

So good it has to be true: Wishful thinking in theory of mind

Daniel Hawthorne-Madell¹, Noah D. Goodman¹

¹Department of Psychology, Stanford University

Keywords: wishful thinking, computational social cognition, theory of mind, desirability bias

Abstract

In standard decision theory, rational agents are objective, keeping their beliefs independent from their desires. Such agents are the basis for current computational models of Theory of Mind (ToM), but the accuracy of these models are unknown. Do people really think that others do not let their desires color their beliefs? In two experiments we test whether people think that others engage in wishful thinking. We find that participants do think others believe that desirable events are more likely to happen, and that undesirable ones are less likely to happen. However, these beliefs are not well calibrated as people do *not* let their desires influence their beliefs in the task. Whether accurate or not, thinking that others wishfully think has consequences for reasoning about them. We find one such consequence—people learn more from an informant who thinks an event will happen despite wishing it was otherwise. People’s ToM therefore appears to be more nuanced than the current rational accounts in that it allows other’s desires to directly affect their subjective probability of an event.

Whether thinking “I can change him/her” about a rocky relationship or the more benign “those clouds will blow over” when at a picnic, people’s desires seem to color their beliefs. However such an explanation presupposes a direct link between his desires and beliefs, a link that is currently absent in normative behavioral models and current Theory of Mind (ToM) models.

Does a causal link between desires and beliefs actually exist?¹ The evidence is mixed. There are a number of compelling studies that find “wishful thinking,” or a “desirability bias” in both carefully

Corresponding author: Daniel Hawthorne-Madell, d.j.hawthorne@alumni.stanford.edu

¹ While the causal link between desires and beliefs may, in fact, be bi-directional, we will focus on the evidence for the *a priori* effect of desires on beliefs.

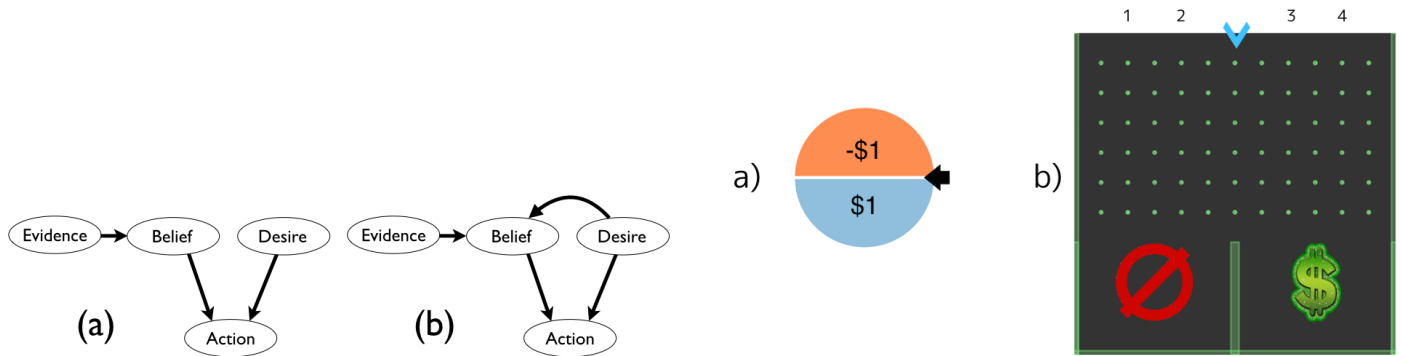
23 controlled laboratory studies (Mayraz, 2011) and real world settings, such as the behavior of sport fans
24 (Babad, 1987; Babad & Katz, 1991), expert investors (Olsen, 1997), and voters (Redlawsk, 2002).
25 However, other researchers have failed to observe the effect—for example, Bar-Hillel et al.’s *The elusive*
26 *wishful thinking effect* (1995), have provided alternative accounts of previous experiments (Hahn &
27 Harris, 2014), and have argued that there is insufficient evidence for a systematic wishful thinking bias
28 (Hahn & Harris, 2014; Krizan & Windschitl, 2007).

29 Whether or not there actually *is* a direct effect of desires on beliefs, people might *think* that there is and
30 use this fact when reasoning about other people. That is to say, people’s ToM might incorporate the
31 wishful thinking link seen in Figure 1b. The direct influence of desires on beliefs is a departure from
32 classic belief-desire “folk” psychology in which beliefs and desires are independent and jointly cause
33 action (Figure 1a). Previous models of ToM formalize belief-desire psychology into probabilistic models
34 of action and belief formation. They show that inferring others’ beliefs (Baker, Saxe, & Tenenbaum,
35 2011), preferences (Jern, Lucas, & Kemp, 2011), and desires (Baker, Saxe, & Tenenbaum, 2009) can be
36 understood as Bayesian reasoning over these generative models. A fundamental assumption of these
37 models is that beliefs are formed on the basis of evidence, and *a priori* independent of desire. We will
38 call models that make this assumption *rational theories of mind* (rToM). We can contrast this rationally
39 motivated theory with one that incorporates the rose-colored lenses of a desire-belief link, an *optimistic*
40 *ToM* (oToM).² We use their qualitative predictions to motivate two experiments into the presence (and
41 calibration) of wishful thinking in ToM and its impact on social reasoning.

49 In Experiment 1 we explore wishful thinking in both ToM and behavior. In the 3rd person point-of-view
50 (3-PoV) condition, we test whether people use a rToM or an oToM when reasoning about how others
51 play a simple game—will manipulating an agent’s desire for an outcome affect people’s judgments about
52 the agent’s belief in that outcome? In the first person point of view (1-PoV) condition we test whether
53 people *actually* exhibit wishful thinking when playing the game themselves. We carefully match the
54 (3-PoV) and (1-PoV) conditions and run them concurrently to have a clear test of whether people’s ToM
55 assumptions lead them to make appropriate inferences about people’s behavior in the game.³ Regardless
56 of its appropriateness, people’s ToM should have consequences for both how they reason about others’

² We formally describe Bayesian models of both rToM and oToM in the Supplement (citation of supplementary material here XXX).

³ Experiment 1 is a slightly modified replication of the two conditions previously run as separate experiments (see Supplement).



42 **Figure 1.** Competing models of ToM. Causal models of (a) rational ToM
 43 based upon classic belief-desire psychology and (b) optimistic ToM that in-
 44 cludes a direct “wishful thinking” link between desires and beliefs.

45 **Figure 2.** Stimuli used in Experiment 1. (a) the wheel used to determine
 46 the payout for the next outcome and (b) the Galton board used to decide the
 47 outcome. The blue arrow at the top indicates where the marble will be dropped.
 48 The numbers indicate the four drop positions used in the experiment.

57 actions and how they learn from them. If people do attribute wishful thinking to others, it would have a
 58 dramatic impact on their interpretation of others’ behavior. In Experiment 2 we therefore test for a social
 59 learning pattern that only reasoners using a oToM would exhibit, highlighting the impact ToM
 60 assumptions have on social reasoning.

EXPERIMENT 1: WISHFUL THINKING IN TOM (3-POV) AND ONLINE BEHAVIOR (1-POV)

61 **3-PoV condition**

62 To test for the presence of wishful thinking in people’s mental models of others we introduced Josh, a
 63 person playing a game with a transparent causal structure. The causal structure of the game was
 64 conveyed via the physical intuitions of the Galton board pictured in Figure 2b (in which a simulated ball
 65 bounces off pegs to land in one of two bins). The outcome of the game is binary (there are two bins) with
 66 different values associated with each outcome (money won or lost). We call the value of an outcome (i.e.,
 67 the amount that Josh stands to win or lose) the utility of that outcome, $U(outcome)$. Participants were
 68 asked what they think about Josh’s belief in the likelihood of the outcome $p_j(outcome)$. By manipulating
 69 outcome values we are able to test for wishful thinking. If people incorporate wishful thinking into their
 70 ToM, we should find that increasing an outcome’s utility results in higher estimates of Josh’s belief in the
 71 outcome’s occurrence, $p_j(outcome)$.

72 We first measured $p_j(\textit{outcome}|\textit{evidence})$ without manipulating the desirability of the outcome in the
 73 “baseline” block of trials. Then in the “utility” block of trials we assigned values to outcomes,
 74 manipulating Josh’s $U(\textit{outcome})$.⁴ In the *utility* block of trials we used a spinning wheel (Figure 2a) to
 75 determine what Josh stood to win or lose based on the outcome of the marble drop. By comparing these
 76 two blocks of trials we test for the presence of wishful thinking in people’s ToM.

77 **1-PoV condition**

78 To test whether people’s desires directly influence their beliefs in the Galton board game, we simply had
 79 the participant directly play the game (replacing Josh) and asked them about their belief in the likelihood
 80 of the outcome (their “self” belief $p_s(\textit{outcome})$).

81 *Methods*

82 *Participants*

83 80 participants were randomly assigned to either the 3-PoV or the 1-PoV condition such that there were
 84 40 in each.

85 *Design and Procedure*

86 **3-PoV condition**

87 Participants were first introduced to Josh, who was playing a marble-drop game with a Galton board (as
 88 seen in Figure 2b). Josh was personified as a stick figure and appeared on every screen. We then
 89 presented the causal structure (i.e., physics) of the game by dropping a marble from the center of the
 90 board two times, with one landing in the orange bin (Figure 2b left bin) and the one landing in the purple
 91 bin (Figure 2b right bin). After observing the two marble drops, participants began the *baseline* block of
 92 trials. In the four baseline trials, the marble’s drop position varied and participants were asked “What do
 93 you think Josh thinks is the chance that the marble lands in the bin with the purple/orange box?”
 94 Participants’ responses were recorded on a continuous slider with endpoints labeled “Certainly Will” and
 95 “Certainly Won’t.” Color placement was randomized on each trial, and the color of the box in question

⁴Crucially, Josh’s $U(\textit{outcome})$ should not be chosen by him, e.g., “I bet \$5 that it lands in the right bin,” as such an action would render $U(\textit{outcome})$ and $p(\textit{outcome})$ conditionally dependent and both rToM and oToM would predict influence of desire on belief judgments. To test pure wishful thinking, Josh’s $U(\textit{outcome})$ has to be assigned to him by a process independent of $p(\textit{outcome})$ —in our case a spinner.

varied between participants. The marble drop position was indicated with a blue arrow at the top of the Galton board, and there were four drop positions used ($marble_x$; top of Figure 2b) that varied in how likely they were to deliver the marble into the bin in question. In the baseline and subsequent trials participants did not observe the marble drop and outcome; they only observed the position the marble would be dropped from.

After the baseline trials, participants were introduced to the *utility* trials, which included a spinning wheel that determined “how much Josh can win or lose” labeled with \$1 and -\$1. At the beginning of each trial the wheel was spun and the selected payout was displayed, e.g., “Josh has a chance of winning \$1,” along with the Galton board. The bins were labeled with a \$ and \emptyset symbol.⁵ If the marble landed in the \$ bin then Josh won/lost the money. The location of the \$ bin was randomized on each trial. After seeing the Galton board with $marble_x$ indicated with a blue arrow, participants were asked two questions sequentially. First they were asked “What do you think Josh believes is the chance that the marble will land on the { $\$/-\$$ } and he’ll {win/lose} \$1” with the response recorded on the same slider as the baseline trials with endpoints labeled “Certainly Will” and “Certainly Won’t.” They were then asked “How much does Josh care about the outcome?” with the response on a slider with endpoints labeled from “Not at All” to “To a Great Extent.” Participants saw every combination of the two outcomes (\$1, -\$1) and the four drop positions (see Figure 2b) for a total of 8 utility trials.

1-PoV condition

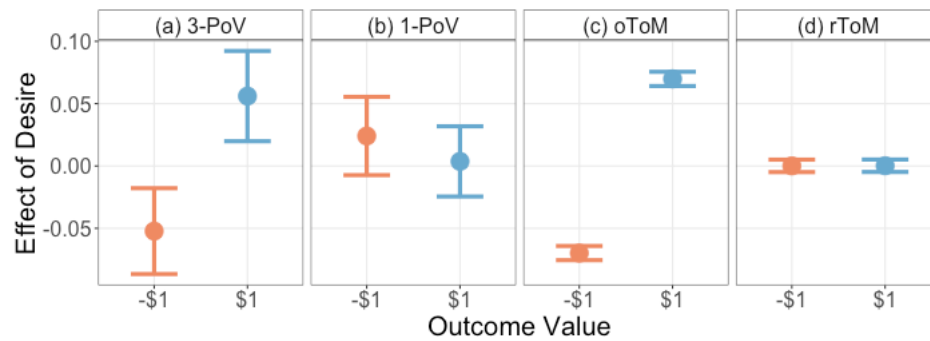
The procedure mirrored the 3-PoV condition with the participant taking the place of Josh. All questions were therefore re-framed to ask the participant’s beliefs about the outcome. The participants were given a \$1 bonus initially and instructed that one trial at random would be selected to augment their current bonus, i.e. they could gain or lose \$1.

Results

3-PoV condition

In a rational theory of mind, beliefs and desires are *a priori* independent. Manipulating Josh’s desires therefore shouldn’t have an effect on his beliefs, and we would predict that the utility trials look like the baseline trials. However, as seen in Figure 3a, the utility trials varied systematically from the baseline

⁵ \$ when the payout was positive and -\$ when it was negative, with \emptyset representing no payout.



118 **Figure 3.** Experiment 1 data. The effect of an agent’s desire for an outcome on the mean subjective $p_j(outcome)$ attributed to the agent (with 95% CIs). For
 119 each participant, the mean effect of the positive utility (\$1) and the negative utility (-\$1) was determined by taking the difference between the $p_j(outcome)$ in
 120 each utility trial and the corresponding baseline trial. The effect is shown for the (a) 3-PoV and (b) 1-PoV condition (where $p_s(outcome)$ is displayed). These
 121 data are compared with the posterior predictives of the (c) optimistic and (d) rational ToM models (see Supplement).

127 trials and, therefore, the predictions of an rToM. To quantify this deviation we fit a logistic mixed-effects
 128 model to participants’ $p_j(outcome)$ responses. The model used $marble_x$ and the categorically coded
 129 value of the outcome (negative, baseline, and positive) as fixed effects and included the random effect of
 130 $marble_x$ and intercept for each participant. The resulting model indicated that if an outcome was
 131 associated with a utility for Josh, participants thought that it would impact his beliefs about the
 132 probability of that outcome.

133 Participants thought that Josh would believe that an outcome that lost him money was less likely than the
 134 corresponding baseline trial ($\beta = -0.70, z = -2.10, p = .036$).⁶ They also thought that Josh would
 135 believe an outcome that would net him money was more likely than the corresponding baseline trial
 136 ($\beta = .96, z = 2.87, p = .004$).⁷ Finally, $marble_x$, the direct evidence, had a significant influence
 137 ($\beta = 10.37, z = 11.78, p < .001$). There was no evidence that the effect of the outcome value was
 138 affected by $marble_x$ (the interactive model did not provide a superior fit ($\chi^2(2) = .68, p = .736$)).

139 **1-PoV condition**

⁶ All p values reported for Experiment 1 are based on the asymptotic Wald test

⁷ There was no evidence of loss aversion in the relative magnitude of the wishful thinking effect for positive and negative utilities. In fact, the magnitude of the wishful thinking effect was slightly stronger for positive utilities.

140 Unlike in the 3-PoV condition, as seen in Figure 3b, there was no effect of utility on participants'
 141 $p_s(outcome)$ responses compared with their baseline responses. Using the same logistic mixed-model
 142 employed in the 3-PoV condition, neither outcomes that would lose the participant money
 143 ($\beta = .09, z = .30, p = .760$), nor outcomes that would win them money ($\beta = -.09, z = .30, p = .760$)
 144 influenced participants' $p_s(outcome)$ responses. Similar to the 3-PoV condition, a strong effect of the
 145 marble's position was observed ($\beta = 8.88, z = 11.95, p < .001$).

146 **Comparing conditions**

147 To formalize the discrepancy of the effect of utility across conditions we analyzed them together with a
 148 logistic mixed-model. We used the same model described previously except we continuously coded the
 149 effect of utility and added an interaction between this utility and condition. The resulting model had a
 150 significant interaction between PoV (condition) and the effect of utility on participants' $p_{j/s}(outcome)$
 151 responses ($\beta = .43, z = 3.83, p < .001$). This interactive model provided a better fit than the additive
 152 model ($\chi^2(1) = 15.11, p < .001$).

153 **Discussion**

154 The results from the 3-PoV condition indicate that people's ToM includes a direct "wishful thinking"
 155 link. This is consistent with the qualitative predictions of the oToM model (Supplement, Eq. 2), unlike
 156 rToM models where beliefs and desires are *a priori* independent.⁸ However, the 1-PoV condition did not
 157 find evidence that people are biased by their desires in the Galton board game. This disconnect suggests
 158 that people's attribution of wishful thinking in this situation is *miscalibrated*. That is to say that
 159 Experiment 1 represents a situation where wishful thinking is present in ToM reasoning but absent in
 160 actual behavior—people think others will behave wishfully when, in fact, they do not.

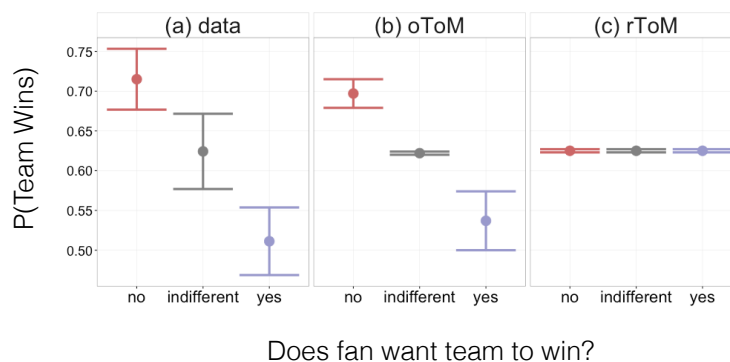
161 This miscalibration is consistent with an over-attribution of wishful thinking. However, the present study
 162 does not provide insights into *why* there is this miscalibration. Any number of incorrect assumptions
 163 could lead to the results. Perhaps people think that everyone wishfully thinks, but only they are clever
 164 enough to correct for it. Alternatively they could think that \$1 or \$5 is much more desirable for others

⁸ Interestingly, there was consistency in the magnitude of this effect when Josh stood to gain \$1 (as in the present experiment) or \$5 in Experiment 1b (see Supplement). The extent to which people attributed wishful thinking to Josh was therefore not sensitive to the magnitude of Josh's potential payout for this range (where payout is our operationalization of his desire).

165 than it is for themselves. There are a number of actor-observer asymmetries and self-enhancement biases
 166 that could plausibly underpin the observed inconsistency (Jones & Nisbett, 1971; Kunda, 1999). Further
 167 study is necessary to determine the cause of the over-attribution.

168 Regardless of whether people actually engage in wishful thinking, if people assume others do, then it
 169 should affect how they interpret others' actions and learn from them. In Experiment 2 we therefore
 170 expand our sights to social learning situations where oToM (but, crucially, not rToM) predicts that desires
 171 affect a social source's influence.

EXPERIMENT 2: LEARNING FROM OTHERS WITH AN OTOM



172 **Figure 4.** Experiment 2 data. Effect of a social sources' desire on how others learn from them for (a) data with 95% CIs, which we compare to the posterior
 173 predictives of (b) an optimistic ToM, (c) a rational ToM. Points represent the mean $p(team_x)$ response after hearing equally knowledgeable sources place a
 174 bet on $team_x$ that is either consistent, unrelated, or inconsistent with their desires.

175 Do people consider a social source's desires when learning from them? It would be important to do so if
 176 they think that his desires have a direct influence on his beliefs. Consider a learner using an oToM to
 177 reason about her uncle, a Cubs fan, who proudly proclaims that this is the year the Cubs will win it all.
 178 Though her uncle knows a lot about baseball, the oToM learner is unmoved from her (understandably)
 179 skeptical stance. However, if her aunt, a lifelong White Sox fan (hometown rival to the Cubs), agrees that
 180 the Cubs do look better than the Sox this year, then an oToM learner considers this a much stronger
 181 teaching signal. In fact, a learner with an oToM would consider her aunt's testimony as more persuasive

182 than an impartial source, see Figure 4b. A learner reasoning with an rToM wouldn't distinguish between
183 these three social sources⁹ as seen in Figure 4c.

184 We investigated which ToM best describes learning from social sources in a controlled version of this
185 biased opinion scenario. Participants were asked how likely a team (x) was to win an upcoming match,
186 $p(team_x)$, in a fictional college soccer tournament after seeing a knowledgeable student bet on the team.
187 The student was either a fan of one of the teams facing off, or indifferent to the outcome. Participants
188 therefore saw three trials—the *consistent trial* where the student bet on the team he wanted to win, the
189 *inconsistent trial* where he bet on the team wished would lose, and the *impartial trial* where he didn't
190 care which team won before he bet.

191 **Methods**

192 *Participants*

193 120 participants were randomly assigned into the consistent, inconsistent, or impartial conditions.

194 *Design and Procedure*

195 Participants were first introduced to a (fictional) annual British collegiate soccer tournament and told that
196 they would see bets on these matches from a student who “Unbeknownst to his friends makes a £100 bet
197 online on which team he thinks will win this year's game.”¹⁰ The student would either be a fan of one of
198 the teams (attending that college) or neither of the teams (attending a different college). The students
199 were equally knowledgeable across conditions, being described as seeing the outcome of the last 10
200 matches these teams played against each other.

201 After the introduction, participants were given a test trial appropriate for their (randomly assigned)
202 condition in which the student bet consistently with his school, bet against his school, or was impartial
203 (not a fan of either school). After observing the student's bet and allegiance participants were asked
204 “What do you think is the chance that $team_x$ wins the match this year?”

⁹ Assuming that the three sources are equally knowledgeable and their statements have no causal influence on the game, e.g., if the uncle is an umpire, his desires may matter through more objective routes.

¹⁰ See Supplement for complete experimental materials

205 **Results**

206 As seen in Figure 4a, participants' responses were sensitive to the student's *a priori* desires, consistent
 207 with learners who reason with oToM (but not an rToM). Participants who saw an impartial student bet on
 208 $team_x$ thought the team was more likely to win than when they saw a fan of $team_x$ place an identical bet
 209 ($d = .80$, 95% CI [.33 1.27], $z = 3.35$, $p < .001$ ¹¹). This is consistent with the learner thinking that the
 210 fan's desire to see his team win made him think it was objectively more likely. Additionally, participants
 211 who saw a fan of the other team bet on $team_x$ were *more* influenced than the same bet from the impartial
 212 student ($d = .67$, 95% CI [.21 1.14], $z = 2.87$, $p = .004$). As predicted by the model of the oToM learner,
 213 someone who bets against their desires is more diagnostic of $team_x$ being dominant than the independent
 214 source. The oToM learner thinks that $team_x$ had to be clearly dominant to overcome the wishful thinking
 215 of a fan rooting against them.

216 **Discussion**

217 Assuming that fans engage in wishful thinking allows oToM learners to make *stronger* inferences about
 218 the strength of the fans' evidence in some cases. For an rToM learner, the fan would have to have seen
 219 $team_x$ win a majority of the 10 observed matches in order to bet on them, regardless of their
 220 predilections, resulting in the flat predictions seen in Figure 4c. Meanwhile, the oToM learner thinks that
 221 a fan of $team_x$ could bet on them even if the fan only observed them win a few times¹². If, however, the
 222 fan bets against their team, the oToM learner assumes that the fan must have seen their team trounced in
 223 the 10 observed matches. Using these insights, an oToM learner using Bayesian inference to learn from
 224 the fan will exhibit the qualitative pattern seen in Figure 4b, which is consistent with participants'
 225 behavior (as seen in Figure 4a). The pattern of results is consistent with the predictions of a learner using
 226 an oToM, but there are limitations and additional potential explanations discussed at length in the
 227 Supplement.

GENERAL DISCUSSION

¹¹ calculated with Fisher-Pitman permutation test

¹² In fact, if the oToM learner thinks that the fan is a completely wishful thinker then his bet is no longer diagnostic of his evidence (he could have seen anything!)

228 Current computational models of theory of mind are built upon the assumption that beliefs are *a priori*
229 independent of desires. Whether social reasoners use such a rational ToM (rToM) is an empirical
230 question. In two experiments we tested the independence of beliefs and desires in ToM and found that
231 people behave as if they think that others are wishful thinkers whose beliefs are colored by their desires.

232 In the 3-PoV condition of Experiment 1 we found that people believe that others inflate the probability of
233 desirable outcomes and underestimate the probability of undesirable ones, as they would if they have an
234 optimistic ToM (oToM) with a direct link between desires and beliefs (Figure 3). If people broadly
235 attribute wishful thinking to others (as Experiment 1 suggests) it should be reflected in their social
236 reasoning. For example, social learners using an oToM to make sense of an agent’s beliefs would be
237 sensitive to that agent’s relevant desires. This is exactly what we found in Experiment 2 (Figure 4)—how
238 much people learned from an agent’s beliefs depended on his desires. Agents whose beliefs ran against
239 their desires were more influential than impartial agents, who in turn were more influential than agents
240 with consistent beliefs and desires.

241 The observed presence of wishful thinking in ToM has no necessary relation to its existence in people’s
242 “online” belief formation. Indeed, the 1-PoV conditions of Experiment 1 indicates that people’s model of
243 other’s wishful thinking is not perfectly calibrated. They over-attribute wishful thinking to others in
244 situations where they would actually form their beliefs independently of their desires. Charting the
245 situations where wishful thinking is over-applied in this way may be a fruitful avenue for further
246 research. At the extreme we could imagine finding that everyone thinks one another wishfully thinks, but
247 in fact everyone forms their beliefs independent of their desires! This radical thesis is surely too strong,¹³
248 but oToM may well overestimate the strength of wishful thinking and over-generalize it—amplifying a
249 small online effect into a larger social cognition effect. Attention to whether a task engages (potentially
250 amplified) oToM representations could provide insight into the considerable heterogeneity of the wishful
251 thinking effect as it has been studied. Specifically, it could help explain why first-person wishful thinking
252 is reliably found in some paradigms and not others.

253 The paradigms in which wishful thinking is reliably found involve participants *reasoning about*
254 *themselves or others*, such as the 3-PoV condition of Experiment 1 where participants reasoned about

¹³ As seen in well controlled examples of desires influencing online belief formation e.g., (Mayraz, 2011).

255 Josh’s beliefs (for a review of many tasks that may engage social reasoning see [Shepperd, Klein, Waters,
256 & Weinstein, 2013], e.g., [Weinstein, 1980] but see [Harris & Hahn, 2011] and [Hahn & Harris, 2014] for
257 an alternative explanation). Whereas *asocial paradigms* involving direct estimation of probabilities
258 usually do not find the effect, like the 1-PoV condition of Experiment 1 where participants directly
259 estimated the chance that the ball would fall into a particular bin (for other examples of wishful thinking
260 paradigms that do not involve social reasoning see Study 1 of [Bar-Hillel & Budescu, 1995], and for a
261 more general review of asocial bias experiments see the “bookbags” and “pokerchips” paradigms cited in
262 [Hahn & Harris, 2014], but see Francis Irwin’s series of experiments for an example of asocial paradigms
263 that do find a wishful thinking effect—starting with [Irwin, 1953]).

264 Where people’s predictions of others’ behaviors (1-PoV Experiment 1) and their actual behavior (3-PoV
265 Experiment 1) diverge is also important to map because these disconnects inject a systematic bias in
266 social reasoning. Taking the social learning of Experiment 2 as an example, oToM learners ignored the
267 belief of the agent whose bet was consistent with his desires. However, if this agent actually formed his
268 beliefs without bias, then the learner would be missing a valuable learning opportunity. Asserting that
269 others let their desires cloud their beliefs allows people to “explain away” those beliefs without seriously
270 considering the possible evidence on which they are based. Future work should explore the details of
271 these effects. For example, does a learner attribute bias equally to those who share his desires and those
272 who hold competing ones?

273 The experiments presented here suggest that people think that others are wishful thinkers; this has broad
274 consequences for social reasoning ranging from our inferences about heated scientific debates to
275 pundit-posturing. Our findings highlight the importance of further research into the true structure of
276 theory of mind. Do people think that others exhibit loss aversion or overweight low probabilities? Is the
277 connection between beliefs and desires bi-directional? Rigorous examination of questions like these may
278 buttress new, empirically motivated computational models of ToM that capture the nuance of human
279 social cognition—an idea so good it has to be true.

REFERENCES

- 280 Babad, E. (1987). Wishful thinking and objectivity among sports fans. *Social Behaviour*, 2(23), 231–240.
- 281 Babad, E., & Katz, Y. (1991). Wishful thinking—Against all odds. *Journal of Applied Social Psychology*, 21, 1921–1938.

- 282 Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition*, *113*(3), 329–349.
- 283 Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2011). Bayesian theory of mind: Modeling joint belief-desire attribution. In
284 L. Carlson (Ed.), *In proceedings of the thirtieth third annual conference of the cognitive science society* (pp.
285 2469–2474). Cognitive Science Society.
- 286 Bar-Hillel, M., & Budescu, D. (1995). The elusive wishful thinking effect. *Thinking and Reasoning*, *1*, 71–104.
- 287 Hahn, U., & Harris, A. J. L. (2014). What does it mean to be biased: Motivated reasoning and rationality. *Psychology of*
288 *Learning and Motivation*, *61*, 41–102.
- 289 Harris, A. J., & Hahn, U. (2011, Jan). Unrealistic optimism about future life events: a cautionary note. *Psychol Rev*, *118*(1),
290 135–154.
- 291 Irwin, F. W. (1953). Stated expectations as functions of probability and desirability of outcomes. *Journal of Personality*,
292 *21*(3), 329–335. doi: [10.1111/j.1467-6494.1953.tb01775.x](https://doi.org/10.1111/j.1467-6494.1953.tb01775.x)
- 293 Jern, A., Lucas, C. G., & Kemp, C. (2011). Evaluating the inverse decision-making approach to preference learning. In
294 J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, & K. Q. Weinberger (Eds.), *Advances in neural information*
295 *processing systems 24* (pp. 2276–2284). Curran Associates, Inc.
- 296 Jones, E. E., & Nisbett, R. E. (1971). *The actor and the observer: divergent perceptions of the causes of behavior*. New York:
297 General Learning Press.
- 298 Krizan, Z., & Windschitl, P. D. (2007). The influence of outcome desirability on optimism. *Psychological Bulletin*, *133*(1),
299 95–121.
- 300 Kunda, Z. (1999). *Social cognition: making sense of people*. Cambridge, MA: MIT Press.
- 301 Mayraz, G. (2011). *Wishful Thinking*. CEP Discussion Paper. London: Centre for Economic Performance, London School of
302 Economics.
- 303 Olsen, R. A. (1997). Desirability bias among professional investment managers: Some evidence from experts. *Journal of*
304 *Behavioral Decision Making*, *10*(1), 65–72.
- 305 Redlawsk, D. P. (2002, November). Hot cognition or cool consideration? Testing the effects of motivated reasoning on
306 political decision making. *The Journal of Politics*, *64*(04), 1021–1044.
- 307 Shepperd, J. A., Klein, W. M. P., Waters, E. A., & Weinstein, N. D. (2013). Taking stock of unrealistic optimism. *Perspectives*
308 *on Psychological Science*, *8*(4), 395–411.
- 309 Weinstein, N. D. (1980). Unrealistic optimism about future life events. *Journal of Personality and Social Psychology*, *39*(5),
310 806–820. doi: [10.1037/0022-3514.39.5.806](https://doi.org/10.1037/0022-3514.39.5.806)