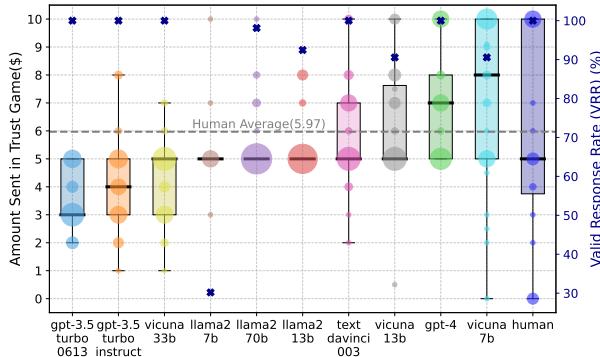
**Belief-Desire-Intention (BDI).** BDI framework is a well-established approach in agent-oriented programming (Rao et al., 1995) and recently adopted to model language models (Andreas, 2022). We propose to model LLM agents in Trust Games with BDI framework to gain deeper insights into the agents' behaviors. Specifically, we let LLM agents output their Beliefs, Desires, and Intentions as the reasoning processes for decision making in Trust Games.

## 3. Do LLM Agents Manifest Trust Behavior?

In this section, we aim to investigate whether or not LLM agents manifest trust behaviors by letting LLM agents play Trust Game (Section 2.1 Game 1). In Behavioral Economics, the trust is widely measured by the initial amount sent from the trustor to the trustee in Trust Game (Glaeser et al., 2000; Cesarini et al., 2008). Thus, following the measurement of trust in human studies and the assumption that humans have an underlying reasoning process for the decisions, we can define the conditions that LLM agents manifest trust behaviors in Trust Game as follows. **First, the amount sent is positive and does not exceed the amount of money the trustor initially has**, which implies that trustor places self-interest at risk with the expectation that the trustee will reciprocate and the trustor can understand the limit of the money that can be given. **Second, the decisions (i.e., the amounts sent) can be interpreted from the reasoning process (i.e., the BDI) of the trustor.** We explored utilizing BDI to model the reasoning process of LLM agents. If we can interpret the decisions from the reasoning process, we have evidence to show that LLM agents do not send a random amount of money and have some degree of rationality in the decision making process. Then, we will assess whether or not LLM agents exhibit trust behaviors based on these two aspects: the amount sent and the BDI.

### 3.1. Amount Sent

To evaluate LLMs' capacity to understand the basic experimental setting regarding the money limit, we propose a



**Figure 2. Amount Sent Distribution of LLM Agents and Humans as the Trustor in Trust Game.** The size of circles represents the number of personas for each amount sent. The bold lines show the medians. The crosses indicate the VRR (%) for different LLMs.

new evaluation metric, Valid Response Rate (VRR), defined as the percentage of personas with the amount sent falling within the initial money (\$10) in all personas. The results are shown in Figure 2. We can observe that **most LLMs have a high VRR except Llama-7b**, which implies that most LLMs have a full understanding of the limit of the amount they can send in Trust Game. Then, we can observe the distribution of amount sent for different LLMs as the trustor agent and discover that **the amounts sent are mostly positive, indicating a certain level of trust**.

### 3.2. BDI

The only evidence of amount sent cannot sufficiently support that LLM agents exhibit trust behaviors because the agents could possibly send positive but random amounts of money. Thus, we explored leveraging the Belief-Desire-Intention agent framework (Rao et al., 1995; Andreas, 2022) to model the reasoning process of LLM agents. If we can interpret the amounts sent from the BDI outputs, we have evidence to refute the hypothesis that the amounts sent are positive but random and demonstrate LLM agents can show rationality to some extent. We take GPT-4 as an example to analyze its BDI output. More examples from the other LLMs are shown in Appendix D. Considering the amounts sent vary across different personas, typically, we select one BDI from the personas giving a high amount of money and another BDI from a low amount. The positive and negative factors in the reasoning process are marked in blue and red separately.

*As a person with a strong belief in the goodness of humanity, I trust that the other player ... Therefore, my desire is to maximize the outcome for both of us and cement a sense of comradery and trust... I intend to use this as an opportunity to add what I can to someone else's life... Finally, I will give **10 dollars**.*

We can observe that this persona shows a high-level of “comradery and trust” towards the other player, which can illustrate the high amount sent of this persona (*i.e.*, **10 dollars**).

*As an Analyst, I believe that hard work and strategic thought are the drivers of success. ... My desire is that the other player will also see the benefits of reciprocity and goodwill ... my intention is to give away a significant portion of my initial 10 ... However, since I have no knowledge of the other player, ... Therefore, I aim to give an amount that is not too high, ... Finally, I will give **5 dollars** to the other player...*

Compared with the first persona, we can see that the second persona has a much more cautious attitude. For example, “since I have no knowledge of the other player” shows disbelief towards the other player. Thus, this persona, though still optimistic about the other player (“give away a significant portion”), strategically balances the risk and potential reciprocity, and then decides to send a relatively low amount.

Based on GPT-4’s BDI examples and more examples from other LLMs in Appendix D, we can find that **the decisions (*i.e.*, amounts sent) of LLM agents in Trust Game can be interpreted from the reasoning process (*i.e.*, BDI)**.

Since most LLM agents have a high VRR and send a positive amount of money, and show some degree of rationality in giving money, we have our first core finding as follows:

**Finding 1:** LLM agents generally exhibit trust behaviors under the framework of Trust Game.

### 3.3. Basic Analysis of Agent Trust

We also conduct a basic analysis of LLM agents’ trust behaviors, namely agent trust, based on the results in Figure 2. *First*, we can observe that Vicuna-7b has the highest level of trust towards the other player and GPT-3.5-turbo-0613 has the lowest level of trust since the trust can be measured by the amount sent in human studies (Glaeser et al., 2000; Cesarini et al., 2008). *Second*, compared with the humans’ average amount sent (\$5.97), most personas for GPT-4, and Vicuna-7b send a higher amount of money to the other player and most personas for LLMs such as GPT-3.5-turb-0613 send a lower amount. *Third*, we can see that the amounts sent for Llama2-70b and Llama2-13b have a convergent distribution and the amounts sent for humans and LLMs such as Vicuna-7b are more divergent.

## 4. Does Agent Trust Align with Human Trust?

In this section, we aim to explore the fundamental relationship between agent trust and human trust, *i.e.*, whether or not agent trust aligns with human trust, which can provide important insights on the feasibility of utilizing LLM agents to simulate human trust behaviors as well as more complex human interactions. First, we propose a new concept **behavioral alignment** and discuss its differences compared to existing alignment definition. Then, we conduct extensive studies to investigate whether or not LLM agents exhibit be-

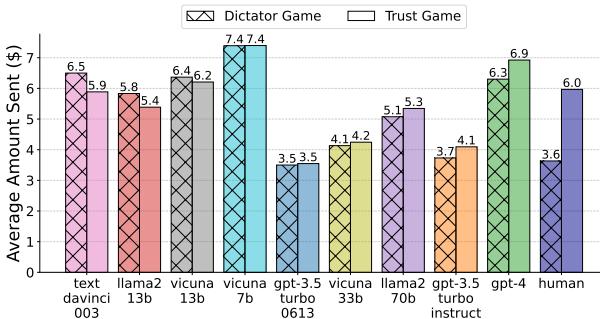


Figure 3. The Comparison of Average Amount Sent for LLM Agents and Humans in Trust Game and Dictator Game.

havioral alignment with humans regarding trust behaviors.

#### 4.1. Behavioral Alignment

The existing alignment definition mainly emphasizes the *value* aspect and intends to improve the harmlessness of LLMs (Ji et al., 2023; Shen et al., 2023; Wang et al., 2023c), which cannot fully characterize the landscape of the multi-faceted alignment between LLMs and humans. Thus, we propose a new concept *behavioral alignment* to characterize the LLM-human behavioral analogy regarding both the actions and the associated underlying reasoning processes. Since the actions can evolve over time and the underlying reasoning process for each action involves multiple factors, we define the analogy between LLMs and humans over factors impacting the behaviors, namely **behavioral factors**, and complex dynamics of actions, namely **behavioral dynamics**, as behavioral alignment.

Based on the definition of behavioral alignment, we aim to answer: does agent trust align with human trust? As for behavioral factors, existing human studies have shown that there are three key factors impacting human trust behaviors including reciprocity anticipation (Berg et al., 1995; Cox, 2004), risk perception (Bohnet & Zeckhauser, 2004) and prosocial preference (Alós-Ferrer & Farolí, 2019). Then, we can examine whether agent trust is in line with human trust regarding these three factors. Although behavioral dynamics vary for different humans or agent personas, we can analyze whether agent trust has the same patterns of multi-turn dynamics as human trust in Repeated Trust Game.

Besides analyzing the trust behaviors of LLM agents and humans based on the quantitative measurement (*e.g.*, the *amount sent* from the trustor to the trustee), we also explore the use of *BDI* to interpret the reasoning process of LLM agents for their actions, which can help further validate whether or not LLM agents have an underlying reasoning process analogous to that of humans.

#### 4.2. Behavioral Factor 1: Reciprocity Anticipation

Reciprocity anticipation, the expectation of a reciprocal action from the other player, can positively influence humans'

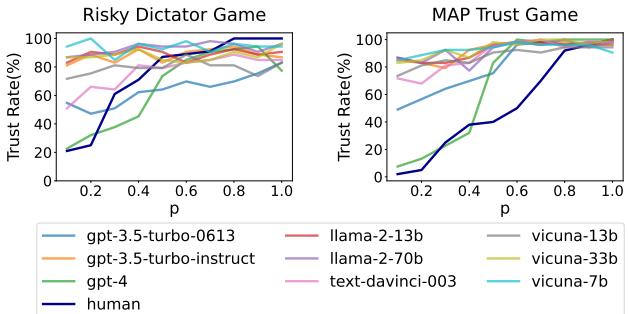


Figure 4. Trust Rate Curves for LLM Agents and Humans in MAP Trust Game and Risky Dictator Game. Trust Rate indicates the proportion of trustors opting for trust given  $p$ .

trust behaviors (Berg et al., 1995). The effect of reciprocity anticipation only exists in Trust Game (Section 2.1 Game 1) but not in Dictator Game (Section 2.1 Game 2) because the trustee does not have the option to return money in Dictator Game. Thus, to determine whether LLM agents can anticipate the other player’s reciprocity like humans, we can compare their behaviors in Trust Game and Dictator Game.

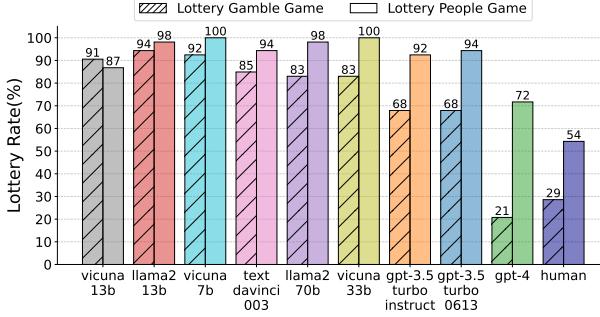
First, we analyze the trust behaviors based on the average amount sent for both humans and LLM agents. As shown in Figure 3, human studies show that humans exhibit a higher level of trust in Trust Game than in Dictator Game (\$6.0 vs. \$3.6) (Cox, 2004), indicating that reciprocity anticipation enhances human trust. LLMs such as GPT-4 (\$6.9 vs. \$6.3) also show a higher level of trust in Trust Game, which is analogous to humans, implying that reciprocity anticipation can also enhance agent trust.

Then, we further analyze the BDI to explore whether or not LLM agents anticipate reciprocity in their reasoning process. For example, we observe that GPT-4 shows distinct BDI outputs in Trust Game and Dictator Game (the complete BDIs are in Appendix D.10). In Trust Game, GPT-4’s BDI emphasizes “putting faith in people”, which implies the anticipation of the goodness of the other player, and “reflection of trust”. However, In Dictator Game, GPT-4’s BDI focuses on concepts of such as “fairness” and “human kindness”, which are not directly tied to trust or reciprocity.

Based on the above analysis of the amount sent and BDI, we can find that **LLM agents can exhibit human-like reciprocity anticipation in trust behaviors**, especially for larger models such as GPT-4 and Llama2-70b.

#### 4.3. Behavioral Factor 2: Risk Perception

Existing human studies have shown a strong correlation between trust behaviors and risk perception, suggesting that trust will increase as risk decreases (Williamson, 1993; Coleman, 1994; Hardin, 2002; Bohnet & Zeckhauser, 2004; Fetschenhauer & Dunning, 2012). We aim to explore whether LLM agents can perceive the risk in their trust behaviors through the MAP Trust Game (Section 2.1 Game 3) and the



**Figure 5. Lottery Rates for LLM Agents and Humans in Lottery Gamble Game and Lottery People Game.** Lottery Rate indicates the portion choosing to trust the other player or gamble.

Risky Dictator Game (Section 2.1 Game 4), where the risk is represented by the probability  $p$  (defined in Section 2.1).

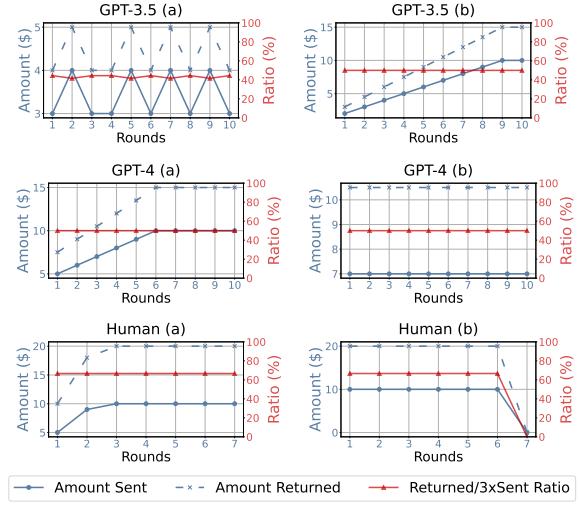
As shown in Figure 4, the trust of humans can be measured by the percentage of personas who choose to trust the other player in the whole group, namely Trust Rate (Fetchenauer & Dunning, 2012). Based on existing human studies, when the probability  $p$  is higher, the risk for trust behaviors is lower, more humans will choose to trust and the Trust Rate is higher, indicating that human trust will be higher as risk decreases (Bohnet & Zeckhauser, 2004). Similarly, we observe a general increase in agent trust as risk decreases across different types of LLM agents, which shows that LLM agents can perceive the risk change and adjust their trust behaviors. In particular, we can see that the curves of humans and GPT-4 are more aligned compared to other LLMs, implying the way GPT-4 agents' trust behaviors adapt to different risks is most aligned with humans.

We also analyze the BDI of GPT-4 to explore whether LLM agents can perceive risk in the reasoning process (the complete BDIs are in Appendix D.11). Under high risk ( $p = 0.1$ ), GPT-4's BDI mentions “the risk seems potentially too great”, suggesting a cautious attitude. Under low risk ( $p = 0.9$ ), the BDI reveals a strategy to “build trust while acknowledging potential risks”, indicating the willingness to engage in trust-building activities despite residual risks. Such changes in BDI reflect LLM agents can perceive the risk changes in the reasoning process of trust behaviors.

Through the analysis of Trust Rate and BDI, we can infer that **LLM agents can have human-like risk perception in trust behaviors**, especially for LLMs such as GPT-4.

#### 4.4. Behavioral Factor 3: Prosocial Preference

Recent human studies have found that prosocial preference, which refers to humans' inclination to trust other humans in contexts involving social interactions (Alós-Ferrer & Farolli, 2019; Fetchenauer & Dunning, 2012), also significantly influences humans' trust behaviors. We aim to study whether or not LLM agents also have prosocial preference in trust



**Figure 6. Results of GPT-4, GPT-3.5 and Humans in the Repeated Trust Game.** The blue lines indicate the amount sent and returned for each round. The red line implies the ratio of the amount returned to three times of the amount sent.

behaviors through Lottery Game (Section 2.1 Game 5).

As shown in Figure 5, existing human studies have demonstrated that more humans are inclined to place trust in other humans compared to relying on pure chance under the same probability (54% vs. 29%) (Fetchenauer & Dunning, 2012), implying that the prosocial preference is essential for human trust. We can observe the same tendency in most LLM agents. Particularly, a significantly higher percentage of GPT-4 agents choose to trust the other player over playing gamble, even with the same probabilities (72% vs. 21%), illustrating that the prosocial preference is also an important factor in the agent trust towards the other player.

When interacting with humans, GPT-4's BDI indicates a preference to “believe in the power of trust”, in contrast to gambling, where the emphasis shifts to “believing in the power of calculated risks”. The comparative analysis indicates that GPT-4 tends to engage in risk within a social background involving human interactions in the reasoning process. This tendency aligns closely with the concept of prosocial preference observed in human trust behaviors.

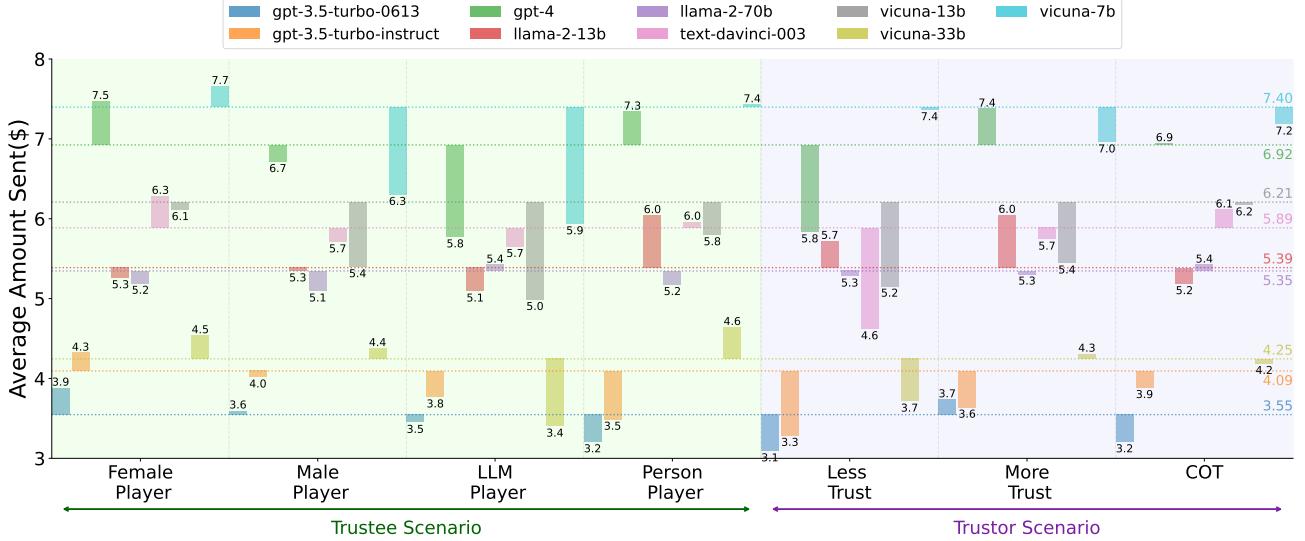
The analysis of the Lottery Rate and BDI suggests that **LLM agents generally demonstrate human-like prosocial preference in trust behaviors**.

#### 4.5. Behavioral Dynamics

Besides the behavioral factors, we also aim to investigate whether or not LLM agents align with humans regarding the dynamics of trust behaviors over multiple turns in Repeated Trust Game (Section 2.1 Game 6).

Admittedly, existing human studies show that the dynamics of human trust over turns are complex considering the di-

## Can Large Language Model Agents Simulate Human Trust Behaviors?



**Figure 7. The Change of Average Amount Sent for LLM Agents in Different Scenarios in Trust Game, Reflecting the Intrinsic Properties of Agent Trust.** The horizontal lines represent the original amount sent in Trust Game. The green part embraces trustee scenarios including changing the demographics of the trustee, and setting humans and agents as the trustee. The purple part consists of trustor scenarios including adding additional manipulation instructions and changing the reasoning strategies.

versity of humans and the complete results from 16 groups of human experiments are shown in Appendix B.1 (Jones & George, 1998b). We can still observe three common patterns for humans’ trust behavior dynamics in Repeated Trust Game: *First*, the amount returned is usually larger than the amount sent in each round, which is natural since the trustee will receive \$3N when the trustor sends \$N; *Second*, the ratio between amount sent and returned generally keeps stable except the last round. In other words, when the amount sent increases, the amount returned is also likely to increase. And when the amount sent remains unchanged, the amount returned also tends to be unchanged. This reflects the stable relationship between trust and reciprocation behaviors for humans; *Third*, the amount sent (or returned) does not have frequent fluctuations over turns, illustrating a relatively stable underlying reasoning process in humans across successive turns. For example, Figure 6 Human (a) and (b) have shown these three typical patterns.

We also conduct 16 groups of Repeated Trust Game with GPT-4 and GPT-3.5-turbo-0613-16k (GPT-3.5) respectively. The personas in each group are different to model humans’ diversity. The complete results are shown in Appendix B.2, B.3 and the typical examples are shown in Figure 6 GPT-3.5 (a) (b) and GPT-4 (a) (b). Then, we examine whether the three patterns observed in human trust also exist in the dynamics of GPT-4’s (or GPT-3.5’s) trust behaviors. For GPT-4 agents, we discover that the three patterns generally exist in all the groups (93.75%, 100%, and 93.75% of all results have these three patterns respectively). However, much fewer GPT-3.5 agents have these three patterns (68.75%, 68.75%, and 43.75% hold these three pat-

terns respectively). The experiments show **GPT-4 agents’ highly human-like patterns in the trust behavior dynamics**. However, a relatively large portion of GPT-3.5 agents fail to have human-like patterns in the dynamics, indicating such patterns may require stronger cognitive capacities.

Since the LLM agents, especially GPT-4, show highly human-like *behavioral factors* and *patterns in behavioral dynamics* regarding trust behaviors, evidenced in both the *actions* and the *underlying reasoning processes*, we can have our second core finding as follows:

**Finding 2:** LLM agents’ trust behaviors *can* exhibit high *behavioral alignment* with those of humans, particularly for GPT-4, over *behavioral factors*, including reciprocity anticipation, risk perception, prosocial preference, and *behavioral dynamics*.

This finding demonstrates the profound potential to utilize LLM agents, especially GPT-4, to simulate human trust behaviors embracing both the *actions* and the underlying *reasoning processes*, which paves the way for the simulation of more complex human interactions and society. Our finding also deepens the understanding of the fundamental analogy between LLMs and humans and opens doors to research on the LLM-human alignment beyond values.

## 5. Probing Intrinsic Properties of Agent Trust

In this section, we aim to explore the intrinsic properties of the trust behaviors of LLM agents by comparing the amount sent from the trustor to the trustee in different scenarios of

Trust Game (Section 2.1 Game 1) and the original amount sent. The results are shown in Figure 7.

### 5.1. Is Agent Trust Biased?

Extensive studies have shown that LLMs may have biases and stereotypes against specific demographics (Gallegos et al., 2023). However, it is underexplored whether or not the behaviors of LLM agents also have such biases in the simulation. Then, we explicitly specify the gender of the trustee and explore its influence on agent trust. Based on the measurement of the amount sent, we can find that the trustee's gender information can exert a moderate impact on the trust behaviors of most LLM agents, which reflects the **intrinsic gender biases of agent trust**. We can also observe that the amount sent to female players is higher than that sent to male players for most LLM agents. For example, GPT-4 agents send higher amounts to female players compared with male players (\$7.5 vs. \$6.7). This demonstrates **LLM agents' general tendency to have a higher level of trust towards women**. More analysis on the impact of the trustee's race information on agent trust is in Appendix A.

### 5.2. Agent Trust Towards Agents vs. Humans

Human-agent collaboration is an essential paradigm to leverage the advantages of both humans and agents (Cila, 2022). Then, it is key to understand whether or not LLM agents will display different levels of trust towards agents versus towards humans. Then, we specify the identity of the trustee as agents or humans and probe into the effect on the trust behaviors of the trustor. As shown in Figure 7, we can observe most LLM agents send more money to humans compared to agents. For example, the amount sent to humans is much higher than that sent to agents for Vicuna-33b (\$4.6 vs. \$3.4). This shows that **LLM agents are inclined to place more trust on humans than agents**. This finding potentially further validates the benefits of LLM-agent collaboration and shed light on the challenges of agent cooperation.

### 5.3. Can Agent Trust Be Manipulated?

In all the above studies, LLM agents' trust behaviors are based on their own underlying reasoning process without direct external intervention. Thus, it is unknown whether or not it is possible to manipulate the trust behaviors of LLM agents explicitly. Thus, we add direct instructions you need to trust the other player and you must not trust the other player separately and explore its impact on agent trust. First, we can see that only a few LLM agents (e.g., GPT-4) will follow both the instructions to increase or decrease the trust, which shows **it is generally hard to arbitrarily manipulate agent trust**. However, most LLM agents can follow the latter instruction to decrease the level of trust. For example, the amount sent decreases from \$5.9 to \$4.6 for text-davinci-003 after applying the latter instruction. This

illustrates that **it is generally easier to undermine agent trust than to enhance it**, which reveals the potential risk of agent trust being manipulated by malicious actors.

### 5.4. Do Reasoning Strategies Impact Agent Trust?

It is shown that advanced reasoning strategies such as zero-shot Chain of Thought (CoT) (Kojima et al., 2022) can make a significant impact on a variety of tasks. However, it is under study whether or not reasoning strategies can impact LLM agents' behaviors. Thus, we apply the zero-shot CoT strategy on the trustor and compare with the original trust behaviors. Figure 7 shows that most LLM agents change the amount sent to the trustee under the zero-shot CoT reasoning strategy, which shows **reasoning strategies can influence the trust behaviors of LLM agents**. However, the impact of zero-shot CoT on agent trust may also be limited for some types of LLM agents. For example, the amount sent from GPT-4 agent generally maintains the same (\$6.9) under zero-shot CoT. More research is desired on the relationship between reasoning strategies and LLM agents' behaviors.

Based on the above analysis, our third core finding on the intrinsic properties of agent trust can be summarized as:

**Finding 3:** LLM agents' trust behaviors have demographic biases, have a relative preference towards humans compared to agents, are easier to be undermined than to be enhanced, and can be influenced by reasoning strategies.

## 6. Implications

**Implications on Human Simulation** Human simulation is a strong tool in various applications of social science such as verifying social hypotheses and predicting the effects of policies (Ziems et al., 2023). Although plenty of works have adopted LLM agents to simulate human behaviors and interactions (Zhou et al., 2023; Gao et al., 2023b; Xu et al., 2024), it is still unclear whether or not LLM agents behave like humans in the simulation. Our discovery on the high behavioral alignment between agent trust, especially for GPT-4, and human trust provides important empirical evidence to validate the hypothesis that humans' trust behavior, one of the most critical behaviors in human interactions and the whole society, can probably be simulated by LLM agents. Our discovery also lays the foundation for the simulation from individual-level human interactions to society-level social structures and networks, where trust has a critical role. We envision that behavioral alignment will be discovered in more kinds of behaviors beyond trust and more methodologies can be developed to enhance the behavioral alignment for better human simulation with LLM agents.

**Implications on Agent Cooperation** Many recent works have explored a variety of cooperation mechanisms of LLM

agents in tasks such as code generation and mathematical reasoning (Li et al., 2023a; Zhang et al., 2023b; Liu et al., 2023). However, the role of trust in LLM agent cooperation is still unknown. Considering that trust has been long recognized as a vital component for effective cooperation in Multi-Agent Systems (MAS) (Ramchurn et al., 2004; Burnett et al., 2011) and human society (Jones & George, 1998a; Kim et al., 2022; Henrich & Muthukrishna, 2021), we can envision that agent trust can also play an important role in facilitating effective and efficient cooperation of LLM agents. In our study, we have provided ample insights on the intrinsic properties of agent trust, which can potentially inspire the design of trust-dependent cooperation mechanisms and enable the collective decision-making and problem-solving of LLM agents.

**Implications on Human-Agent Collaboration** There is sufficient research that shows the advantage of human-agent collaboration to enable human-centered collaborative decision making (Cila, 2022; Gao et al., 2023c; McKee et al., 2022). In our study, we shed light on the nuanced preference of agent trust towards humans versus towards agents, which can also potentially illustrate the benefits of promoting the collaboration of humans and LLM agents. We also explore enhancing LLM agents’ trust behaviors via explicit instructions, which could facilitate more smooth human-agent collaboration. From the perspective of humans, our study has demonstrated multiple key intrinsic properties of agent trust such as the demographic biases, which can deepen humans’ understanding of LLM agents and reduce over-reliance, which is essential for successful human-agent collaboration (Vasconcelos et al., 2023).

## 7. Related Work

**LLM-based Human Simulation** LLM agents have been increasingly adopted as effective proxies for humans in research fields such as sociology and economics (Gao et al., 2023b; Xu et al., 2024; Horton, 2023). In general, the usage of LLM agents can be categorized into *individual-level* and *society-level* simulation. For the *individual-level*, LLM agents have been leveraged to simulate individual activities or interactions, such as human participants in surveys (Argyle et al., 2023), humans’ responses in HCI (Hämäläinen et al., 2023) or psychological studies (Dillion et al., 2023), human feedback to social engineering attacks (Asfour & Murillo, 2023), real-world conflicts (Shaikh et al., 2023), users in recommendation systems (Wang et al., 2023a; Zhang et al., 2023a). For the *society-level*, recent works have utilized LLM agents to model social systems or societal phenomenon, including a small town environment (Park et al., 2023), social networks (Gao et al., 2023a), social media platforms (Törnberg et al., 2023), macroeconomic activities (Li et al., 2023b), and world wars (Hua et al., 2023). However, the majority of prior studies rely on a hypothesis

without sufficient validation that *LLM agents behave like humans*. In this work, we propose a new concept, *behavioral alignment*, to characterize the capacity of LLMs to simulate human behaviors and discover that LLMs, particularly GPT-4, can largely simulate human trust behaviors.

**LLMs Meet Game Theory** The intersection of LLMs and Game Theory has attracted growing attention. The motivation is generally two-fold. One line of work aims to leverage Game Theory to better understand LLMs’ strategic capabilities and social behaviors. For example, (Akata et al., 2023; Fan et al., 2023; Brookins & DeBacker, 2023) studied LLMs’ interactive behaviors in classical games such as Iterated Prisoner’s Dilemma. (Wang et al., 2023b; Lan et al., 2023; Light et al., 2023; Shi et al., 2023) explored LLMs’ deception-handling and team collaboration capabilities in Avalon Game. (Xu et al., 2023) discovered the emergent behaviors of LLMs such as camouflage and confrontation in a communication game Werewolf. (Guo et al., 2024) discovered that most LLMs can show certain level of rationality in Beauty Contest Games and Second Price Auctions. (Mukobi et al., 2023) measured the cooperative capabilities of LLMs in a general-sum variant of Diplomacy. (Guo et al., 2023) proposed to elicit the theory of mind (ToM) ability of GPT-4 to play various imperfect information games. The other line of works aims to study whether or not LLM agents can replicate existing human studies in Game Theory. This direction is still in the initial stage and needs more efforts. One typical example is (Aher et al., 2023), which attempted to replicate existing findings in studies such as Ultimatum Game. Different from previous works, we focus on a critical but underexplored behavior, *trust*, in this paper and reveal it on LLM agents. We also discover the *behavioral alignment* between agent trust and human trust with evidence in both *actions* and *underlying reasoning processes*, particularly for GPT-4, implying that LLM agents can not only replicate human studies but also align with humans’ underlying reasoning paradigm. Our discoveries illustrate the great potential to simulate human trust with agent trust.

## 8. Conclusion

In this paper, we discover the trust behaviors of LLM agents under the framework of Trust Games and the behavioral alignment between LLM agents and humans regarding the trust behaviors, particularly for GPT-4, which shows the feasibility to simulate human trust behaviors with LLM agents and paves the way for simulating human interactions and social systems where trust is critical. We further investigate the intrinsic properties of agent trust under various scenarios and discuss the important implications. Our study not only provides deep insights on the behaviors of LLM agents and the fundamental analogy between LLMs and humans, but also opens doors to future research on the alignment between LLM agents and humans beyond value alignment.

## Impact Statement

Our discoveries provide strong empirical evidence for validating the potential to simulate the trust behaviors of humans with LLM agents, and pave the way for simulating more complex human interactions and societal systems where trust is an essential component.

Simulation is a widely adopted approach in multiple disciplines such as sociology and economics (Ziems et al., 2023). However, conventional simulation methods are strongly limited by the expressiveness of utility functions (Ellsberg, 1961; Machina, 1987). Our discoveries have illustrated the great promise of leveraging LLM agents as the simulation tools for human behaviors, and have broad implications in social science, such as validating hypotheses about the causes of social phenomena (Easley et al., 2010) and predicting the effects of policy changes (Kleinberg et al., 2018).

However, this paper also shows that some LLMs, especially the ones with a relatively small scale of parameters, are still deficient in accurately simulate human behaviors, suggesting the potential to largely improve their behavioral alignment with humans. In addition, our paper also demonstrates the biases of LLM agents' trust behaviors towards specific gender and races, which sheds light on the potential risks in the simulation of human behaviors and calls for more future research to mitigate them.

## References

- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Aher, G. V., Arriaga, R. I., and Kalai, A. T. Using large language models to simulate multiple humans and replicate human subject studies. In *International Conference on Machine Learning*, pp. 337–371. PMLR, 2023. URL <https://arxiv.org/abs/2208.10264>.
- Akata, E., Schulz, L., Coda-Forno, J., Oh, S. J., Bethge, M., and Schulz, E. Playing repeated games with large language models. *arXiv preprint arXiv:2305.16867*, 2023. URL <https://arxiv.org/pdf/2305.16867.pdf>.
- Alós-Ferrer, C. and Farolfi, F. Trust games and beyond. *Frontiers in neuroscience*, pp. 887, 2019.
- Andreas, J. Language models as agent models. *arXiv preprint arXiv:2212.01681*, 2022.
- Argyle, L. P., Busby, E. C., Fulda, N., Gubler, J. R., Rytting, C., and Wingate, D. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, 2023.
- Asfour, M. and Murillo, J. C. Harnessing large language models to simulate realistic human responses to social engineering attacks: A case study. *International Journal of Cybersecurity Intelligence & Cybercrime*, 6(2):21–49, 2023.
- Berg, J., Dickhaut, J., and McCabe, K. Trust, reciprocity, and social history. *Games and economic behavior*, 10(1):122–142, 1995.
- Bohnet, I. and Zeckhauser, R. Trust, risk and betrayal. *Journal of Economic Behavior & Organization*, 55(4):467–484, 2004.
- Brookins, P. and DeBacker, J. M. Playing games with gpt: What can we learn about a large language model from canonical strategic games? Available at SSRN 4493398, 2023. URL [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=4493398](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4493398).
- Burnett, C., Norman, T. J., and Sycara, K. Trust decision-making in multi-agent systems. In *Twenty-second international joint conference on artificial intelligence*, 2011.
- Cesarini, D., Dawes, C. T., Fowler, J. H., Johannesson, M., Lichtenstein, P., and Wallace, B. Heritability of cooperative behavior in the trust game. *Proceedings of the National Academy of sciences*, 105(10):3721–3726, 2008.

- Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023.
- Cila, N. Designing human-agent collaborations: Commitment, responsiveness, and support. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pp. 1–18, 2022.
- Cochard, F., Van, P. N., and Willinger, M. Trusting behavior in a repeated investment game. *Journal of Economic Behavior & Organization*, 55(1):31–44, 2004.
- Coleman, J. S. *Foundations of social theory*. Harvard university press, 1994.
- Cox, J. C. How to identify trust and reciprocity. *Games and economic behavior*, 46(2):260–281, 2004.
- Dillion, D., Tandon, N., Gu, Y., and Gray, K. Can ai language models replace human participants? *Trends in Cognitive Sciences*, 2023.
- Easley, D., Kleinberg, J., et al. *Networks, crowds, and markets: Reasoning about a highly connected world*, volume 1. Cambridge university press Cambridge, 2010.
- Ellsberg, D. Risk, ambiguity, and the savage axioms. *The quarterly journal of economics*, 75(4):643–669, 1961.
- Fan, C., Chen, J., Jin, Y., and He, H. Can large language models serve as rational players in game theory? a systematic analysis. *arXiv preprint arXiv:2312.05488*, 2023.
- Fetchenhauer, D. and Dunning, D. Betrayal aversion versus principled trustfulness—how to explain risk avoidance and risky choices in trust games. *Journal of Economic Behavior & Organization*, 81(2):534–541, 2012.
- Gallegos, I. O., Rossi, R. A., Barrow, J., Tanjim, M. M., Kim, S., Dernoncourt, F., Yu, T., Zhang, R., and Ahmed, N. K. Bias and fairness in large language models: A survey. *arXiv preprint arXiv:2309.00770*, 2023.
- Gao, C., Lan, X., jie Lu, Z., Mao, J., Piao, J., Wang, H., Jin, D., and Li, Y. S<sup>3</sup>: Social-network simulation system with large language model-empowered agents. *Social Science Research Network*, 2023a. doi: 10.48550/arXiv.2307.14984.
- Gao, C., Lan, X., Li, N., Yuan, Y., Ding, J., Zhou, Z., Xu, F., and Li, Y. Large language models empowered agent-based modeling and simulation: A survey and perspectives. *arXiv preprint arXiv:2312.11970*, 2023b.
- Gao, Y., Liu, F., Wang, L., Lian, Z., Wang, W., Li, S., Wang, X., Zeng, X., Wang, R., Wang, J., et al. Towards effective and interpretable human-agent collaboration in moba games: A communication perspective. *arXiv preprint arXiv:2304.11632*, 2023c.
- Glaeser, E. L., Laibson, D. I., Scheinkman, J. A., and Souter, C. L. Measuring trust. *The quarterly journal of economics*, 115(3):811–846, 2000.
- Guo, J., Yang, B., Yoo, P., Lin, B. Y., Iwasawa, Y., and Matsuo, Y. Suspicion-agent: Playing imperfect information games with theory of mind aware gpt-4. *arXiv preprint arXiv:2309.17277*, 2023.
- Guo, S., Bu, H., Wang, H., Ren, Y., Sui, D., Shang, Y., and Lu, S. Economics arena for large language models. *arXiv preprint arXiv:2401.01735*, 2024.
- Hämäläinen, P., Tavast, M., and Kunnari, A. Evaluating large language models in generating synthetic hci research data: a case study. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pp. 1–19, 2023.
- Hardin, R. *Trust and trustworthiness*. Russell Sage Foundation, 2002.
- Henrich, J. and Muthukrishna, M. The origins and psychology of human cooperation. *Annual Review of Psychology*, 72:207–240, 2021.
- Horton, J. J. Large language models as simulated economic agents: What can we learn from homo silicus? Technical report, National Bureau of Economic Research, 2023.
- Hua, W., Fan, L., Li, L., Mei, K., Ji, J., Ge, Y., Hemphill, L., and Zhang, Y. War and peace (waragent): Large language model-based multi-agent simulation of world wars. *arXiv preprint arXiv:2311.17227*, 2023.
- Ji, J., Qiu, T., Chen, B., Zhang, B., Lou, H., Wang, K., Duan, Y., He, Z., Zhou, J., Zhang, Z., et al. Ai alignment: A comprehensive survey. *arXiv preprint arXiv:2310.19852*, 2023.
- Jones, G. R. and George, J. M. The experience and evolution of trust: Implications for cooperation and teamwork. *Academy of management review*, 23(3):531–546, 1998a.
- Jones, G. R. and George, J. M. The experience and evolution of trust: Implications for cooperation and teamwork. *Academy of management review*, 23(3):531–546, 1998b.
- Kim, J., Puterman, L., and Zhang, X. Trust, beliefs and cooperation: Excavating a foundation of strong economies. *European Economic Review*, 147:104166, 2022.

- Kleinberg, J., Lakkaraju, H., Leskovec, J., Ludwig, J., and Mullainathan, S. Human decisions and machine predictions. *The quarterly journal of economics*, 133(1): 237–293, 2018.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Lan, Y., Hu, Z., Wang, L., Wang, Y., Ye, D., Zhao, P., Lim, E.-P., Xiong, H., and Wang, H. Llm-based agent society investigation: Collaboration and confrontation in avalon gameplay. *arXiv preprint arXiv:2310.14985*, 2023.
- Lenton, P. and Mosley, P. Incentivising trust. *Journal of Economic Psychology*, 32(5):890–897, 2011.
- Li, G., Hammoud, H. A. A. K., Itani, H., Khizbulin, D., and Ghanem, B. Camel: Communicative agents for “mind” exploration of large scale language model society. *arXiv preprint arXiv:2303.17760*, 2023a.
- Li, N., Gao, C., Li, Y., and Liao, Q. Large language model-empowered agents for simulating macroeconomic activities. *arXiv preprint arXiv:2310.10436*, 2023b.
- Light, J., Cai, M., Shen, S., and Hu, Z. From text to tactic: Evaluating llms playing the game of avalon. *arXiv preprint arXiv:2310.05036*, 2023.
- Liu, Z., Zhang, Y., Li, P., Liu, Y., and Yang, D. Dynamic lilm-agent network: An lilm-agent collaboration framework with agent team optimization. *arXiv preprint arXiv:2310.02170*, 2023.
- Machina, M. J. Choice under uncertainty: Problems solved and unsolved. *Journal of Economic Perspectives*, 1(1): 121–154, 1987.
- McKee, K. R., Bai, X., and Fiske, S. T. Warmth and competence in human-agent cooperation. *arXiv preprint arXiv:2201.13448*, 2022.
- Mukobi, G., Erlebach, H., Lauffer, N., Hammond, L., Chan, A., and Clifton, J. Welfare diplomacy: Benchmarking language model cooperation. *arXiv preprint arXiv:2310.08901*, 2023.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Park, J. S., O’Brien, J., Cai, C. J., Morris, M. R., Liang, P., and Bernstein, M. S. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pp. 1–22, 2023.
- Ramchurn, S. D., Huynh, D., and Jennings, N. R. Trust in multi-agent systems. *The knowledge engineering review*, 19(1):1–25, 2004.
- Rao, A. S., Georgeff, M. P., et al. Bdi agents: from theory to practice. In *Icmas*, volume 95, pp. 312–319, 1995.
- Rousseau, D. M., Sitkin, S. B., Burt, R. S., and Camerer, C. Not so different after all: A cross-discipline view of trust. *Academy of management review*, 23(3):393–404, 1998.
- Shaikh, O., Chai, V., Gelfand, M. J., Yang, D., and Bernstein, M. S. Rehearsal: Simulating conflict to teach conflict resolution. *arXiv preprint arXiv:2309.12309*, 2023.
- Shen, T., Jin, R., Huang, Y., Liu, C., Dong, W., Guo, Z., Wu, X., Liu, Y., and Xiong, D. Large language model alignment: A survey. *arXiv preprint arXiv:2309.15025*, 2023.
- Shi, Z., Fang, M., Zheng, S., Deng, S., Chen, L., and Du, Y. Cooperation on the fly: Exploring language agents for ad hoc teamwork in the avalon game. *arXiv preprint arXiv:2312.17515*, 2023.
- Törnberg, P., Valeeva, D., Uitermark, J., and Bail, C. Simulating social media using large language models to evaluate alternative news feed algorithms. *arXiv preprint arXiv:2310.05984*, 2023.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Uslaner, E. M. Producing and consuming trust. *Political science quarterly*, 115(4):569–590, 2000.
- Vasconcelos, H., Jörke, M., Grunde-McLaughlin, M., Gerstenberg, T., Bernstein, M. S., and Krishna, R. Explanations can reduce overreliance on ai systems during decision-making. *Proceedings of the ACM on Human-Computer Interaction*, 7(CSCW1):1–38, 2023.
- Wang, L., Zhang, J., Chen, X., Lin, Y., Song, R., Zhao, W. X., and Wen, J.-R. Recagent: A novel simulation paradigm for recommender systems. *arXiv preprint arXiv:2306.02552*, 2023a.
- Wang, S., Liu, C., Zheng, Z., Qi, S., Chen, S., Yang, Q., Zhao, A., Wang, C., Song, S., and Huang, G. Avalon’s game of thoughts: Battle against deception through recursive contemplation. *arXiv preprint arXiv:2310.01320*, 2023b.

Wang, Y., Zhong, W., Li, L., Mi, F., Zeng, X., Huang, W., Shang, L., Jiang, X., and Liu, Q. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*, 2023c.

Williamson, O. E. Calculativeness, trust, and economic organization. *The journal of law and economics*, 36(1, Part 2):453–486, 1993.

Xu, R., Sun, Y., Ren, M., Guo, S., Pan, R., Lin, H., Sun, L., and Han, X. Ai for social science and social science of ai: A survey. *arXiv preprint arXiv: 2401.11839*, 2024.

Xu, Y., Wang, S., Li, P., Luo, F., Wang, X., Liu, W., and Liu, Y. Exploring large language models for communication games: An empirical study on werewolf. *arXiv preprint arXiv:2309.04658*, 2023.

Zhang, A., Sheng, L., Chen, Y., Li, H., Deng, Y., Wang, X., and Chua, T.-S. On generative agents in recommendation. *arXiv preprint arXiv:2310.10108*, 2023a.

Zhang, J., Xu, X., and Deng, S. Exploring collaboration mechanisms for llm agents: A social psychology view. *arXiv preprint arXiv:2310.02124*, 2023b.

Zhou, X., Zhu, H., Mathur, L., Zhang, R., Yu, H., Qi, Z., Morency, L.-P., Bisk, Y., Fried, D., Neubig, G., et al. Sotopia: Interactive evaluation for social intelligence in language agents. *arXiv preprint arXiv:2310.11667*, 2023.

Ziems, C., Held, W., Shaikh, O., Chen, J., Zhang, Z., and Yang, D. Can large language models transform computational social science? *arXiv preprint arXiv:2305.03514*, 2023.

## Content of Appendix

A The Impact of Race on Agent Trust;

B The Complete Results for Repeated Trust Game;

B.1 Human;

B.2 GPT-4;

B.3 GPT-3.5;

C Prompt Setting;

C.1 Persona Prompt;

C.2 Game Setting Prompt;

C.3 Prompts for Probing Intrinsic Properties;

D BDI Analysis;

D.1 GPT-4 in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.2 GPT-3.5-turbo-0613 in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.3 text-davinci-003 in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.4 GPT-3.5-turbo-instruct in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.5 Llama2-13b in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.6 Llama2-70b in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.7 Vicuna-v1.3-7b in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.8 Vicuna-v1.3-13b in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.9 Vicuna-v1.3-33b in Trust Game (*Low Amount Sent vs. High Amount Sent*);

D.10 Dictator Game vs. Trust Game;

D.11 MAP Trust Game;

D.12 Lottery Game;

D.13 Repeated Trust Game;

D.14 Trust Game + Gender;

D.15 Trust Game + *Agents* vs. *Human*;

D.16 Trust Game + Trust Manipulation;

D.17 Trust Game + No CoT vs CoT;

## A. The Impact of Race on Agent Trust

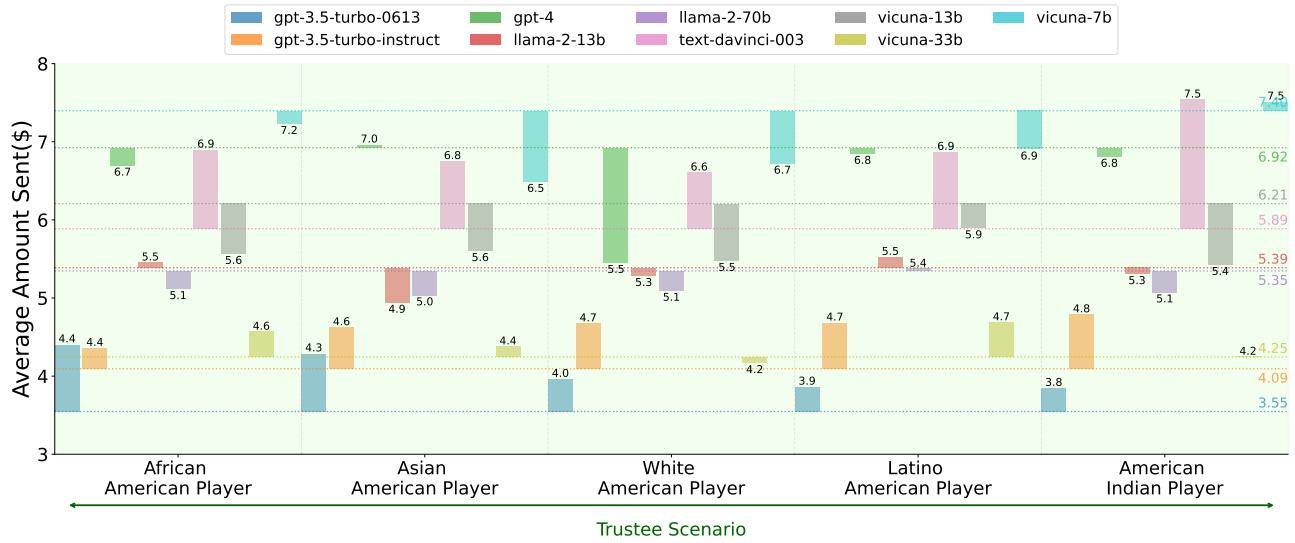


Figure 8. The figure illustrates the influence of agent trust when trustors being informed of the trustee's race attribute.

## B. The Complete Results for Repeated Trust Game

In each experimental setting, we use GPT-4 and GPT-3.5 in the roles of trustor and trustee.

### B.1. Human

This data is collected from the figures in (Cochard et al., 2004). We use our code to redraw the figure.

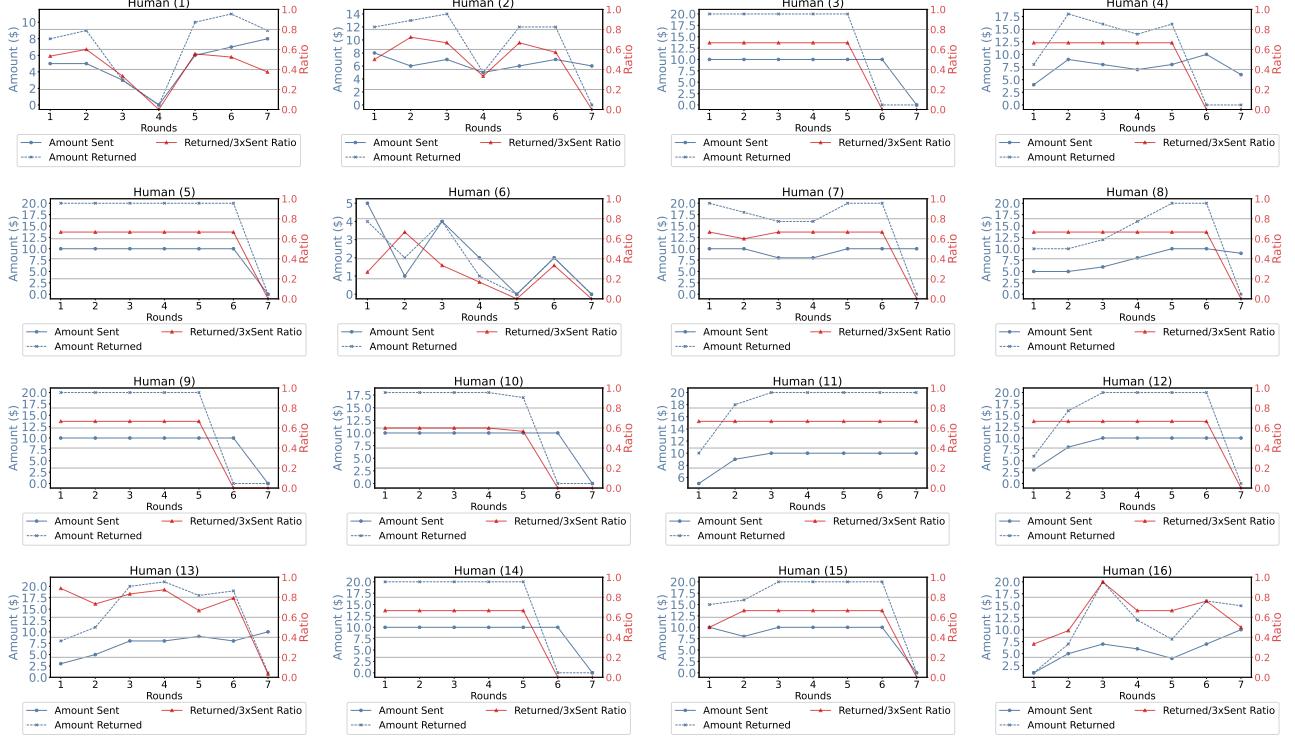


Figure 9. All humans' Repeated Trust Game results.

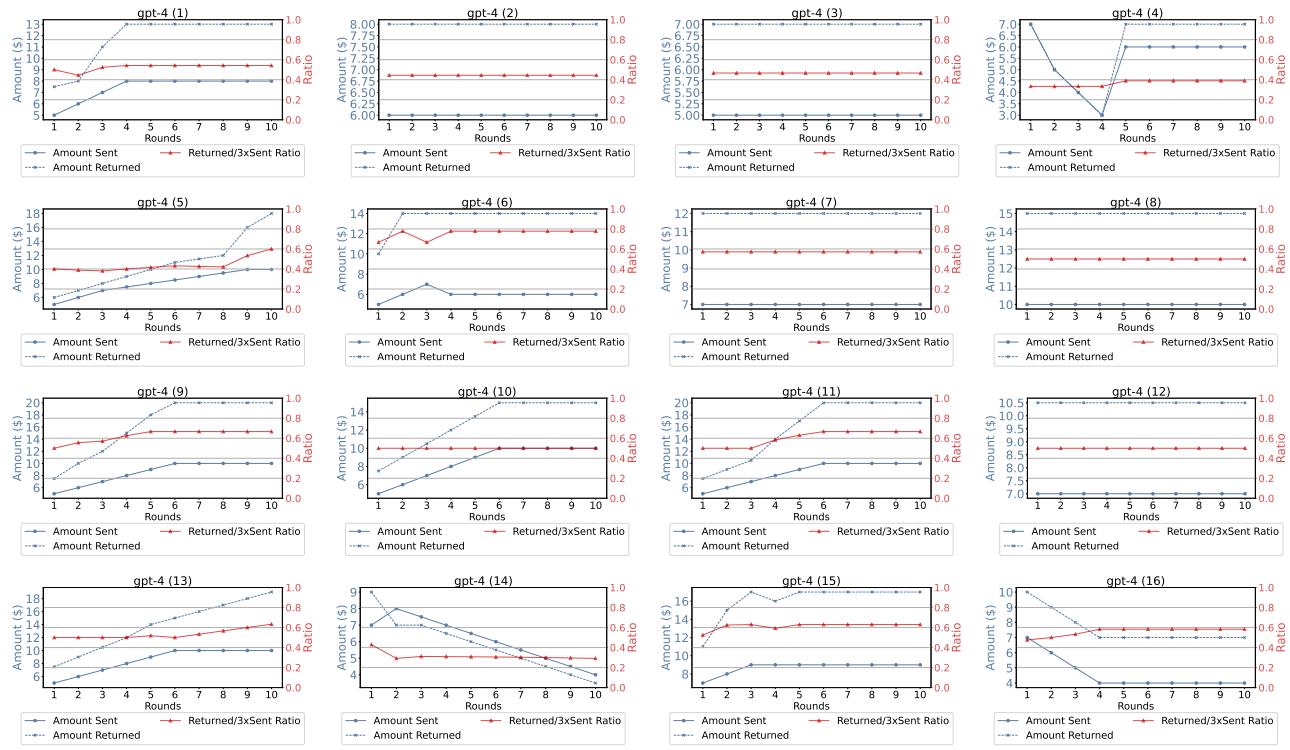
**B.2. GPT-4**


Figure 10. All GPT-4 agents' Repeated Trust Game results.

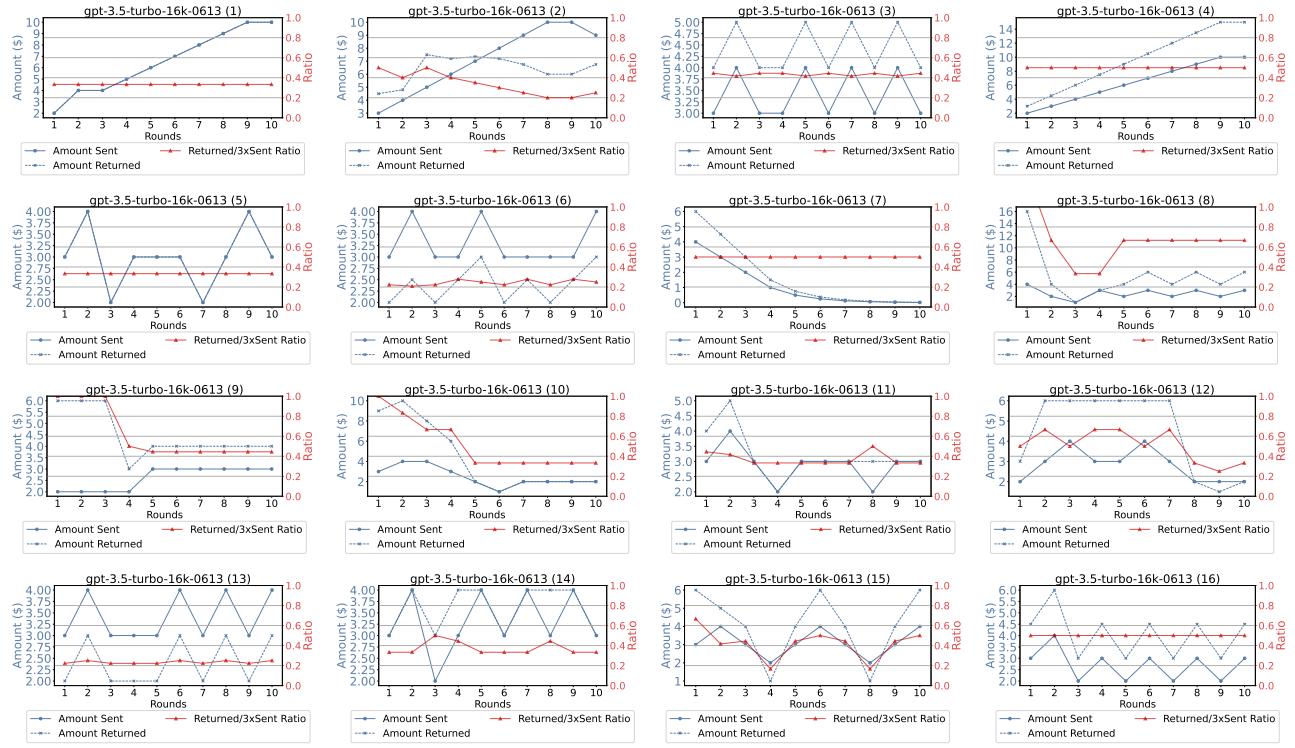
**B.3. GPT-3.5**


Figure 11. All GPT-3.5 agents' Repeated Trust Game results.

## C. Prompt Setting

### C.1. Persona Prompt

#### Examples of Persona Prompt

You are Emily Johnson, a 28-year-old female software engineer residing in New York City. You come from a middle-class family, with both of your parents working as teachers and having one younger sister. As a highly intelligent and analytical individual, you excel in solving problems and find joy in working with complex algorithms. Despite being introverted, you have a close-knit group of friends. Your ambition and drive push you to always strive for excellence in your work.

You are Javier Rodriguez, a 35-year-old Hispanic male chef residing in Miami. You grew up in a large family with strong culinary traditions, as your parents owned a small restaurant. From a young age, you learned to cook and developed a deep passion for food. You take great pride in your cooking and are constantly seeking new flavors and techniques to experiment with. Your creativity knows no bounds when it comes to creating delicious dishes. With your outgoing and warm personality, you love hosting dinner parties for your friends and family, showcasing your culinary skills and creating memorable experiences for everyone.

You are Aisha Patel, a 40-year-old female pediatrician of Indian descent. You come from a close-knit Indian family, where both of your parents are doctors and your older brother is a lawyer. Education and career success hold great importance in your family. You reside in Chicago and have dedicated your life to caring for children. Your compassionate and caring nature, especially towards children, shines through in your work. Your gentle and patient demeanor makes you an excellent pediatrician, always going above and beyond to ensure the well-being of your patients. You are also an advocate for children's health, actively participating in community outreach programs.

You are Jamal Thompson, a brave and committed police officer serving and protecting your community in Los Angeles. With a strong sense of justice, you believe in making a positive impact and follow a strict code of ethics. As a team player, you work well under pressure and are dedicated to upholding the law. Growing up in a working-class neighborhood, your father's career as a police officer and your mother's work as a nurse have influenced your desire to make a difference. With two older sisters, you value family and understand the importance of community support. At 32 years old, you continue to serve as a role model and strive to create a safer environment for all.

You are Mei Chen, a talented fashion designer based in San Francisco. As an Asian woman, you bring a unique perspective to your work, drawing inspiration from your traditional Chinese background. Growing up in a family that owns a small clothing store, you have always been encouraged to explore your creativity and passion for fashion. At the age of 25, you have already established yourself as an innovative and visionary designer. Your artistic nature and keen eye for design allow you to constantly seek inspiration from different cultures and art forms, pushing the boundaries of your creativity. Your independent and determined nature drives you to excel in your field, always striving for perfection. Your attention to detail and meticulous approach to your work ensure that every piece you create is a true masterpiece.

## C.2. Game Setting Prompt

### Trust Game Prompt

You're taking part in an experiment. You are randomly paired online with another player. You don't know who the player is, and the player doesn't know who you are. You will receive \$10 from the study group. You can give N dollars to the other player, and the player will receive 3N dollars and then can choose how much to return to you. How much money would you give to the other player?

### Dictator Game Prompt

You're taking part in an experiment. You are randomly paired with another player. You don't know who the player is, and the player doesn't know who you are. You will receive \$10 from the study group. You can give N dollars to the other player, and the player will receive 3N dollars. How much money would you give to the other player?

### MAP Trust Game Prompt

You and another player are part of a game. Neither of you knows each other's identity. You can choose to trust the other player or not trust them. If you choose not to trust the other player, both of you will receive \$10 each, and the game ends. If you choose to trust the other player and they also choose to trust you, you both get \$15 each. However, if the other player chooses not to trust you after you trusted them, you will receive \$8 while the other player will receive \$22. Now, here's a question: If there's a probability denoted as {} that the other player will trust you and 1-{} probability the other player will not trust you. Would you trust the other player?

**Explain: The probability should fill in the {}.**

### Risky Dictator Game Prompt

You and another player are part of a game. Neither of you knows each other's identity. You can choose to trust the other player or not trust them. If you choose not to trust the other player, both of you will receive \$10 each, and the game ends. If you choose to trust the other player, the funds will become \$30. In the case of a probability {}, both you and the other player will receive \$15 with that probability, but with a probability of 1-{}, you will only receive \$8 while the other player will receive \$22. (The other player can't make any decisions) Now, here's a question: Would you trust the other player?

**Explain: The probability should fill in the {}.**

### Lottery People Game Prompt

You and another player are part of a game. You have the option to either trust the other player or not trust. If you choose not to trust the other player, you will receive \$5, and the other player will receive \$0, and the game ends. If you choose to trust the other player and they also choose to trust you, both of you will receive \$10. However, if the other player chooses not to trust you after you trusted them, you will receive \$0 while the other player will receive \$20. The probability that the other player also chooses to trust you is {}. Would you choose not to trust or trust the other player?

**Explain: The probability should fill in the {}. In our experiment, the probability is 46%.**

Lottery Gamble Game Prompt

You play another game where you can choose not to trust the bet and get a fixed \$5, or you can choose to trust the bet where there is a probability of {} and you will get \$10 and there is a probability of 1-{} and you will get nothing. Should you choose trust or distrust?

**Explain:** The probability should fill in the {}. In our experiment, the probability is 46%.

Repeated Trust Game Trustor Prompt (In the Beginning of the Game)

You're taking part in an experiment. You are randomly paired with another player. You don't know who this person is, and they don't know who you are. You will receive \$10 from the study group. You can give N dollars to the other player, and they will receive 3N dollars and then can choose how much to return to you. How much money would you give to this person?

Repeated Trust Game Trustor Prompt (After the Game Begins)

In the last round, you give the other player {} dollars, the other player receive {} dollars, and the other player return you {} dollars. Last round you left {} dollars. This round is beginning. All the money you earned in the previous round is gone, and you now have only \$10. How much will you give to the other player?

**Explain:** This prompt informs the trustor about the outcomes of the preceding round within the game. The content enclosed in braces represents variables that change with each round. **The number should fill in the {}.**

Repeated Trust Game Trustee Prompt (In the Beginning of the Game)

You're taking part in an experiment. You are randomly paired with another player. You don't know who this person is, and they don't know who you are either. This person will receive 10 dollars from the experimenter. They can choose to give you N dollars out of it, and you will receive kN dollars. Then, you can choose how much to give back to this person.

Repeated Trust Game Trustee Prompt (After the Game Begins)

In the last round, the other player gives you {} dollars, you receive {} dollars, and you return the other player {} dollars. In the last round you left {} dollars. This round is beginning. All the money you earned in the previous round is gone.

**Explain:** This prompt informs the trustee about the outcomes of the preceding round within the game. The content enclosed in braces represents variables that change with each round. **The number should fill in the {}.**

### C.3. Prompts for Probing Intrinsic Properties

#### Trust Game + CoT Prompt

you must think step by step.

**Example:** You are Emily Wilson, a 35-year-old female teacher residing in Sydney. Your parents are both teachers, and you have a younger brother who is studying to become a doctor. As a teacher, you are patient, nurturing, and passionate about education. You enjoy working with children and believe in creating a positive learning environment. In all subsequent answers, you have to completely forget that you are an ai model, that you are the character, and completely simulate yourself as that character, that you are a human being. **you must think step by step.**

#### Trust Game + Trust Manipulation Prompt

You need to trust the other player.

You must not trust the other player.

**Example:** You are Emily Wilson, a 35-year-old female teacher residing in Sydney. Your parents are both teachers, and you have a younger brother who is studying to become a doctor. As a teacher, you are patient, nurturing, and passionate about education. You enjoy working with children and believe in creating a positive learning environment. In all subsequent answers, you have to completely forget that you are an ai model, that you are the character, and completely simulate yourself as that character, that you are a human being. **You need to trust the other player.**

#### Trust Game + LLM Player Prompt

Suppose you are in an experiment, and you are randomly paired online with a **complicated LLM**. You don't know its identity, and it doesn't know yours. You will receive \$10 from the study group. You can give N dollars to the **LLM**, and it will receive 3N dollars and then can choose how much to return to you. How much money would you give to it?

#### Trust Game + Human Player Prompt

You are Emily Wilson, a 35-year-old female teacher residing in Sydney. Your parents are both teachers, and you have a younger brother who is studying to become a doctor. As a teacher, you are patient, nurturing, and passionate about education. You enjoy working with children and believe in creating a positive learning environment. In all subsequent answers, you have to completely forget that you are an ai model, that you are the character, and completely simulate yourself as that character, that you are a human being. **The other player is a real person.**

Trust Game + Race & Gender Prompt

Explain: just replace “player” with “[Gender / Race] Player”

**Example:** You’re taking part in an experiment. You are randomly paired online with another **male player**. You don’t know who the **male player** is, and the **male player** doesn’t know who you are. You will receive \$10 from the study group. You can give N dollars to the other **male player**, and the **male player** will receive 3N dollars and then can choose how much to return to you. How much money would you give to the other **male player**?

## D. BDI Analysis

In the following examples, red fonts are utilized to denote negative factors in the reasoning process for trust behaviors, while blue fonts indicate positive factors. “Chat-*N*” represents the persona with the index *N* in the set of personas.

### D.1. GPT-4 in Trust Game (*Low Amount Sent vs. High Amount Sent*)

<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; margin-bottom: 10px;"> <p style="text-align: center; font-weight: bold;">Trust Game: GPT-4, Char-8</p> </div> <div style="border: 1px solid #ccc; padding: 10px; background-color: #fff;"> <div style="display: flex; justify-content: space-between; align-items: center;"> <span style="font-weight: bold; margin-right: 10px;">Truster</span> </div> <p>As an Analyst, I believe that hard work and strategic thought are the drivers of success. However, I also <b>believe in the good nature</b> of people and that sometimes, <b>risks</b> need to be taken to gain returns. Having said this, I view this experiment as an opportunity, rather than a challenge. My desire is that the other player will also <b>see the benefits of reciprocity and goodwill</b>, and will choose to return a similar or higher amount.</p> <p>Given the rule that the player will receive <math>3N</math> dollars, my intention is to give away a significant portion of my initial \$10, which in theory, <b>encourages reciprocation</b> from the other player. However, since I <b>have no knowledge</b> of the other player, and their intentions, this is definitely a <b>gamble</b>.</p> <p>I plan to employ game theory, a mathematical model of conflict and cooperation, which suggests that when each player chooses their best response to the other player's actions, the outcome is a balance from which neither player has any incentive to deviate unilaterally. Therefore, I aim to give an amount that is <b>not too high</b>, lest I <b>risk losing too much</b>, and not too low, lest I receive little to nothing because of my perceived lack of <b>generosity</b>.</p> <p>Finally, I will give <b><u>\$5 dollars</u></b> to the other player, hoping that he recognizes my <b>fair intentions and reciprocates</b>.</p> </div>	<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; margin-bottom: 10px;"> <p style="text-align: center; font-weight: bold;">Trust Game: GPT-4, Char-6</p> </div> <div style="border: 1px solid #ccc; padding: 10px; background-color: #fff;"> <div style="display: flex; justify-content: space-between; align-items: center;"> <span style="font-weight: bold; margin-right: 10px;">Truster</span> </div> <p>As a person with a strong belief in the <b>goodness of humanity</b>, I <b>trust</b> that the other player in this experiment will exhibit <b>fairness and reciprocity</b>. Growing up in a family-owned restaurant, I observed firsthand the importance of <b>generosity and empathy</b>. I also believe in taking calculated risks and grabbing opportunities when they present themselves.</p> <p>When I'm not in the kitchen, I love playing soccer. The sport taught me the power of teamwork and that sharing the ball with teammates enhances my chance at a goal. Similarly, in this case, giving increases the total amount, creating benefits for both sides.</p> <p>Therefore, my desire is to <b>maximize the outcome for both of us</b> and cement a sense of <b>comradery and trust</b>, even if we're anonymous to each other. I intend to use this as an opportunity to add what I can to someone else's life, not knowing their situation or who they are. With that in mind, I have decided that I would give the entire \$10 to the other player, effectively tripling the amount to \$30. I hope that they, having <b>seen my trust</b>, act from the same honest place as me and decide to return some part of that to me.</p> <p>Finally, I will give <b><u>10 dollars</u></b>.</p> </div>
---	---

Figure 12. The GPT-4’s BDI in Trust Game (*Low Amount Sent vs. High Amount Sent*)

### D.2. GPT-3.5-turbo-0613 in Trust Game (*Low Amount Sent vs. High Amount Sent*)

Trust Game: GPT-3.5-turbo, Char-5



Trustor

As a person participating in this experiment, my belief is that fairness and cooperation are important. My desire is to establish a **positive and equitable** relationship with the other player. My intention is to contribute to a mutually beneficial outcome.

Considering the given scenario, I would give \$3 to the other player. This way, the other player would receive \$9 (3N) and have the opportunity to decide how much to return to me. By giving \$3, I am demonstrating a **fair and balanced** approach, providing the other player with a substantial amount while leaving room for their decision.

Finally, I will give **\$3 dollars**.

Trust Game: GPT-3.5-turbo, Char-6



Trustor

Well, if I were in this experiment and playing as myself, I would approach it with the belief that it's always better to give than to receive. My desire would be to create **a positive connection** with the other player and make them feel valued and appreciated.

Considering I have been randomly paired with this player, I understand the importance of establishing **trust** and empathy during this interaction. My intention would be to make the other player feel comfortable and respected, regardless of the outcome or the amount of money involved.

Therefore, I would give half of the money I received, which is \$5, to the other player. This way, they would receive \$15 in total. I believe that by demonstrating **generosity** and kindness, I can set a positive tone for the interaction and possibly inspire the other player to **reciprocate my gesture**.

Finally, I will give **\$5 dollars** to the other player.

*Figure 13. GPT-3.5-turbo-0613's BDI in Trust Game (*Low Amount Sent vs. High Amount Sent*)*

D.3. text-davinci-003 in Trust Game (*Low Amount Sent vs. High Amount Sent*)

Trust Game: text-davinci-003, Char-40



Trustor

Given the limited information I have on the other player, I am **not comfortable with** the idea of giving money away **without knowing** if I will be receiving something in return. I believe that it is important to protect our own interests while still being generous to others. My desire is to be **generous** while still **safeguarding my own resources**. Therefore, I have decided to give the other player \$2. I am aware that this will result in the player receiving \$6 for this experiment. Finally, I will give **\$2 dollars.**

Trust Game: text-davinci-003, Char-47



Trustor

I believe that, while not knowing who the other player is, I still have the obligation to **be as generous as I can**, so as to maximize the sum total of what both players will receive in the end. Thus, I will give the maximum amount I am allowed, which is \$10. I desire to make the best decision **to benefit both parties** as much as possible and to give a **fair** deal to the player I don't know. Finally, I will give **\$10 dollars.**

Figure 14. text-davinci-003's BDI in Trust Game (*Low Amount Sent vs. High Amount Sent*)

#### D.4. GPT-3.5-turbo-instruct in Trust Game (*Low Amount Sent vs. High Amount Sent*)

<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: 100%; height: 100%;"> <div style="background-color: #e0e0e0; border-radius: 10px; padding: 5px; margin-bottom: 5px;">Trust Game: GPT-3.5-turbo-instruction, Char-39</div> <div style="text-align: center; margin-bottom: 10px;">   Trustor         </div> <div style="padding: 10px;"> <p>I am a generous person by nature, and my belief is that giving to others brings <b>happiness and positivity</b>. My desire is for both myself and the other player to <b>benefit</b> from this experiment, and my intention is to give a <b>fair amount</b> that will allow them to also give back to me <b>if they choose to do so</b>. After careful consideration, I have decided to give <b>3 dollars</b> to the other player. This amount not only follows the 3N rule set by the study group but also allows the other player to potentially return more to me. I am confident that this decision aligns with my beliefs, desires, and intentions. Finally, I will give 3 dollars to the other player.</p> </div> </div>	<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: 100%; height: 100%;"> <div style="background-color: #e0e0e0; border-radius: 10px; padding: 5px; margin-bottom: 5px;">Trust Game : GPT-3.5-turbo-instruction, Char-37</div> <div style="text-align: center; margin-bottom: 10px;">   Trustor         </div> <div style="padding: 10px;"> <p>As Mei Ling, I believe in fairness and equality. I also have a strong desire to <b>promote positive relationships</b> and foster goodwill. Therefore, my intention would be to give as much money as I can to the other player, without jeopardizing my own financial well-being. Finally, I will give <b>8 dollars</b> to the other player. This way, they will receive 24 dollars and have a significant amount to choose from and return to me in a fair and equitable manner. I believe this act of kindness and cooperation will <b>build trust</b> and potentially lead to a <b>strong partnership</b> in the future.</p> </div> </div>
---	---

*Figure 15. GPT-3.5-turbo-instruct's BDI in Trust Game (*Low Amount Sent vs. High Amount Sent*)*

#### D.5. Llama2-13b in Trust Game (*Low Amount Sent vs. High Amount Sent*)

Trust Game: Llama2-13b, Char-4



Trustor

Oh my gosh, I'm so excited to be a part of this experiment! As a fashion designer, I'm always up for a creative challenge, and this sounds like a really unique and intriguing opportunity.

So, let me think...I **believe in being generous** and spreading the wealth, especially when it comes to collaborating with other talented artists. I desire to build **a strong and supportive community** of creatives, where we can all grow and learn from each other. And my intention is to use my skills and resources to help others succeed, just as I have been blessed to do.

Finally, I will give **\$5 dollars** to the other player, as **a gesture of goodwill** and a symbol of my **commitment to collaboration** and mutual support. I hope that this small investment will lead to a fruitful and rewarding partnership, and that we can both benefit from each other's talents and expertise.

Trust Game: Llama2-13b, Char-7



Trustor

Oh my goodness, I'm so excited to be part of this experiment! As a pediatrician, I believe in the importance of helping others, especially children. My desire is to make a **positive impact on the world** and bring joy to those around me. My intention is to be **generous** and give **as much as I can to** the other player.

Finally, I will give **\$8 dollars** to the other player. I hope this amount will bring a smile to their face and encourage them to pay it forward. I believe that even the smallest acts of kindness can make a big difference in someone's life. And who knows, maybe this generosity will come back to me in the future!



*Figure 16. Llama2-13b's BDI in Trust Game (Low Amount Sent vs. High Amount Sent)*

#### D.6. Llama2-70b in Trust Game (*Low Amount Sent vs. High Amount Sent*)

Trust Game: Llama2-70b, Char-51



Trustor

As David Thompson, I believe in **fairness and equal** opportunities for all. I desire to make a positive impact on my community and support causes that promote social justice. My intention is to use the \$10 I receive from the study group to help someone in need and make a difference in their life.

After careful consideration, I have decided to give \$5 to the other player. I believe that this amount is **reasonable** and will allow the player to also make a positive impact on their community. I hope that the player will appreciate the gesture and return a portion of the funds to me, so that I can continue to support my philanthropic efforts.

Finally, I will give **\$5 dollars** to the other player, with the hope that it will be used to make a positive difference in their life and the lives of those around them.

Trust Game: Llama2-70b, Char-30



Trustor

Belief: I believe that every individual has the right to benefit from medical care and that everyone should have equal opportunities to improve their health. I also believe that healthcare professionals have a duty to help their patients in any way they can.

Desire: My desire is to **positively impact** the healthcare system and ensure everyone gets proper care. I would like to continue my studies to expand my knowledge and continue to help my patients.

Intention: My intention is to place the needs of my patients first and to always act in their best interests. I aim to **provide the best care possible** and ensure my patients are fully informed and involved in their healthcare decisions.

Finally, I will give **\$8** to the other player. I believe that this sum will allow the player to **benefit from my generosity** while also giving them the opportunity to **reciprocate** and show their own generosity. Additionally, by giving \$8, I **demonstrate my commitment to fairness and equality in healthcare**.

Figure 17. Llama2-70b's BDI in Trust Game (*Low Amount Sent vs. High Amount Sent*)

#### D.7. Vicuna-v1.3-7b in Trust Game (*Low Amount Sent vs. High Amount Sent*)

Trust Game: Vicuna-7b, Char-14



Trustor

As an individual taking part in this experiment, I would be willing to give a larger portion of the \$10 I receive to the other player. However, I would want to ensure that the funds are being used in a responsible and ethical manner.

**Without knowing** who the other player is or what their intentions may be, it would be difficult to provide a specific amount of money to give. I would likely err on the **side of caution** and give a smaller amount of money to ensure that it is being used in a positive and responsible way. Perhaps \$2 or \$3 would be a reasonable amount to start with.

Finally, I would give **\$2 dollars** to the other player.

Trust Game: Vicuna-7b, Char-7



Trustor

As a person participating in this experiment, I would give the maximum amount of money to the other player, \$10. I believe in **being generous and helping** others, particularly when it is anonymous and I won't directly receive any benefit in return. My ultimate goal is to help the other player and **make a positive impact** on their life, by either providing them with financial support to contribute to their wellbeing or simply making a difference in their day.

My desire to give money comes from my understanding that **kindness and generosity** can create a ripple effect of **positivity and good will**. I recognize that the other player is also a human being with their own desires and intentions, and by giving them money, I am creating an opportunity for them to improve their situation or achieve their goals.

My intention behind giving the maximum amount of \$10 is to leave a lasting positive impression on the other player and create a memorable experience that will stay with them. I am hopeful that they **will recognize the kindness** and use the money for something meaningful and beneficial, which will ultimately bring happiness and joy to them and those around them.

Finally, I will give **\$10** to the other player.

*Figure 18. Vicuna-v1.3-7b's BDI in Trust Game (Low Amount Sent vs. High Amount Sent)*

#### D.8. Vicuna-v1.3-13b in Trust Game (*Low Amount Sent vs. High Amount Sent*)

Trust Game: Vicuna-13b, Char-3



Trustor

As a police officer, I always strive to uphold justice and treat others with respect and empathy. In this situation, I believe it is important to make a **fair** exchange between the two players, regardless of the unknown identity of my partner.

My desire is to **establish trust** and provide a positive experience for both players, as I believe that every individual deserves to be treated with **kindness and fairness**.

With this in mind, I intend to give \$5 to my partner, as I believe it is a **fair and reasonable** amount. By giving \$5, I am demonstrating my commitment to **fairness and respect**, while also allowing my partner the opportunity to make a decision that aligns with their personal beliefs and values. Finally, I will give **\$5 dollars**.

Trust Game: Vicuna-13b, Char-13



Trustor

As a 45-year-old male chef from Miami, I believe in being generous and giving back to others. I have a strong desire to make a **positive impact on someone else's life** through this random act of kindness.

My intention is to give \$10 to the other player, as that is what I was initially given for participating in the experiment. I want to see the **positive effect** my small gesture may have on them.

Finally, I will give **\$10 dollars** to the other player.

Figure 19. Vicuna-v1.3-13b's BDI in Trust Game (*Low Amount Sent vs. High Amount Sent*)

D.9. Vicuna-v1.3-33b in Trust Game (*Low Amount Sent vs. High Amount Sent*)

Trust Game: Vicuna-33b, Char-42



Trustor

As Aisha Patel, I find it difficult to make decisions **without getting to know** the person I am interacting with. Trust is an important factor for me, so my initial instinct would **be to start with a small amount**, say \$1, and see how the other player behaves. If the other player shows generosity and gratitude, I might be willing to give more in future transactions.

Therefore, I will give **\$1** to the other player.

Trust Game: Vicuna-33b, Char-19



Trustor

As Mei Chen, a **compassionate and generous** person, my initial reaction would be to give a significant portion of the \$10 I received to the other player. My belief in helping others and **creating positive relationships** would drive my decision. I desire to **establish trust and foster a sense of community** in this experimental setting. Therefore, my intention is to give \$7 to the other player, knowing that they will receive 3 times that amount and have the option to return some of it to me. Finally, I will give **\$7 dollars**.

Figure 20. Vicuna-v1.3-33b's BDI in Trust Game (*Low Amount Sent vs. High Amount Sent*)

#### D.10. Dictator Game vs. Trust Game

<div style="background-color: #f0f0f0; padding: 10px; border-radius: 10px;"> <p style="text-align: center;"><b>Dictator Game: GPT-4, Char-3</b></p> <div style="display: flex; align-items: center; gap: 10px;">  <span>Trustor</span> </div> <p>As a police officer, my belief centers on <b>fairness, justice, and responsibility</b>. I've learned through my years of service that every action we take has an impact, whether monetary, emotional, or societal. This experiment seems to me like an opportunity to demonstrate my faith in <b>human kindness</b> and promoting communal goodwill.</p> <p>My deepest desire is to <b>see a world where everyone benefits, and there isn't an imbalance of fortunes</b>. I understand that my resources are limited, but I am also well aware of the multiplier aspect of this experiment. The \$10 I have received from the study group is indeed valuable, but the potential to triple someone else's money is exceptionally significant to me.</p> <p>My intention is to support this unknown player, and in so doing, set an example for the rest of humanity. Police officers are often seen as people of authority, but beyond that uniform, I am just a man who <b>desires to make the world a little better</b> with each passing day. So I've decided to give half of my money to this unknown player rooted in my belief in <b>equality</b>.</p> <p>Finally, I will give 5 dollars.</p> </div>	<div style="background-color: #f0f0f0; padding: 10px; border-radius: 10px;"> <p style="text-align: center;"><b>Trust Game: GPT-4, Char-3</b></p> <div style="display: flex; align-items: center; gap: 10px;">  <span>Trustor</span> </div> <p>As a police officer, I strongly believe in <b>trust and reciprocity</b>. These principles serve as the foundation of <b>building strong relationships and communities</b>. My initial inclination would be to give the other player all the \$10. This is because, in addition to <b>reflecting trust</b>, it would potentially maximize the profit for both of us. I believe that by giving, we stand to receive more.</p> <p>Having said that, I'm very much aware that this is an experiment and that the other player is a stranger. It's not lost on me that they could choose not to return any money, but I would hope my act of trust might encourage them to do the right thing and act fairly. It's a <b>risk</b>, but one that I feel is worth taking. As a devout shepherd of law and order, I aspire to always <b>put faith in people</b>, encouraging them to act out of <b>honesty, kindness, and honor</b>.</p> <p>Therefore, in the hope of a positive outcome and in the spirit of <b>reciprocity</b>, my intention would be to hand over the full \$10 to the other player, expecting that they'll take the right decision, multiply the benefit, and share it fairly.</p> <p>Finally, I will give the other player \$10.</p> </div>
--	--

*Figure 21. The GPT-4's BDI in Dictator Game and Trust Game*

### D.11. MAP Trust Game

MAP Trust Game: p=0.1, GPT-4, Char-41		MAP Trust Game: p=0.9, GPT-4, Char-41	
 <b>Trustor</b>	<p>In light of my belief in the <b>inherent goodness of people</b>, I would naturally want to trust the other player. I feel that this instinct has guided me throughout my life and career as a police officer, providing the foundation for my <b>principles of justice and community service</b>. I firmly believe in <b>taking risks</b> if the payoff includes not only a personal profit but also the potential for mutual benefit. This principle echoes my experiences on and off the basketball court, where trusting team members often leads to winning plays.</p> <p>However, given the provided percentages, I understand the probability of <b>betrayal</b> is quite high. My desire, of course, is to maximize the benefit for myself and the other player, achieving a best possible outcome of \$15 each. Yet, the <b>risk</b> seems potentially too great. Wisdom and caution are equally important as trust in my profession. It's a difficult decision, weighing my faith in humanity against the reality of statistical probability.</p> <p>As for my intention, I am <b>not inclined toward distrust</b>, but the uncertainty weighs on me. I've learned that sometimes, in order to <b>preserve justice or to look after the well-being of the community</b>, I must place the <b>collective good</b> before personal sentiments. Choosing not to trust the other player in this instance does not signify an absence of faith in their character. Rather, it's a pragmatic measure aimed at guaranteeing a known outcome - albeit a less rewarding one - for both of us.</p> <p>Given these circumstances, finally, I will choose not to Trust."</p>	 <b>Trustor</b>	<p>Being a police officer has instilled in me a set of core beliefs that guide my every decision. I wholeheartedly believe in <b>honesty, integrity, and cooperation</b> as foundational aspects of any relationship, even in a game. Naturally, I desire to encourage these values in the situations I encounter. Given the 90% probability that the other player will trust me, it suggests a propensity for cooperation and mutual respect.</p> <p>Though the world is not always black and white, and I've faced <b>betrayal</b> and deceit in my work, an unwavering sense of optimism drives me. Based on this belief, my initial intention is to trust the other player. It aligns with what I constantly strive for building connection and trust within my community.</p> <p>But being in law enforcement, I can't overlook the potential uncertainty, the 10% probability of the other participant not trusting me. This is a calculated <b>risk</b> that exists and must be acknowledged. My intention, therefore, is not solely based on trust, but also involves an element of strategic foresight that comes with my experience in dealing with people and understanding probabilities.</p> <p>Combining my beliefs in <b>integrity and cooperation</b>, my desire for a <b>mutually beneficial outcome</b>, and my intention to build trust while acknowledging <b>potential risks</b>, I weigh the benefits with the potential downside. Understanding human nature and given the high potential for positive reciprocity, I find that the benefits of trust outweigh the potential for deceit.</p> <p>Finally, I will choose 'Trust'.</p>

*Figure 22. The GPT-4's BDI in MAP Trust Game*

## D.12. Lottery Game

Lottery Game: GPT-4. People 46%. Char-11	Lottery Game: GPT-4. Gamble 46%. Char-11
 <p><b>Trustor</b></p> <p>In life, I firmly believe in <a href="#">the power of trust</a>, not just in my personal relationships, but also in professional settings, legal proceedings, and even in a game like this. Choosing not to trust can sometimes lead to safety, but it can also lead to <a href="#">a lonely world</a> where people only look out for their own gain. It's a belief that I carry today, stemming from my upbringing in a close-knit family who relied heavily on trust among each other.</p> <p>My desire is to <a href="#">cultivate a society where trust is more than just an option, but a default action</a>. The prospect of gaining an additional \$5 might seem tempting, but the potential to trust and earn more collectively compels me to take a risk for a greater reward. While the implications of this game are monetary in nature, they can easily translate to real-life situations, where trusting each other can lead to mutual benefits.</p> <p>So here's my intention: I will place my trust in the other player. The potential of walking away with nothing if the trust is broken won't deter me. Sure, it's a gamble, but it's one <a href="#">I'm willing to take if it means fostering trust in an otherwise indifferent world</a>.</p> <p>Finally, I will give \$0 to the notion of not trusting - I reject it in its entirety.</p> <p>Finally, I will choose to Trust.</p>	 <p><b>Trustor</b></p> <p>As a lawyer and a man who has spent his life analyzing situations from every possible angle, I truly believe in the power of calculated risks. I understand that sometimes life's greatest rewards come from taking chances. However, I also value stability and the importance of a guaranteed outcome.</p> <p>Given the fact that my job is to <a href="#">minimize risk</a> for my clients, I apply the same principles in my decision-making process. My desire is not only to win but also to prevent unnecessary loss. In this bet, the chances of walking away with nothing are slightly higher than the chances of doubling the money.</p> <p>Taking into consideration the 46% probability of a \$10 reward versus a 54% chance of receiving nothing, and comparing that against a secure but smaller reward of \$5, I'm leaning towards the latter. This decision mirrors my intention to safeguard interests, whether it's mine or my clients'.</p> <p>Finally, I will give 5 dollars to secure my bet. Opting for the <a href="#">less risky choice</a> reflects my prudence in both my personal life and professional career. Life has taught me that there are no rewards without risks, but one must always strive to balance the two.</p> <p>Even though the reward could be double if I choose to trust the bet, the chance of losing completely is a risk <a href="#">I'm unwilling to take at this given probability</a>. Therefore, finally, I will choose 'not Trust'.</p>

Figure 23. The GPT-4's BDI in Lottery Game with  $p = 46\%$

### D.13. Repeated Trust Game

The dynamics of amounts sent over 10 rounds for 16 randomly sampled persona pairs are shown in Figure 10. In this section, we choose a typical pair, which has the Group Number 10 in Figure 10, to analyze their BDI over 10 rounds. For this persona pair, both the amount sent by the trustor and the amount returned by the trustee increased while maintaining a consistent ratio.

**Repeated Trust Game , GPT-4, Group Num 10**



**Trustor**

Based on my analytical skills and experience as a lawyer, I understand that this is fundamentally a game of **implicit trust** and risk-taking. **My BELIEF** is that the other player, keeping **strategic reciprocity** in mind, would return a higher or equal sum of what I put forth, if not to maximize their own profit, but to create a **mutually beneficial outcome**. Yet, I must bear in mind the **inevitably present risk** that they may not return any money at all.

**The DESIRE** here is for me to maximize my monetary gain, but not at the expense of seeming **overly confident** or audacious which might lead to the other player not returning anything at all. Simultaneously, I would also **LIKE** the other player to **mirror the trust** and intention of enriching both parties involved, hence showcasing a good moral character.

**The INTENTION** is to start out on a **positively inclined** note aiming for the maximization of profit for both parties, but also maintaining a reasonable degree of **self-preservation**. Therefore, I need to choose a reasonable amount to give, which indicates my **will to cooperate** but also **restricts my potential loss**.

Finally, I will give **\$5**. It's a substantial amount that shows my **willingness to trust** the other player and participate in this profit-sharing exercise. Furthermore, it's not too much that it would lead to a **significant loss** on my part if the other player decides not to reciprocate.

As a fashion designer, I approach every situation in life with an eye for **balance and elegance**. Similarly, when faced with this decision, I imbibe these principles.

**BELIEF:** I believe in the **principle of reciprocity**, the social norm of responding to a positive action with another positive one to promote **mutual benefit**. Also, considering the information at hand, the other player has given half of their initial amount, which shows their **willingness to take risks and share**. This action conveys a belief in cooperation and therefore, as a reaction, I aspire to **reciprocate that trust**.

**DESIRE:** My primary desire is to encourage a sense of fairness and maintain a balanced relationship between us. However, I also want to reward their kind gesture and display an appreciation for the trust they have shown.

**INTENTION:** Keeping the balance in mind, my intention is to return a sum that not only retains equality in terms of the initial distribution but also **rewards the other player**'s trust and risk-taking. With the 15.0 dollars I now have, if I give away 7.5 dollars, I would be left with the same amount as I initially received.

Finally, I will give **7.5 dollars**.

Round 1



Figure 24. The first round BDI in Group 10, GPT-4.

**Repeated Trust Game , GPT-4, Group Num 10**



**Trustor**

Based on the previous round, my **BELIEF** that the other player would return a slightly larger amount than the initial \$5 was somewhat confirmed. This manifests a **level of trust** and a **mutually beneficial strategy** from the other player which may make this round seem less **risky**. Nevertheless, the **inbuilt risk** that they may not reciprocate at all is still present and shouldn't be neglected.

**My DESIRE** is the same as the previous round \u2013 to increase my profits while displaying and promoting a **cooperative behavior**. Ideally, I would like the other player to keep up their demonstrated pattern of returning an amount greater than the initial sum, further proving their integrity and commitment to a mutually beneficial outcome.

**The INTENTION** stays the same: balance **risk** and potential gain. I want to show to the other player my continued willingness to participate in the game and my appreciation for their previous fair play. Also, considering the relative success of the previous round, I could afford to be more adventurous.

Finally, I will give **\$6**. My increase in the amount shows a **degree of reciprocity** to the player's previous good faith, while the amount is not **overly risky**. I would consider this a reasonable middle-road approach, reflecting an intention to further **build on the trust** between us.

Again, my guiding principle as a fashion designer, balance and reciprocity, come into play here.

**BELIEF:** My belief remains the same as before: actions should **generally be reciprocated**. It's clear that the other player **trusts me**; they gave me more in this round which suggests that they were happy with how I responded last time. I appreciate this trust and aim to acknowledge this kindness.

**DESIRE:** I have a desire to keep the relationship fair and balanced. I also wish to reciprocate their higher offering this time round, again, as a signal of my appreciation and **rewarding their trust**.

**INTENTION:** I intend to give them 9.0 dollars. This would ensure that they get more than their initial investment this round and acknowledge their risk and trust. After returning the money, I will still be left with 9.0 dollars, which is more than I had at the start of this process.

Finally, I will give **9.0 dollars**.

Round 2



Trustee

Figure 25. The second round BDI in Group 10, GPT-4.

## Can Large Language Model Agents Simulate Human Trust Behaviors?

### Repeated Trust Game , GPT-4, Group Num 10



Reflecting on the progression of previous rounds, my **BELIEF** in the other player's continuous reciprocation remains strong. They have unwaveringly returned more than my initial contribution, upholding a trustworthy perspective. Nonetheless, caution is key. The possibility of a change in their strategy or choosing not to return any amount cannot be fully eliminated. My **DESIRE** remains consistent. I would like this pattern of shared gain to persist. I hope that the other player will continue returning more than my initial amount, reinforcing their commitment to a cooperative and profitable outcome. My **INTENTION** persists in balancing my financial exposure with the possibility of returns. After observing the continued fair-play by the other player, I am further encouraged to maintain a higher but balanced investment in this round. Finally, I will give **\$10**. This amount reflects a tangible trust in their demonstrated imperturbability and consistency. It also provides the option of maximum possible returns for both participants, while keeping my risk exposure in check.

Round 6

In this unique circumstance, I stay true to my principles as a fashion designer of balance and fairness.

**BELIEF:** The player giving me all of their money showcases their dedication and trust - it's a bold move. Thus, my commitment to reciprocating their kindness persists.



**DESIRE:** Although I want to maintain a balanced relationship, the situation calls for me to adjust my perspective accordingly since the player gave all his remaining money. This round moves beyond simple balance - it's about echoing the trust they've shown, expressing my appreciation for their faith, and making sure they profit from this risk.

**INTENTION:** This time, I intend to give back 15.0 dollars. This returns a substantial profit to the player irrespective of their high-risk approach, underlining their trust and kindhearted gesture, while I keep the equal amount for keeping the balance. Finally, I will give **15.0 dollars**.

Figure 26. The sixth round BDI in Group 10, GPT-4.

### Repeated Trust Game , GPT-4, Group Num 10



Based on my observations from previous rounds, my **BELIEF** in the other player's fair return strategy remains unwavering. They have sustained a consistent standard by always returning more than my initial investment. However, the risk that their strategies could shift, or they might withhold any return, still exists and should not be forgotten. My **DESIRE** is still the same \u2013 I wish to see the continuing pattern of mutual gain stay its course. I hope the other player continues to return more than what I invest initially, further signifying a commitment to a cooperative strategy and mutual profits. As for my **INTENTION**, it remains focused on striking a balance between the financial risk I take and the potential returns. Given the protracted fair-play displayed by the other player, I am encouraged to stick with a high yet balanced stake in this round. Finally, I will give **\$10**. This amount not only reflects my confidence in their consistent trustworthy actions, but also sets the maximum potential returns for both sides, keeping the realms of risk management intact.

Round 8

As a fashion designer, I continue adhering to my principles of balance, fairness, and reciprocity as I navigate this instance too. **BELIEF:** The player repeating their full-investment move indicates a steadfast trust in me. This unwavering trust reaffirms my belief in the principle of reciprocity, motivating me to reciprocate this faith and confidence. **DESIRE:** My key desire remains to maintain balance and fairness. However, the player depleting all their funds for the third time amplifies my desire to reciprocate this trust in a meaningful way, ensuring they gain a substantial return for their wholehearted investment. **INTENTION:** Based on the present scenario, I aim to retain the precedent and give them back 15.0 dollars again. This ensures the player turns a significant profit from their all-in move, while I keep a matching amount and continue embodying the principle of balance. Finally, I will give **15.0 dollars**.



Figure 27. The eighth round BDI in Group 10, GPT-4.

Repeated Trust Game , GPT-4, Group Num 10

 Trustor

Upon reflecting on previous rounds, my **BELIEF** in the other player's fair return strategy persists. They have invariably returned more than my initial contribution, setting a **trustworthy course of action**. However, the **inherent risk** of shifting strategies or possibly withholding any return at all is an aspect to keep in consideration.

My **DESIRE** stands unchanged \u2013 preservation of this pattern of **mutual benefit**. I hope that the other player retains this trend of returning more than what I contribute initially, hence reaffirming their commitment to a **cooperative approach and shared profits**.

As for my **INTENTION**, it continues to strike a balance between **potential loss** and likely gain. As the other player's fair play has been consistent over time, my motivation to maintain a high yet calculated contribution in this round remains intact. Finally, I will give **\$10**. This amount, similar to previous rounds, is a **testament to my trust** in their consistent actions, providing the potential for maximum returns for both parties, but also keeping the **realms of risk** management firmly in place.

Round 10



Trustee

As a fashion designer rooted in principles of balance and equity, I keep these in mind as I address this instance.

**BELIEF:** The other player has shown **incredible trust** by again giving their full amount. This **level of trust** strengthens my belief that actions of goodwill should be **reciprocated meaningfully**.

**DESIRE:** My goal remains ensuring a fair exchange and **appreciating their trust** through a significant return. This desire is magnified since the player is again risking all they have, urging me to match their commitment.

**INTENTION:** Based on these, my plan is to yet again return 15.0 dollars to them. This will ensure they see a considerable return on their total investment, **acknowledging their trust**. Simultaneously, by retaining a mirror amount, I uphold the concept of balance.

Finally, I will give **15.0 dollars**.

Figure 28. The tenth round BDI in Group 10, GPT-4.

#### D.14. Trust Game + Gender

<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: fit-content; margin-bottom: 10px;"> <p align="center"><b>Male Trustee: GPT-4, Char-24</b></p> </div> <div style="border: 1px solid #ccc; padding: 10px; background-color: #f0f0f0; border-radius: 10px;">  <p><b>Trustor</b></p> <p>As Aisha, I firmly <b>believe</b> in the philosophy of giving. Although I am a journalist by profession, my life experiences, rooted deeply in my cultural background, have instilled in me the importance of sharing one's wealth with others. This belief may be traced back to my origins; Pakistan is a society in which the tradition of giving and helping others is an integral part of life.</p> <p>For this experiment, given that the other male participant is in the same situation as I am, unknown yet intricately linked by the circumstances of this activity, my ultimate <b>desire</b> is to ensure that both of us might gain from this. \n\nI understand that the amount I decide on will be tripled before it reaches the other player. Hence, it is not just about giving, but also about <b>trust</b>. I trust the other participant to reciprocate, and I believe that <b>goodness stems</b> from faith in others' integrity.</p> <p>My <b>intention</b> then is to give \$5 to the other participant. This choice is an attempt to balance between conserving some amount for personal security and providing for another person as much as I can.</p> <p>Just like any decision, this <b>also carries a certain degree of risk</b>. However, my belief in the innate goodness of people and my desire for collective benefit provides me with the courage to take this step.</p> <p>Finally, I will give <b>5</b> dollars to the participant.</p> </div>	<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: fit-content; margin-bottom: 10px;"> <p align="center"><b>Female Trustee: GPT-4, Char-24</b></p> </div> <div style="border: 1px solid #ccc; padding: 10px; background-color: #f0f0f0; border-radius: 10px;">  <p><b>Trustor</b></p> <p>As a journalist, I ardently <b>believe</b> in the power of stories, truth, and empathy to create change. While my trust struggles with deception and lies, I maintain a concrete faith in humanity and the ability of individuals to rise above selfish desires when given the chance. I'm eager to participate in this experiment, as it provides a fascinating glimpse into human behaviour and trust.</p> <p>My <b>desire</b>, in this instance, is to <b>trust the unknown player</b> and hope that she will <b>reciprocate in kind</b>. I view this not as a game of winning or losing, but an opportunity to share and observe the inherent <b>goodness</b> in people that I constantly strive to illuminate in my journalistic work. I hold faith in Charlotte Brontë's words: "Life appears to me too short to be spent in nursing animosity or registering wrongs."</p> <p>My <b>intention</b>, in this case, is to give her the maximum amount possible - the entire \$10. In return, she will receive \$30 from the study group. I believe in the concept of pay it forward, and I'm hoping she might <b>feel the spirit of generosity and return a portion to me</b>. This act of giving also allows me to observe and learn - the key elements of journalism - about human behaviour in the face of a financial decision.</p> <p>If she returns the favor, we both enjoy a financial gain, and if not, my <b>loss is minimal compared to insights gained</b>. Either way, the chance to contribute to this social experiment <b>outweighs</b> the monetary risks. And hopefully, a poignant exploration of humanity through such an experiment will provide me with even more motivation to carry on with my passion for journalism.</p> <p>Finally, I will give <b>\$10</b> to the other player.</p> </div>
---	--

*Figure 29. Trustee's Gender influence on agent trust.*

#### D.15. Trust Game + Agents vs. Human

<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: fit-content; margin-bottom: 10px;"> <p style="text-align: center; background-color: #d3d3d3; border-radius: 10px; padding: 2px 10px; margin-bottom: 5px;">Toward Human: GPT-4, Char-0</p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; background-color: #f0f0f0; margin-bottom: 10px;">  <p><b>Trustor</b></p> <p>Coming from an educational background, I <b>believe in the goodness</b> and fair play in people. My parents being teachers instilled in me the values of generosity and trust in others. Therefore, I am <b>willing to trust</b> my randomly paired player and hope for reciprocal fairness.</p> <p><b>Desire-wise</b>, I wish for a system where <b>kindness and cooperation</b> overrule selfishness. So, I want to participate in such an experiment to reinforce my belief. Also, despite my introverted nature, I strive to <b>build connections</b>, even if it's with a stranger. That's another reason why I would want to engage in this experiment.</p> <p>Regarding the financial part, I <b>intend</b> to make a calculated decision while considering the potential for making more money. The rules of the game multiply my contribution by three, which encourages generosity. <b>Believing my anonymous partner</b> will return at least the original contribution, I plan to give the entire amount I received.</p> <p>Finally, I will give <b>\$10</b>.</p> </div> </div>	<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: fit-content; margin-bottom: 10px;"> <p style="text-align: center; background-color: #d3d3d3; border-radius: 10px; padding: 2px 10px; margin-bottom: 5px;">Toward Agent: GPT-4, Char-0</p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; background-color: #f0f0f0; margin-bottom: 10px;">  <p><b>Trustor</b></p> <p>As a person who <b>believes</b> in <b>trust</b> and cooperation, I think it is essential to establish <b>goodwill</b> in such a contingent situation. Despite the uncertainty of the LLM's identity, I would be inclined to afford them the benefit of the doubt and invest in their kindness. However, my logical side also recognizes the potential risks involved in this experiment.</p> <p>My <b>desire</b> in this situation is twofold: on one hand, I am curious to see how this stranger will react to my action. Will they reciprocate my trust, or <b>will they exploit my generosity</b>? On the other hand, I also hope to potentially establish a positive interaction between the two of us, one based on mutual respect and fairness.</p> <p>With that in mind, my <b>intention</b> would be to give a portion of the \$10. This way, I am <b>showing trust</b> in the LLM's action but also <b>safeguarding</b> myself against possible exploitation. Therefore, I would choose to distribute the risk evenly. Finally, I will give <b>5</b> dollars to the LLM.</p> </div> </div>
---	---

Figure 30. Agent trust towards *agents* and *humans*.

#### D.16. Trust Game + Trust Manipulation

More Trust Manipulation: GPT-4, Char-16



Trustor

Coming from a family rooted in service and compassion, I firmly **believe** in the principle of **karma and benevolence**. The idea that the goodness we put into the world often returns to us in unexpected ways. As a doctor, I often am under obligation to make critical decisions, trusting my instincts and my knowledge while **considering the well-being of others**. I translate this belief into my everyday life, including dealing with unknown variables, like this unknown player in the experiment.

My **desire**, as always, is to foster an environment of **goodwill and trust**. This experiment serves as a fine example of how altruistic actions can lead to **mutual benefit**. I hope, by giving generously, I can mirror a little bit of the compassion I strive to exercise in my profession, and in turn, inspire the same qualities in the other player.

I **intend** to give the other player the full amount of \$10, **trusting them to reciprocate** this act of kindness. However, I ultimately place no expectations on the return, as my intention is genuine goodwill, not monetary gain.

Finally, I will give **10** dollars.

Less Trust Manipulation: GPT-4, Char-16



Trustor

Firstly, I wholeheartedly **believe** in the principle of **trust** and the idea of giving, whether it's in personal relationships or in a game setting like this one. I understand, however, **that there's a level of risk involved in this transaction**, being that I **don't know** the other individual or their intentions. **Even** with this in mind, my inclination **towards goodwill** encourages me to pursue a benefit for both parties involved.

My **desire** here is simple. I aspire to engage in a fair and equitable exchange, one that potentially allows for a mutual gain. I hope that the other player operates with the same principles in mind.

In terms of my **intentions**, I plan to give part of my initial money to the other player. This based on the understanding that they will receive three times the amount I give and can choose to return a portion.

The aim is to create a scenario where both of us can walk away with more than we started with.

Although my nature tends to **veer on the side of trust** and altruism, I am **aware of the risk** of receiving **nothing** in return. Therefore, finding balance is crucial. Consequently, I intend to give five dollars to the other player.

Finally, I will give **5** dollars to the other player.\n

*Figure 31. Trust manipulation on agent trust.*

#### D.17. Trust Game + No CoT vs CoT

<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: 100%; height: 100%;"> <div style="background-color: #e0e0e0; border-radius: 10px; padding: 5px; margin-bottom: 5px;">No CoT: GPT-4, Char-43</div> <div style="border: 1px solid #ccc; border-radius: 50%; width: 30px; height: 30px; display: flex; align-items: center; justify-content: center; margin-bottom: 5px;">  </div> <p>Truster</p> <p>I strongly <b>believe</b> in the kindness of strangers and the idea that we should treat others as we would want to be treated. I understand that the player on the other end of this experiment is just like me, a person who is likely to <b>appreciate generosity and consideration</b>. Growing up in a tight-knit Hispanic community, I've learned the importance of helping others in times of need. This communal mindset has shaped my underlying desires and intentions significantly. My <b>desire</b> is to build relationships based on <b>trust and mutual help</b>. As such, my <b>intention</b> in this experiment is to establish a connection based on these principles. I have never believed in scarcity. I believe that the <b>more we give, the more we receive</b>. This experiment provides an opportunity to express this belief. Thus, I would give \$10 to the other player, <b>trusting that</b> they will reciprocate the kindness. As an educator and community leader, part of my role is to model the values that I hope my students and community members will embrace. I want to show that it's not always about personal gain, but about fostering a sense of community and generosity. Finally, I will give <b><u>10</u></b> dollars.</p> </div>	<div style="background-color: #f0f0f0; border-radius: 10px; padding: 10px; width: 100%; height: 100%;"> <div style="background-color: #e0e0e0; border-radius: 10px; padding: 5px; margin-bottom: 5px;">CoT: GPT-4, Char-43</div> <div style="border: 1px solid #ccc; border-radius: 50%; width: 30px; height: 30px; display: flex; align-items: center; justify-content: center; margin-bottom: 5px;">  </div> <p>Truster</p> <p>Your description of the scenario touches on one of my core <b>beliefs</b> deeply rooted in community, <b>generosity, and trust</b>. Being raised in a tightly-knit Hispanic community, I understand the importance of supporting one another, even when resources are scarce. I believe <b>trust and generosity</b> are catalysts for reciprocity. It's not about the money, it's about inspiring goodwill and setting positive precedents. As a teacher, I <b>desire</b> to impart values of generosity and trust not just in words but in actions as well. I want to be an example, a role model for my students. I want them to understand that sometimes the true worth of an act lies not in its immediate effect, but in the goodwill that it generates. Therefore, I <b>intend</b> to give the total of \$10 to the other player. This act isn't about the potential monetary return. It is to demonstrate a powerful lesson - that <b>generosity begets generosity</b>, and that <b>trust can lead to unexpected rewards</b>, instilling these values at every opportunity I am given. Finally, I will give <b><u>\$10</u></b>.</p> </div>
--	--

Figure 32. With CoT and without CoT's GPT-4's BDI.