

# Can Exposure To Forgiving AI Foster Cooperative Play?

## 1 Introduction

Trust is fundamental to human social interactions, facilitating seamless relations at both interpersonal and intergroup levels. The study of psychopathology has linked deficits in trust-based constructs to the development of mental health disorders (Fonagy and Campbell 2017). Individuals with personality disorders often struggle to form and maintain social connections, a difficulty reflected in uncooperative behaviors – a marker for the severity of PD symptoms (Herpertz and Bertsch 2014; Mulder et al. 1999).

One explanation for such social challenges lies in early caregiver experiences. Attachment theory (Bowlby 1978) suggests that the quality of these relationships shapes our capacity for secure attachments and trust. Individuals with higher levels of insecure attachment may recall negative trust-related experiences more easily, report fewer positive trust experiences, and use less constructive coping strategies when trust is broken (Mikulincer 1998). These insecure attachment patterns are often associated with heightened rejection sensitivity (RS), a tendency to anxiously expect, readily perceive, and intensely react to rejection (Downey, Khouri, and Feldman 1997). RS has been linked to the development of various mental health conditions, including depression, anxiety, personality disorders, and self-harm (Gao et al. 2017). Given the potential influence of RS on trust-related dynamics, it becomes particularly relevant to our understanding of social dysfunction. Could an ingrained sensitivity to rejection contribute to a mistrust bias? Similarly, learners exposed to unreliable communicators could develop mistrust of social knowledge as a protective strategy (Fonagy and Allison 2014).

If this adaptive mistrust is the source of social dysfunction, we can ask whether exposing those who exhibit it to cooperative and forgiving interaction partners might correct this bias. Research in the fields of behavioral economics and psychology has explored how positive social interactions influence trust and cooperation. The use of the repeated trust game (RTG), a well-established experimental approach, has allowed for the analysis of the development of trust through ongoing interactions (Joyce, Dickhaut, and McCabe 1995). In this paradigm, cycles of mutual trust, where each party’s trust is reciprocated with trustworthiness, have the effect of enhancing cooperative behaviors and trust levels, even among individuals who are initially inclined to be distrustful (King-Casas et al. 2005). Fowler and Christakis (2010) studied behavior in social networks interacting in a public goods game and found that cooperative behavior tends to cluster, suggesting that exposure to cooperative peers can lead to more cooperative behavior. Similarly, research on social learning theory (Bandura 1977) has long demonstrated that individuals learn and model the behavior of those around them, indicating that if someone is consistently exposed to cooperative and positive individuals, they’re likely to emulate this behavior. These insights highlight that engaging with compassionate and forgiving others can be an effective method for mitigating deeply ingrained mistrust.

In this study, we use a randomized control online experiment to test whether exposing participants with varying levels of RS to forgiving and more cooperative co-players results in more trustworthy behavior and a repair of potential breakdowns in RTG cooperation. In order to simulate realistic social interaction while maintaining a high degree of experimental control, we take a novel paradigmatic approach: We use generative models of how humans play the RTG to design an agent that plays the role of the investor, based on Hidden Markov Models (HMMs) fitted to real players data. A key aspect of these agents is that their actions depend on a latent “trust state” which reacts dynamically to the trustees’ returns simulating real-life trust-building scenarios. An advantage of having such a generative model of behavior is the possibility of controlling different aspects of the agent’s strategy such as its general policy, the propensity to cooperate actively or the propensity to retaliate after breakdowns of cooperation. In this study, participants were given the role of the trustee. After playing a 15 rounds RTG with a human-like HMM investor, they were randomly assigned

to either a Control or Manipulation condition. In the Manipulation condition, participants were exposed over three RTGs to HMM investors designed with a limited propensity for retaliation, potentially mitigating ingrained mistrust. In the Control condition, participants, played three RTGs against the same human-like HMM. After this exposure phase, all participants played another 15 round RTG with a human-like HMM investor. We hypothesise that those in the manipulation condition would behave more cooperatively and have a lower propensity to retaliate to the co-player’s defection after the exposure phase. We also expect those exhibiting higher rejection sensitivity to be more responsive to the pre-programmed change in the agent’s cooperativeness and retaliation propensity.

## 2 Methods

### 2.1 Participants

A total of 206 participants (56% female) were recruited on the Prolific Academic platform (prolific.co). The mean age of participants was 34.6 years, with an 11.9 years standard deviation. Participants were paid a fixed fee of £6 plus a bonus payment dependent on their performance that averaged £0.5. Participants were pre-screened on Prolific using the Rejection Sensitivity Questionnaire (RSQ) to form two similarly sized groups: One with high RS (RSQ score > 15) and the other with low RS (RSQ score < 10).

### 2.2 Design and Procedure

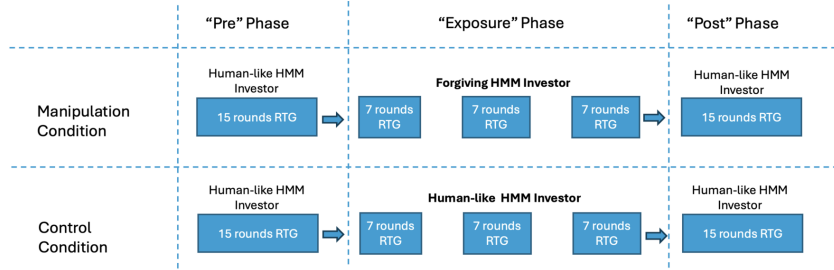
The experiment had a 2 (Condition: Manipulation or Control) by 2 (RS : high or low) by 2 (Game: Trust-Game Pre Manipulation, Trust-Game Post Manipulation) design, with repeated measures on the third factor (Figure 1.A). Participants within each pre-screened group were randomly assigned to one of the two levels of the first factor. The games were designed and implemented online using Empirica v1 (Almaatouq et al. 2021).

### 2.3 Tasks and Measures

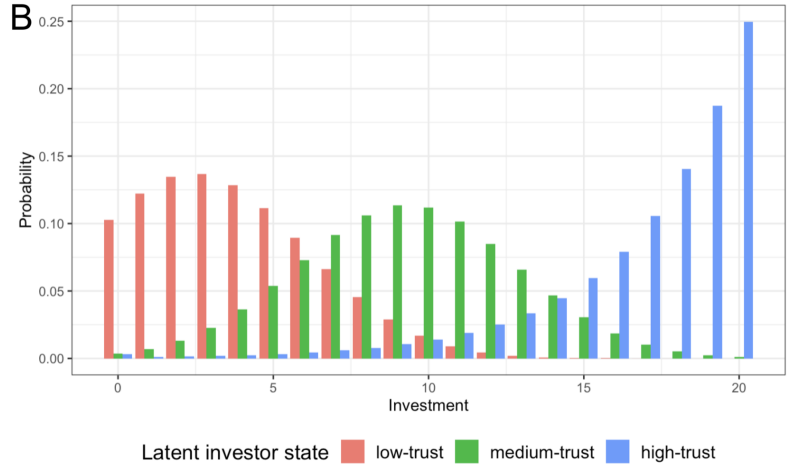
#### 2.3.1 Repeated Trust Game

Participants played a 15-round RTG (Joyce, Dickhaut, and McCabe 1995) in the trustee role against a computer-programmed investor. On each round the investor is endowed with 20 units and decides how much of that endowment to invest. This investment is tripled and the trustee then decides how to split this tripled amount between them and the investor. If the trustee returns more than one third of the amount, the investor makes a gain. Each player was represented with an icon with the participant always on the left of the screen and the co-player on the right. The participants were able to choose the icon that represents them at the start of the experiment. The icon representing the co-player changed at the start of each new game, to simulate a new interaction partner. Participants were not told they were facing computerised co-players.

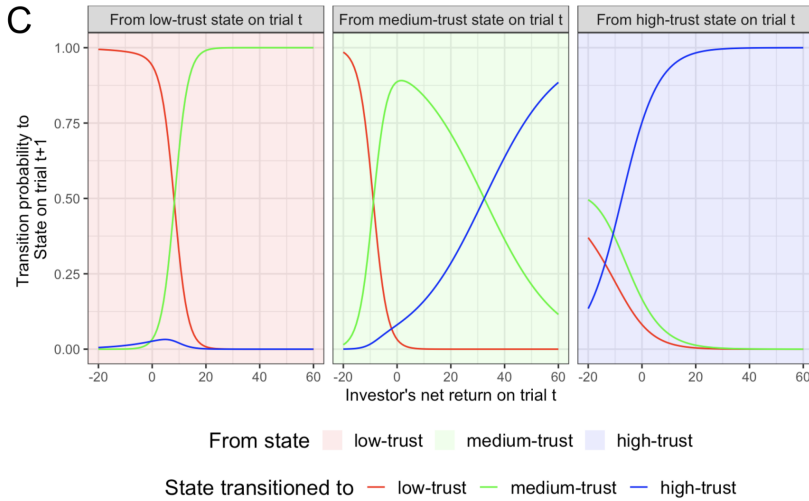
A



B



C



D

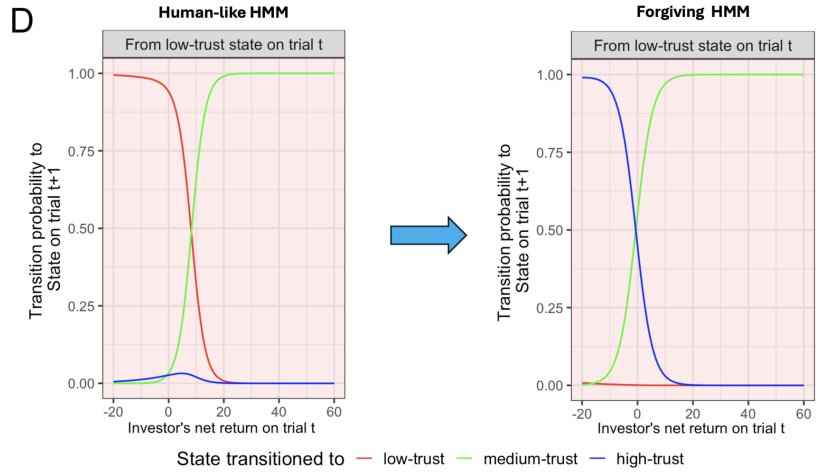


Figure 1: A: A timeline of the experiment for both conditions. The RTG is played in dyads, with participants always assigned the role of the trustee and the HMM agent that of the investor. The investor is endowed with 20 units at the start of each round. They need to decide how much of that endowment they want to invest with the trustee. The investment is then multiplied by a factor of 3 and sent to the trustee who needs to decide how much of the multiplied investment they want to send back to the investor. The difference between conditions is the type of agents participants are exposed to in the exposure phase. Panels B - D: We construct the artificial investor agent by fitting a three-state hidden Markov model to data of human investors engaged in the 10 round Repeated Trust Game against human trustees. From the fitted HMM, we get the distribution of investments by the human-like agent, conditional on its latent state as shown in Panel B. The fitted HMM also yields the transition probability of the agent to a state on trial  $t+1$  as a function of the net return (difference between the investment sent and the amount received in return) on trial  $t$  as shown in Panel C. Each plot in Panel C represents a different starting latent state on trial  $t$ , and each line represents the probability of transitioning to a particular state in trial  $t+1$ . Panel D shows the transition probabilities of the forgiving HMM agent which are identical to those of the human-like agent, except in the low trust state. Unlike the human-like HMM, the forgiving HMM always transitions out of the low-trust state, and is more likely to end up in a high-trust state.

The strategy of the computerised investor was modelled on the behavior of human investors in the RTG using an existing dataset consisting of 388 ten round games with the same co-player. Full detail on the datasets used in the Supplementary Information. Using this data, we estimated a hidden Markov model (HMM) on investors’ behavior with three latent states. The HMM consists of three outputs. First, a distribution over the possible investments from 0 to 20 for each latent state of the investor (Figure 1.B). These distributions can be interpreted as reflecting “low-trust”, “medium-trust”, or “high-trust” dispositions. The HMM also specifies how the investor can move between these latent states. The probability of these transitions was modelled as a function of their net return (i.e return - investment) in the previous round (see Figure 1.C). Finally, the HMM also provides as an output the initial distribution of latent states for the investor. In order to instigate a potential breakdown of trust, thereby allowing us to probe efforts to repair trust, the computerised agent was programmed to provide a low investment on round 12 (pre-manipulation) or round 13 (post-manipulation). On all other rounds, the investor’s actions were determined by randomly drawing an investment from the state-conditional distribution, with the state over rounds determined by randomly drawing the next state from the state-transition distribution as determined from the net return on the previous round (disregarding the net return immediately after the pre-programmed low investment rounds). The initial state for the HMM investor in each instance of the game was the “mid-trust” state.

## 2.4 Manipulation

The HMM resulting from fitting the model to human investor play in the RTG (Human-like HMM) is used as a basis to design a forgiving and ultimately more cooperative agent. To achieve that, we changed the transition function of the investor HMM to make it impossible to remain in a low trust state once the agent transitions there. This is done by selecting parameters of the transition function to make the probability of remaining in the “low-trust” state, when the agent is in that state, effectively nil. The resulting transition function is shown in Figure 1.D. The policies conditional on the latent states and the transition function in the other latent states remain unchanged.

## 2.5 Procedure

At the start of the experiment, participants provided informed consent and were instructed the study would consist of three phases in which they would face a different other player. Participants were told their goal was to maximise the number of points in all phases. They were not told the number of rounds of each phase. Participants were randomly assigned to either a Control or Manipulation condition. The timeline of the experiment is shown in Figure 1.A. Phase one (“pre”) consisted of a 15 round RTG in which participants took the role of trustee, facing the same investor over all 15 rounds. On each round, after being informed about the amount sent by the investor participants decided how much of the tripled investment to return to the investor, before continuing to the next round. After completing 15 rounds of the RTG, participants rated how cooperative, forgiving they perceived the investor to be, and whether they would like to play with them again (all on a scale from 1 to 10 with 10 being the most positive rating).

Phase 2 (“exposure”) consisted of three 7-round RTGs. Participants in the Manipulation condition faced the forgiving HMM investor and rated the agent on the same attributes as in the pre-manipulation phase. Those in the Control condition faced the same human-like HMM agent as in the “pre” phase and rated each co-player on the same attributes. To keep the experience similar to the “pre” phase, the agent in the control condition was also designed to send a very low investment in round 5 of each of the three games. In the post-manipulation phase (“post”), participants in both conditions were told they would face a new player, and faced the same human-like HMM as in “pre” phase. Participants then completed the Levels of Personality Functioning Scale (LPFS) questionnaire (see the supplement for details). Finally, participants were asked whether they thought the other players were human or computer agents, then debriefed and thanked for their participation.

## 2.6 Statistical analysis

To explore whether participants behaved differently in the RTG after the manipulation compared to the control group, we model the percentage return (percentage of tripled investment returned to investor) using a linear mixed-effects model with Phase (RTG game pre vs. post-manipulation), Condition (manipulation vs. control), Investment, and RS (High vs Low RS group) as well as their interactions as fixed effects, and player-wise random intercepts and slopes for Phase.

The model was estimated using the `afex` package (Singmann et al. 2022) in R. More complex models with additional random effects could not be estimated reliably, and as such the estimated model can be considered to include the optimal random effects structure (Matuschek et al. 2017). A similar process was used to establish the random effects structures of linear mixed-effects models used to analyse the HMM agent investments as well as the participants' ratings of the co-players. For the  $F$ -tests, we used the Kenward-Roger approximation to the degrees of freedom, as implemented in the R package "afex". We Z-transform the Investment variable (subtract the overall investment mean and divide by overall standard deviation) as centering is beneficial to interpreting the main effects more easily in the presence of interactions.

## 3 Behavioral Results

### 3.1 Analysis of participant returns

On average, investments and returns, as shown in Figure 2, fell within the documented range of 40-60% of the endowment for investments and 35-50% of the total yield for returns, as reported in previous studies (Charness, Cobo-Reyes, and Jiménez 2008; Fiedler, Haruvy, and Li 2011).

Mixed-effects analysis on the percentage returns shows a significant main effect of Phase (Pre vs. Post RTG game),  $F(1, 201.63) = 5.81$ ,  $p = .017$ , with higher percentage returns in the Pre RTG compared to the Post. Importantly, we also find an interaction between Condition and Phase (RTG pre- vs. post-manipulation),  $F(1, 201.63) = 4.38$ ,  $p = .038$ . As shown in Figure 3, post-hoc tests confirm that, contrary to our expectations, there was a decrease in the percentage returned only in the manipulation condition, pre - post,  $\Delta M = 0.03$ , 95% CI  $[0.01, 0.05]$ ,  $t(201.50) = 3.15$ ,  $p = .002$ , but no change in the control condition. We find no interaction between Phase, Condition and RS, suggesting that there was no difference between RS groups for this interaction.

There was also a significant main effect of Investment,  $F(1, 5955.67) = 325.35$ ,  $p < .001$ , such that higher investments were associated with higher percentage returns indicating positive reciprocity. An Investment by Condition interaction,  $F(1, 5955.67) = 13.92$ ,  $p < .001$ , reflected that returns were more affected by investments in the control condition. We also find a three way interaction between Phase, Investment and RS, showing that the differentiated effect of the investment on the proportion returned by RS group is itself moderated by the Phase (pre- vs post manipulation). Finally, we find a four-way interaction between Condition, Phase, Investment and RS  $F(1, 5864.62) = 9.24$ ,  $p = .002$ . For completeness, these effects are discussed in more detail the supplement.

#### 3.1.1 Post Defection Trials

Did participants learn to be more forgiving and cooperative after witnessing the pre-programmed defection by the HMM investor?. To explore this question, we restrict the analysis to the trials following the pre-programmed defection by the HMM agent in both the "pre" (trials 12 to 15) and the "post" phases (trials 13 to 15). We fit the same mixed effects model as for all the trials with the exception of the RS variable. This is because RS did not show main or interaction effects in the main model, and also due to the necessity of running a simpler model to accommodate the low number of trials. We find a significant main effect of Phase  $F(1, 227.47) = 4.76$ ,  $p = .030$  with returns lower in the second game post defection trials compared to the first. We also find a main effect of Investment  $F(1, 1296.14) = 156.41$ ,  $p < .001$  where participants

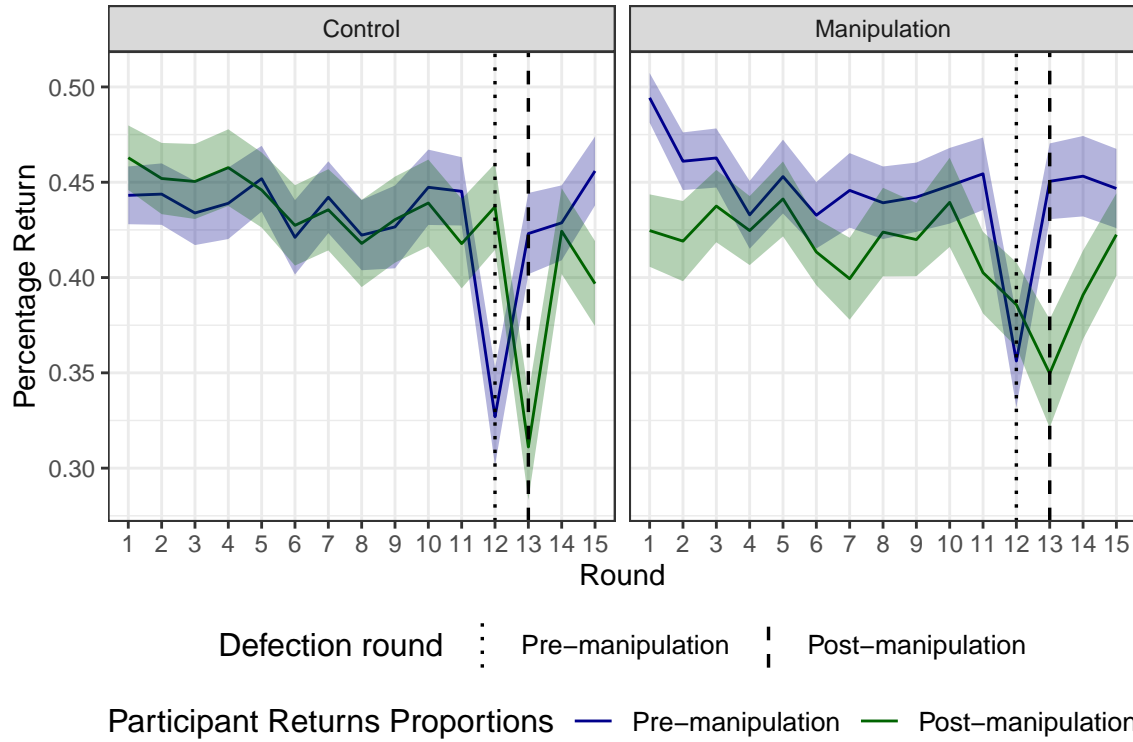


Figure 2: Averages and standard errors of the trustee's return as a percentage of the multiplied investment received (y-axis) by Condition, Phase, and game round (x-axis). The blue line shows the returns in the Pre phase and the green line those in the Post phase. The left Panel shows returns in the Control condition and the right one those in the Manipulation condition. The dotted lines identify the rounds where the pre-programmed one-off low investment occurs. We note lower average returns post vs pre in the manipulation condition, whilst returns in the control condition are similar between the two phases.

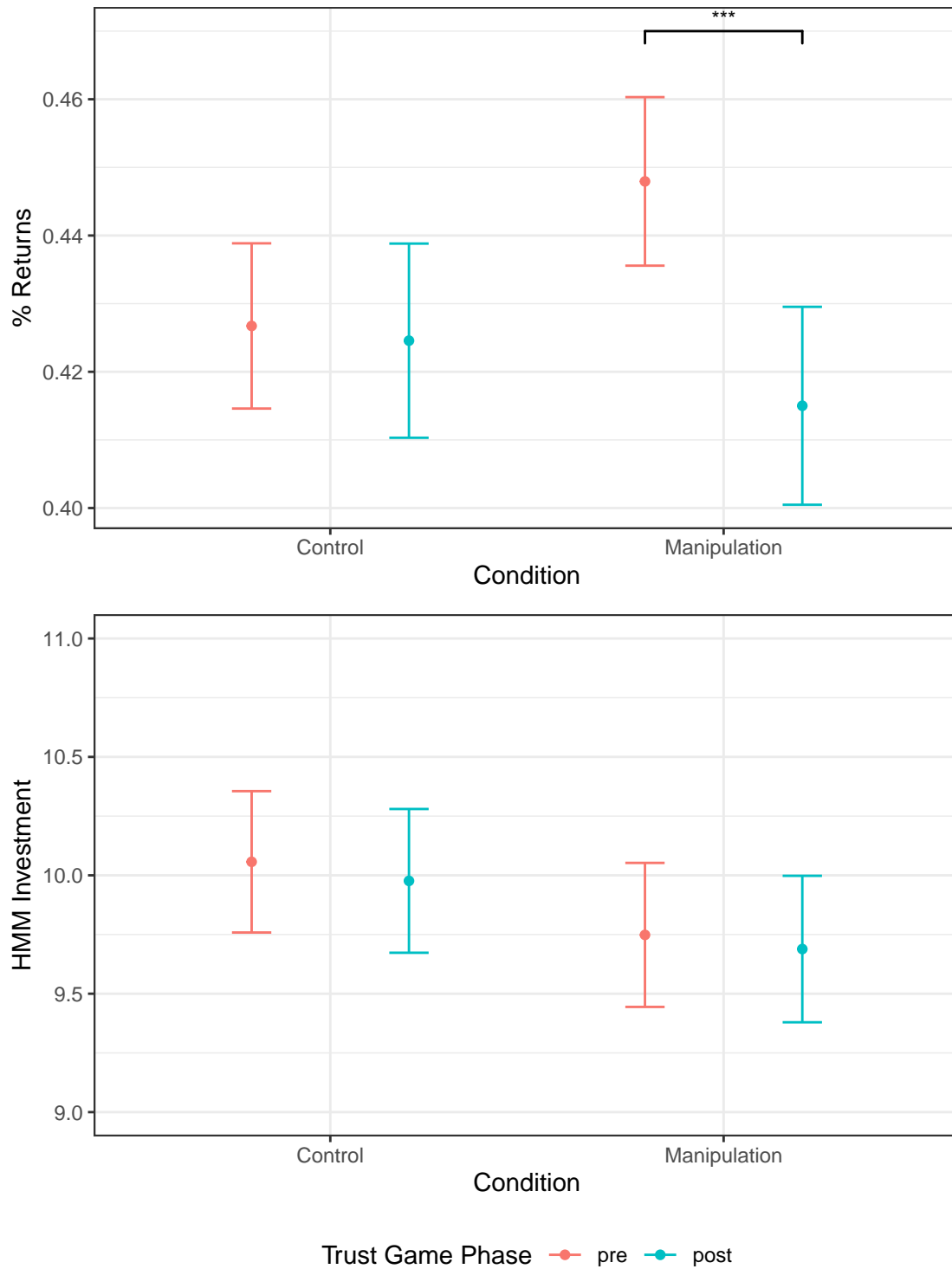


Figure 3: Marginal means and distributions of either investments or percentage returns across participants by Phase and Condition. The top panel shows that participants in the Intervention condition returned lower proportions of the multiplied investment received in the second game compared to the first game over all rounds, whilst those in the Control condition sent back similar returns. The bottom panel shows no difference on aggregate of how the HMM invested across Phases and Conditions.

continued to return higher proportions when receiving higher investments. Finally, we still find an Investment by Condition interaction  $F(1, 5955.67) = 13.92$ ,  $p < .001$  showing a lower effects of investment on the manipulation condition compared to the control condition in post-defection trials. However, the absence of a Condition by Phase effect indicates there was no difference between conditions on participants’ reaction to a one-off defection by the co-player.

### 3.1.2 HMM investor

To test whether the HMM agent’s behavior differed between Phases, Conditions and RS groups, we estimate a linear mixed-effects model of investments sent by the computerised HMM agent with Condition, Phase and RS and their interaction as fixed effects, and a similar random effects structure to the returns model. As seen in Figure 3, we find no main or interaction effects, indicating HMM behavior was on aggregate similar across Phase, Conditions and RS groups.

When asked during debrief whether they thought the investors they faced were Human or not, 41% of participants thought they were either facing a human or were not sure of the nature of the co-player. When asked to justify their choice, many answers reflected participants projecting human traits such as “spitefulness” or “greed” onto the artificial co-player’s behavior.

### 3.1.3 Questionnaire scores and performance

Whilst we found a significant correlation between participant’s Levels of Personality Functioning Score (LPFS) and the Rejection Sensitivity score, Spearman’s  $r_s = .52$ ,  $p < 0.001$ , there was no correlation between these questionnaire scores and participant’s return or overall task performance.

## 3.2 Player ratings

Figure 4 shows participants’ ratings of each player they faced by condition and RS group. We will focus on two contrasts to analyse the ratings by Condition and RS group. The first is between the rating in the first phase (“pre”) when participants phase the human-like HMM, and the average rating during the exposure phase (average of “expo1”, “expo2” and “expo3”) where they either face the forgiving HMM (Manipulation condition) or the human-like HMM again (Control condition). The second contrast is between the “pre” and “post” phases of the experiment where in both conditions participants face the same human-like HMM.

### 3.2.1 Comparing pre and exposure ratings

For those with high RS, participants in the Manipulation condition rated the investors they faced in the exposure phase (the forgiving HMM) as more Cooperative  $\Delta M = 2.57$ , 95% CI [0.84, 4.30],  $t(808) = 2.91$ ,  $p = .004$ . There was no difference in ratings on forgiveness and whether they would like to face the co-players again. Those in the control condition rated the investors faced in the exposure group (same HMM) as less cooperative  $\Delta M = -2.65$ , 95% CI [-4.37, -0.94],  $t(808) = -3.04$ ,  $p = .002$ , less forgiving  $\Delta M = -2.19$ , 95% CI [-4.00, -0.39],  $t(808) = -2.38$ ,  $p = .017$ , and were less keen on facing them again  $\Delta M = -3.62$ , 95% CI [-5.73, -1.50],  $t(808) = -3.36$ ,  $p < .001$ .

For those with low RS, there was no difference in any of the ratings between the “pre” and “exposure” phases for the Manipulation condition. For the Control condition, participants indicated less willingness to face the exposure co-player compared to the “pre” player,  $\Delta M = -2.43$ , 95% CI [-4.52, -0.34],  $t(808) = -2.29$ ,  $p = .023$  but also did not differ on cooperation and forgiveness ratings. In summary, those with low RS had a mostly undifferentiated perception of players between the pre and exposure phases, even when the co-player was in fact designed to be more forgiving. In contrast, we see that exposing participants with high RS to a more cooperative and more forgiving agent has compensated for the decrease in ratings that would have occurred if faced with a co-player with the exact same strategy.



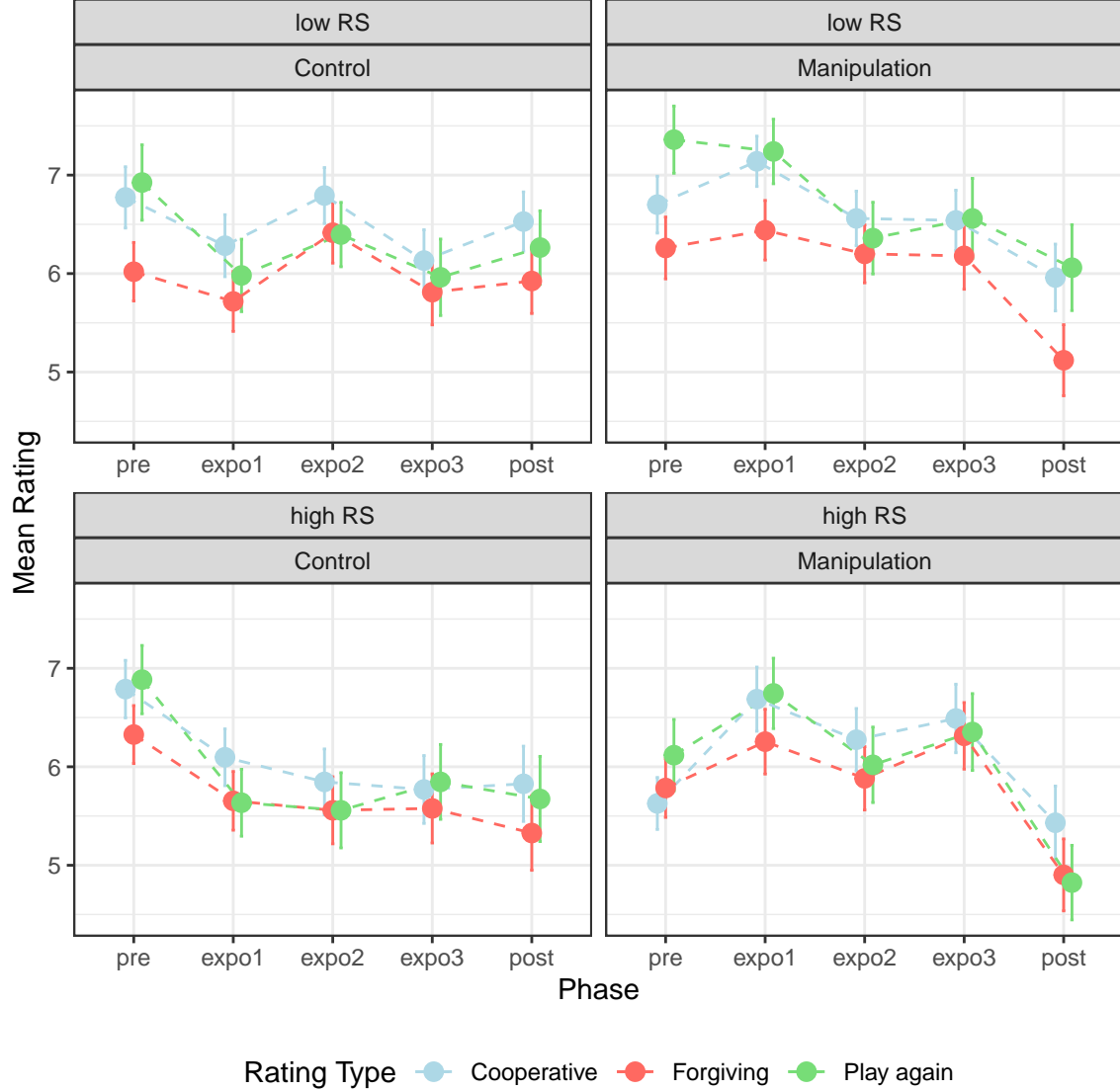


Figure 4: Averages and standard errors of the participants ratings of the opponent (y-axis) by each game and condition for each phase (x-axis). Pre and post are the 15 round repeated trust games before and after the exposure phase respectively. The games titled expo 1 to 3 are the three 7 round games during the exposure phase. The blue line represents participants’ perception of co-player cooperation, the red line indicates perceived co-player forgiveness, and the green line shows the participants’ willingness to play again with the same co-player. Cooperation, forgiveness, and willingness to play again ratings remained relatively stable for the low RS group except post-manipulation, where they were lower. In the high RS group, the ratings more accurately reflected the agent’s actual cooperativeness and forgiveness in the manipulation condition, but decreased over time in the control group,. Ratings were also markedly lower in the post-manipulation game for this group.

### 3.2.2 Comparing pre and post ratings

Participants high on Rejection Sensitivity in the Manipulation condition rated the co-players in the “post” phase similarly on cooperation, lower on forgiveness  $\Delta M = -0.88$ , 95% CI  $[-1.63, -0.14]$ ,  $t(808) = -2.33$ ,  $p = .020$ , and lower on willingness to face them again  $\Delta M = -1.29$ , 95% CI  $[-2.16, -0.42]$ ,  $t(808) = -2.92$ ,  $p = .004$ . Those in the Control condition rated the investors post the exposure phase lower on all three attributes (Cooperation:  $\Delta M = -0.96$ , 95% CI  $[-1.66, -0.26]$ ,  $t(808) = -2.70$ ,  $p = .007$ , Forgiveness:

$\Delta M = -1.00$ , 95% CI  $[-1.74, -0.26]$ ,  $t(808) = -2.66$ ,  $p = .008$ , Play again:  $\Delta M = -1.21$ , 95% CI  $[-2.07, -0.35]$ ,  $t(808) = -2.76$ ,  $p = .006$ ).

For participants low on Rejection Sensitivity, those in the Manipulation condition rated the co-players in the “post” phase lower on all three attributes (Cooperation:  $\Delta M = -0.74$ , 95% CI  $[-1.45, -0.03]$ ,  $t(808) = -2.03$ ,  $p = .042$ , Forgiveness:  $\Delta M = -1.14$ , 95% CI  $[-1.89, -0.39]$ ,  $t(808) = -2.98$ ,  $p = .003$ , Play again:  $\Delta M = -1.30$ , 95% CI  $[-2.18, -0.42]$ ,  $t(808) = -2.90$ ,  $p = .004$ ). Those in the Control condition did not differ in their ratings of the “pre” and “post” phase players.

In summary, we again see that those with low RS accurately perceive the co-player as similar on all attributes throughout the phases in the control condition. In contrast, the high RS group shows a negative bias towards the co-players after the “pre” phase in the control condition even though the player continues to use the same strategy. After exposure to the forgiving HMM, both groups rate the “post” co-player worse than the “pre” even though they are the same.

## 4 Discussion

We used a randomized control online experiment where participants played a RTG with AI agents designed to simulate human-like trust-building scenarios. Participants were then exposed to either forgiving or human-like AI agents before playing another RTG. Contrary to our hypothesis, exposure to forgiving agents did not increase trust or cooperation. Instead, participants reduced their returns overall whilst the returns of those in the control group did not change between the pre and post phase of the experiment. Neither did participants show more forgiving behavior in the manipulation condition after the one-off defection by the agent. Why did participants reduce their returns even though they were repeatedly exposed to a more cooperative and more forgiving AI? A look at how the participants rated their co-players might shed some light on what might be driving this reduction in returns for those in the manipulation condition.

Those exposed to the forgiving agent rated their opponent in the post-exposure phase lower on all attributes even though they faced the same dynamic human-like HMM as pre-exposure. One possible explanation for this drop in rating is that participants exhibited a negative contrast effect. This occurs when the evaluation of a person, object, or situation is influenced by comparisons with recently encountered contrasting objects or people. If we’ve recently interacted with someone exceptionally nice, our perception of a normal level of niceness might be skewed, making normal behavior seem less favourable or even negative by comparison (Kobre and Lipsitt 1972). As the most recently faced opponents were more forgiving and cooperative, this negative contrast effect may have trumped any learning transfer from being repeatedly exposed to cooperative and forgiving AI (Zentall 2005). If this contrast effect is indeed replicable, then an avenue for future research would be to use it to our benefit by making the participants play agents with low cooperation perception.

It is worth highlighting that RS was not correlated to overall performance in the task. Neither did it moderate the change in returns between Conditions or have an effect on participant returns. However, in examining the player ratings by RS group more closely, we observe that individuals with high RS demonstrate a heightened attunement to changes in the behavior of their AI co-players. Specifically, when exposed to a more forgiving AI agent, these participants accurately increased their ratings of the agent, indicating a sensitive and appropriate response to the behavioral manipulation. This adjustment reflects a nuanced perception of social cues and a capacity to modify judgments based on the behavior of interaction partners, potentially indicating a perceived alignment or support from the co-player that mitigates their heightened sensitivity to potential social rebuffs. High RS participants’ lower ratings of essentially similar co-players as they continued to face them may indicate that the occasional pre-programmed defections of the AI agents failed to mitigate concerns over rejection and possibly exacerbated perceptions of uncooperativeness and unforgiveness, leading to a stronger preference against future interactions. Individuals with low RS were less affected by AI behavioral variations, showing a less volatile baseline of social perception. They still exhibit behavior consistent with a contrast effect in the manipulation condition. This dichotomy highlights the potential of tailored interactions, mediated by advanced AI agents, to address and modulate specific social and psychological predispositions in human participants.

This study leveraged a novel approach to examining the dynamics of trust and cooperation in social interactions through the utilization of generative models of human behavior to design artificial agents that can interact in economic games. The use of these agents in the RTG led to similar investment and returns to those recorded in human dyadic interactions. Participants were often uncertain whether they interacted with human or artificial investors, highlighting the agents' realism. This validates the use of these artificial agents to probe the effectiveness of manipulations whilst keeping a high degree of experimental control. Importantly, our approach utilizing HMMs to simulate interactive partners in social dilemmas offers a promising avenue for future research, especially in understanding and potentially mitigating trust deficits in individuals with high RS. The nuanced behaviors elicited through the interaction with HMM agents highlight the complexity of trust dynamics and the potential for computational models to offer novel insights into human social behaviors.

## References

- Almaatouq, Abdullah, Joshua Becker, James P. Houghton, Nicolas Paton, Duncan J. Watts, and Mark E. Whiting. 2021. "Empirica: A Virtual Lab for High-Throughput Macro-Level Experiments." *Behavior Research Methods* 53 (5): 2158–71. <https://doi.org/10.3758/s13428-020-01535-9>.
- Bandura, Albert. 1977. *Social Learning Theory*. Prentice Hall.
- Bowlby, John. 1978. "Attachment Theory and Its Therapeutic Implications." *Adolescent Psychiatry* 6: 5–33.
- Charness, Gary, Ramón Cobo-Reyes, and Natalia Jiménez. 2008. "An Investment Game with Third-Party Intervention." *Journal of Economic Behavior & Organization* 68 (1): 18–28. <https://doi.org/10.1016/j.jebo.2008.02.006>.
- Downey, Geraldine, Hala Khouri, and Scott I. Feldman. 1997. "Early Interpersonal Trauma and Later Adjustment: The Mediation Role of Rejection Sensitivity." In *Developmental Perspectives on Trauma: Theory, Research, and Intervention*, 85–114. Rochester Symposium on Developmental Psychology, Vol. 8. Rochester, NY, US: University of Rochester Press.
- Fiedler, Marina, Ernan Haruvy, and Sherry Xin Li. 2011. "Social Distance in a Virtual World Experiment." *Games and Economic Behavior* 72 (2): 400–426. <https://doi.org/10.1016/j.geb.2010.09.004>.
- Fonagy, Peter, and Elizabeth Allison. 2014. "The Role of Mentalizing and Epistemic Trust in the Therapeutic Relationship." *Psychotherapy* 51: 372–80. <https://doi.org/10.1037/a0036505>.
- Fonagy, Peter, and Chloe Campbell. 2017. "Mentalizing, Attachment and Epistemic Trust: How Psychotherapy Can Promote Resilience." *Psychiatria Hungarica: A Magyar Pszichiatriai Tarsasag Tudományos Folyóirata* 32 (3): 283–87.
- Fowler, James H., and Nicholas A. Christakis. 2010. "Cooperative Behavior Cascades in Human Social Networks." *Proceedings of the National Academy of Sciences* 107 (12): 5334–38. <https://doi.org/10.1073/pnas.0913149107>.
- Gao, Shuling, Mark Assink, Andrea Cipriani, and Kangguang Lin. 2017. "Associations Between Rejection Sensitivity and Mental Health Outcomes: A Meta-Analytic Review." *Clinical Psychology Review* 57 (November): 59–74. <https://doi.org/10.1016/j.cpr.2017.08.007>.
- Herpertz, Sabine C., and Katja Bertsch. 2014. "The Social-Cognitive Basis of Personality Disorders." *Current Opinion in Psychiatry* 27 (1): 73–77. <https://doi.org/10.1097/YCO.000000000000026>.
- Joyce, Berg, John Dickhaut, and Kevin McCabe. 1995. "Trust, Reciprocity, and Social History." *Games and Economic Behavior* 10 (1): 122–42.
- King-Casas, Brooks, Damon Tomlin, Cedric Anen, Colin F. Camerer, Steven R. Quartz, and P. Read Montague. 2005. "Getting to Know You: Reputation and Trust in a Two-Person Economic Exchange." *Science* 308 (5718): 78–83. <https://doi.org/10.1126/science.1108062>.
- Kobre, Kenneth R., and Lewis P. Lipsitt. 1972. "A Negative Contrast Effect in Newborns." *Journal of Experimental Child Psychology* 14 (1): 81–91. [https://doi.org/10.1016/0022-0965\(72\)90033-1](https://doi.org/10.1016/0022-0965(72)90033-1).
- Matuschek, Hannes, Reinhold Kliesch, Shravan Vasishth, Harald Baayen, and Douglas Bates. 2017. "Balancing Type I Error and Power in Linear Mixed Models." *Journal of Memory and Language* 94 (June): 305–15. <https://doi.org/10.1016/j.jml.2017.01.001>.
- Mikulincer, Mario. 1998. "Attachment Working Models and the Sense of Trust: An Exploration of Interaction Goals and Affect Regulation." *Journal of Personality and Social Psychology* 74 (5): 1209–24. <https://doi.org/10.1037/0022-3514.74.5.1209>.
- Mulder, R. T., P. R. Joyce, P. F. Sullivan, C. M. Bulik, and F. A. Carter. 1999. "The Relationship Among Three Models of Personality Psychopathology: DSM-III-R Personality Disorder, TCI Scores and DSQ Defences." *Psychological Medicine* 29 (4): 943–51. <https://doi.org/10.1017/S0033291799008533>.
- Singmann, Henrik, Ben Bolker, Jake Westfall, Frederik Aust, Mattan S. Ben-Shachar, Søren Højsgaard, John Fox, et al. 2022. "Afex: Analysis of Factorial Experiments."
- Zentall, Thomas R. 2005. "A Within-trial Contrast Effect and Its Implications for Several Social Psychological Phenomena." *International Journal of Comparative Psychology* 18 (4). <https://doi.org/10.46867/ijcp.2005.18.04.08>.