

# Socially smart software agents entice people to use higher-order theory of mind in the Mod game

Kim Veltman<sup>1</sup>, Harmen de Weerd<sup>1,2</sup>, and Rineke Verbrugge<sup>1</sup>

<sup>1</sup> Institute of Artificial Intelligence, University of Groningen

<sup>2</sup> Professorship User Centered Design, Hanze University of Applied Sciences

**Abstract.** In social settings, people often need to reason about unobservable mental content of other people, such as their beliefs, goals, or intentions. This ability helps them to understand, to predict, and even to influence the behavior of others. People can take this ability further by applying it recursively. For example, they use second-order theory of mind to reason about the way others use theory of mind, as in ‘Alice believes that Bob does not know about the surprise party’. However, empirical evidence so far suggests that people do not spontaneously use higher-order theory of mind in strategic games. Previous agent-based modeling simulations also suggest that the ability to recursively apply theory of mind may be especially effective in competitive settings. In this paper, we use a combination of computational agents and Bayesian model selection to determine to what extent people make use of higher-order theory of mind reasoning in a particular competitive game, the Mod game, which can be seen as a much larger variant of the well-known rock-paper-scissors game.

We let participants play the competitive Mod game against computational theory of mind agents. We find that people adapt their level of theory of mind to that of their software opponent. Surprisingly, knowingly playing against second- and third-order theory of mind agents entices human participants to apply up to fourth-order theory of mind themselves, thereby improving their results in the Mod game. This phenomenon contrasts with earlier experiments about other strategic one-shot and sequential games, in which human players only displayed lower orders of theory of mind.

## 1 Introduction

Theory of mind, the ability to reason about unobservable mental content such as beliefs and goals of others [17], plays an important role in many social interactions. Indeed, many researchers have claimed that a higher proficiency in theory of mind is related with pro-social behavior [11], social competences [13], and negotiation skills [24]. But while some experiments show impressive theory of mind skills in adults, people are typically slow to take advantage of their theory of mind ability in strategic settings [10,2,25,9]. In this paper, we let participants

interact with artificial theory of mind agents to determine to what extent these agents can encourage and train people in their use of theory of mind.

The human ability for theory of mind is not limited to reasoning about the goals, desires, and beliefs concerning world facts of others. People can also use their theory of mind to reason about the way others use theory of mind. For example, people use *second-order theory of mind* to understand a sentence such as “Alice *knows* that Bob *knows* that Carol is throwing him a birthday party”, by reasoning about what Alice knows about what Bob knows. Experimental results from story comprehension tasks show that people are able to use orders of theory of mind beyond second-order theory of mind as well. In these tasks, adults perform much better than chance on story comprehension questions that explicitly involve theory of mind reasoning up to the fourth order [12,19].

In this article, we let participants play the Mod game against artificial theory of mind agents, and estimate their level of theory of mind reasoning using random-effects Bayesian model selection [18]. Simulation studies show that the ability to make use of higher-order theory of mind can be particularly effective in competitive settings [8,23,4,6]. In particular, results from agent simulations in a variant of the Mod game suggest that both first-order and second-order theory of mind can greatly benefit players, while the use of orders of theory of mind beyond the second hardly provide additional benefits [23]. However, results in the matching pennies game show that in this simple competitive game, many people rely on simpler, behavior-based strategies when engaging with artificial agents [20]. In this paper, we consider an extension of the matching pennies game known as the Mod game. The Mod game has a structure that is similar to that of matching pennies, but due to the larger number of possible actions, the Mod game should make it easier for participants to distinguish between strategies. This may help people to detect the use of theory of mind in their opponent, and encourage them to make use of higher orders of theory of mind themselves.

For this research, we let human participants play against virtual agents that we developed to determine the effectiveness of making use of increasingly higher orders of theory of mind [22]. This allows us to precisely control and monitor the mental content of the opponents faced by the participants, including their application of theory of mind. By letting the participants play against virtual agents, we can analyze the participant data in a more controlled setting. This allows us to diagnose the human behavior. But that’s not the only benefit of using virtual agents: By letting human participants train with virtual agents that are programmed to reason in a certain fashion, we can stimulate people to improve their own reasoning. For example, when a participant plays against a higher-order theory of mind agent, which makes it harder for the participant to win the game, this might be an incentive for the participant to actually try harder and also apply higher-order theory of mind, since the agent is consistent in using the higher orders of theory of mind.

The remainder of this paper is structured as follows. Section 2 and Section 3 introduce the Mod game and a range of strategies that agents and humans could use. In Section 4, we describe how random-effects Bayesian model selection may

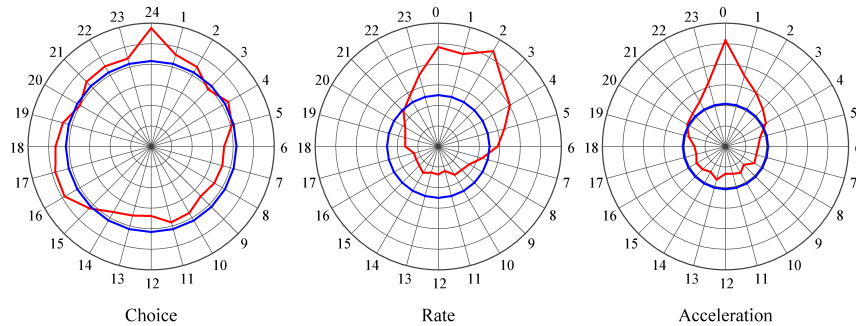


Fig. 1: Histograms over 24 choices, rates, and accelerations of human behaviour in the Mod game. In each graph, the blue curve shows the expected results from random behaviour, while the red curve shows the participant behaviour (reconstructed from [7]).

help to gauge agents’ and participants’ reasoning strategies from their behavior. Section 5 delineates our experiment in which human participants played the Mod game against agents that used different orders of theory of mind. In Section 6, the experimental results are presented and in Section 7 we close the article and conclude that virtual agents can indeed be used to support people in using higher orders of theory of mind in a competitive game such as Mod.

## 2 Mod Game

The Mod game is an  $n$ -player generalization of rock-paper-scissors [23]. It has been introduced by Frey and Goldstone [8] as a way to reveal patterns in individual theory of mind strategies. In the Mod game,  $n$  players simultaneously choose a number in the range  $\{1, \dots, m\}$ , for  $m > n > 1$ . Players gain one point for every opponent that has chosen the number that is exactly one lower than their own choice. For example, a player that has chosen the number 4 gains a point for every opponent that has chosen number 3. The only exception to this rule is that players that have chosen number 1 gain one point for every player that has chosen number  $m$ , hence the name ‘mod’ game; for the goal ‘ $+1 \bmod m$ ’. To visualize the rules of the game to human participants, actions are arranged in a circle (see Figure 2 for  $m = 24$ ).

Note that in the Mod game, each action is dominated by some other actions, similar to games such as rock-paper-scissors. In fact, the Mod game is equivalent to a non-zero-sum version of rock-paper-scissors for  $n = 2$  and  $m = 3$ . The Mod game has a mixed-strategy Nash equilibrium in which each action is chosen with equal probability. When all players play according to this randomizing strategy, none of the players has an incentive to change his or her strategy.

However, participant behavior in repeated Mod games deviates from the Nash equilibrium, as depicted in Figure 1 (reconstructed from [7]). The figure shows the aggregate participant data (red line in left graph) and the idealized randomizing behavior (blue line) over 100 rounds of play and  $m = 24$ . Participant choices (red line in left-most graph) appear to be approximately random, with a slight bias towards 24. However, *participant rates*, defined as the difference in choice between two subsequent rounds, shows a clear deviation from the Nash equilibrium. As Figure 1 shows, participants (red line in the middle graph) are less likely to select numbers that are 7 to 21 ahead of their previous choice. Instead, they are most likely to choose a number that is 0 to 4 higher than their previous choice. If participants were to play according to the Nash equilibrium, however, each rate should be equally likely (blue line). *Participant acceleration*, which is defined as the change in participant rate, shows a similar effect. Figure 1 (red line in right-most graph) shows that participants tend to vary little in their rate. That is, a participant who chose a number in the last round that was 2 higher than the number in the round before that is mostly likely to choose the number that is 2 higher in the current round than his choice in the previous round. In addition, Figure 1 also shows that participants that vary their acceleration do so by a small amount. Again, Nash equilibrium play would result in each acceleration being equally likely (blue line).

In our experiment, participants play a specific variant of the Mod game with  $m = 24$  actions and  $n = 2$  player. We follow [7] in this choice, as choosing  $m$  relatively large allows us to better distinguish reasoning strategies behind participants' behavior. For the remainder of the paper, we will only consider this specific variant of the Mod game.

### 3 Strategies in the Mod game

The Mod game, as outlined in Section 2, can be played using a variety of strategies. In this section, we describe a number of these strategies. These include strategies based on the use of theory of mind, as well as simple behavior-based strategies. Table 1 shows several possible Mod game strategies, which we selected for the following reasons. The *ToM* strategies correspond to the main focus of this research. The Win-Stay-Lose-Shift strategy appears to be used often and has been selected because of its popularity.

The self- and other-regarding strategies were chosen to see whether participants may use a simpler strategy, based on either their own or the opponent's recent behavior. It seems that some participant data could be explained purely by choosing a rate (differences between the players' previous and current number) by which they change their number based on either their own or the opponent's previous choice. These self- and other-regarding strategies include the strategies where a participant consistently chooses 1, 2, 3 or 4 higher than he/she or the opponent did in the previous round. In this section, we describe these strategies in detail. To avoid confusion, in the remainder, we will refer to focal agents as if they were male, while we will refer to their opponents as if they were female.

Table 1: We consider eight possible strategies for playing the Mod game, including four behavior-based strategies and four theory of mind strategies.

Strategies	Behavior-based	ToM
Other regarding	✓	
Self regarding	✓	
Win-Stay-Lose-Shift	✓	
$ToM_0$	✓	
$ToM_1$		✓
$ToM_2$		✓
$ToM_3$		✓
$ToM_4$		✓

### 3.1 Behavior-based strategies

While participants may benefit from using theory of mind in the Mod game, the Mod game can be played using strategies that do not model unobservable mental content of others. In our Bayesian RFX-BMS analysis, we therefore also consider a number of behavior-based strategies. A player that uses a behavior-based strategy responds to actions observed in the last round of play only.

**Self-regarding strategy** An agent that follows a *self-regarding* strategy reacts to the action it has performed in the previous round, while it ignores all actions of the opponent. The self-regarding strategy depends on two free parameters. The drift parameter  $k$  determines the amount of drift a player experiences in every round, so that an agent that follows a self-regarding strategy with drift  $k$  tends to choose the number that is  $k$  higher (modulo 24) than the action he performed in the previous round.

In addition, the choice probability  $p$  determines the strength of the self-regarding tendency. For example, an agent that follows a self-regarding strategy with  $k = 2$  selects the action that is 2 higher than his previous choice with probability  $p$ , while each other action has a probability  $\frac{1}{23}(1-p)$  of being selected.

**Other-regarding strategy** The *other-regarding* strategy is similar to the self-regarding strategy, except that an agent that follows the other-regarding strategy reacts to the previous action of his opponent rather than his own previous action. Like the self-regarding strategy, the other-regarding strategy relies on a drift parameter  $k$  and choice probability  $p$ . An agent that follows an other-regarding strategy with drift  $k$  selects the action that is exactly  $k$  higher (modulo 24) than his opponent’s action in the previous round with probability  $p$ . Each other action is selected with probability  $\frac{1}{23}(1-p)$ . Note that if an agent plays according to an other-regarding strategy with  $k = 1$ , the agent tends to play the action that would have won in the previous round.

**Win-Stay, Lose-Shift (WSLS) strategy** An agent that follows the *win-stay, lose-shift (WSLS)* strategy bases his current decision on the outcome of the

previous round. If the agent won the previous round, he will repeat his previously chosen action with probability  $p$ , while each of the other 23 action is selected with probability  $\frac{1}{23}(1-p)$ . However, if the agent did not win the previous round, he will repeat his previously chosen action with probability  $1-p$ , while each of the other 23 action is selected with probability  $p/23$ . The single parameter  $p$  is a free parameter.

### 3.2 Theory of mind strategies

In addition to behavior-based strategies, we include agents that are capable of using theory of mind while playing the Mod game. Each additional order of theory of mind provides an agent with an additional prediction of his opponent’s behavior, so that a  $ToM_k$  agent has  $k+1$  models of opponent behavior [23]. These agents have been previously used to determine the extent of participants’ use of higher-order theory of mind in the simpler matching pennies setting [20]. Below, we briefly describe these agents. A full mathematical model of these agents can be found in [23].

**Zero-order theory of mind** A zero-order theory of mind ( $ToM_0$ ) agent has no theory of mind at all, and is therefore unable to attribute mental content to others. In particular, a  $ToM_0$  agent cannot form the belief that his opponent is trying to obtain a high score. Instead, the  $ToM_0$  agent forms zero-order beliefs about the actions the opponent will play in future rounds of the game.

In our agent model, a  $ToM_0$  agent forms beliefs  $b^{(0)}$  about the actions of the opponent. For each number  $i = 1, \dots, 24$ , the  $ToM_0$  agent has a belief  $b^{(0)}(i)$  that represents what he believes to be the likelihood that most of his opponents will select to play that number. The  $ToM_0$  agent acts on these beliefs by choosing the number that maximizes his score. For example, if a  $ToM_0$  agent strongly believes that number 4 will be selected by his opponent, the agent should choose to play number 5 himself.

After every round, the  $ToM_0$  agent updates his zero-order beliefs to reflect the actual outcome. An agent-specific learning speed  $\lambda \in [0, 1]$  determines the relative influence of the current observation on the agent’s beliefs. For example, a  $ToM_0$  agent with zero learning speed ( $\lambda = 0$ ) does not update his beliefs at all. Such an agent selects the same action in every round. A  $ToM_1$  agent with the maximal learning speed ( $\lambda = 1$ ), on the other hand, completely replaces his zero-order beliefs after each observation, and forgets all information obtained from previous rounds <sup>3</sup>.

To account for small deviations between participant choices and the  $ToM_0$  agent strategy, we make use of the so-called ‘softmax’ probabilistic policy [3,20]. That is, in addition to the learning speed  $\lambda$ , the  $ToM_0$  agent strategy has an addition parameter  $\beta$  that controls the magnitude of behavioral noise.

---

<sup>3</sup> Note that a  $ToM_0$  agent with learning speed  $\lambda = 1$  is identical to an agent following an other-regarding strategy with  $k = 1$ .

**First-order theory of mind** Unlike the  $ToM_0$  agent, a first-order theory of mind ( $ToM_1$ ) agent reasons about the goals of others, and therefore believes that his opponent may be trying to maximize her score as well. To predict the behavior of his opponent, the  $ToM_1$  agent attributes his own thought process to her. That is, a  $ToM_1$  agent considers the possibility that, while the agent reacts to the actions of his opponent, the opponent is reacting to the actions of the agent. For example, if the agent has played 4 in the previous round, he believes that his opponent is more likely to play 5. After all, that is what the  $ToM_1$  agent would do himself if he has observed his opponent playing 4 in the previous round.

Although the  $ToM_1$  agent models his opponent as being able to use zero-order theory of mind, agents in our setup do not know the extent of the abilities of their opponent for certain. Rather, a  $ToM_1$  agent has two models of opponent behavior, one based on zero-order theory of mind and one on first-order theory of mind. Through repeated interaction, a  $ToM_1$  agent learns which of his models best describes the behavior of his opponent. Based on this information, a  $ToM_1$  agent may therefore choose to play as if he were a  $ToM_0$  agent, and ignore the predictions of his first-order theory of mind.

Note that the  $ToM_1$  agent strategy does not introduce any new free parameters. Like the  $ToM_0$  agent strategy, there are two parameters: the learning speed  $\lambda$  and the magnitude of behavioral noise  $\beta$ .

**Higher orders of theory of mind** For each additional order of theory of mind  $k$ , an agent generates an additional prediction of opponent behavior, by attributing his own  $(k - 1)$ st-order theory of mind thought process to his opponent. For example, a  $ToM_2$  agent models his opponents as  $ToM_1$  agents, in addition to his zero-order and first-order theory of mind models of opponent behavior. As a result, a  $ToM_k$  agent has  $k + 1$  hypotheses for the action that will be chosen by his opponent, with corresponding predictions. Based on the accuracy of these predictions, the  $ToM_k$  agent can therefore choose to behave according to  $k + 1$  patterns of behavior.

## 4 Random-effects Bayesian model selection

In this paper, we attempt to encourage participants in their use of theory of mind through interactions with artificial theory of mind agents. To determine what level of theory of mind reasoning a participant is engaging in at different points throughout the experiment, we make use of a technique known as group-level random-effects Bayesian model selection (RFX-BMS), introduced by Stephan and colleagues [18]. Unlike fixed-effects Bayesian model selection, which assumes that there is one strategy that best describes the actions of all participants, random-effects Bayesian model selection models participants as individuals who may differ in the strategy they use while playing the game. Strategies are treated as random effects that occur with an unknown but fixed probability

in the population. A group of participants represents a random sample drawn from these strategies.

Random-effects Bayesian model selection estimates what distribution of strategies best fits the experimental data. Each strategy  $s$  generates pieces of evidence  $p(y|s)$  representing the probability that choosing actions according to strategy  $s$  will result in some observed data  $y$ . Using these model evidences and the participant data, random-effects Bayesian model selection estimates the relative frequencies of strategies in the general population.

To determine to what extent a participant makes use of theory of mind while playing the Mod game, we compare the observed behavior of each participant with the predicted behavior of computational agents following the strategies described in Section 3. That is, the model evidence  $p(y|s)$  generated by a given model is the probability that that model will perform the same action as the participant, given the history of moves observed by the participant. This combination of our theory of mind agents and RFX-BMS has been previously used to accurately recover the level of theory of mind reasoning of Devaine’s Bayesian theory of mind agents [3], which indicates that this method can overcome biases in the designer’s choice of how to implement theory of mind.

## 5 Experimental Setup

To determine whether interacting with artificial theory of mind agents encourages the use of theory of mind, we let human participants play the Mod game against artificial theory of mind agents. Participants played the two-player Mod game with 24 actions, as described in Section 2.

### 5.1 Participants

Sixteen participants were included in this study, of which eight were male and eight were female, all students and all over the age of 18 ( $M = 21.5$ ,  $SD = 2.3$ ). The experiment was conducted in English, and all participants were sufficiently skilled in reading and understanding the English language.

Before starting with the experiment, all participants gave informed consent about partaking and about the use of the data obtained by the experiment for the purpose of this study.

### 5.2 Experimental design

Each participant played the Mod game against four different opponents: a  $ToM_1$  agent, a  $ToM_2$  agent, a  $ToM_3$  agent, and an agent whose order was randomized each round. This randomizing agent would randomly select to respond as if it were a first-order, second-order, or third-order theory of mind agent in each round. That is, during a block of twenty rounds, the randomizing agent would randomly select a reasoning strategy twenty times.



Note that participants never played against a  $ToM_0$  agent. During a pilot study, we discovered that  $ToM_0$  agents and  $ToM_1$  agents exhibit the same behavior when playing the Mod game against a human participant. The  $ToM_0$  agent believes that the best predictor for a participant’s future behavior is this/her behavior in the previous round. As a result, the  $ToM_0$  agent tends to select the number that is 1 higher than the number last chosen by the participant.

The  $ToM_1$  agent, on the other hand, believes that the opponent wants to win the game. By taking the perspective of the opponent, the  $ToM_1$  agent believes that the opponent will choose the number that is 1 higher than his own last choice. For example, suppose that the agent chose 23 in the last round and the participant played 24. In this case, the  $ToM_1$  agent will believe that the opponent is going to play 24 again, since the  $ToM_1$  agent believes that the participant is a  $ToM_0$  agent who believes that the agent is going to play 23 again. Following this reasoning, the agent decides to play 1, which is exactly 1 higher than the participant’s previous choice (24). However, this yields the same result as the  $ToM_0$  agent, who also chose to play the number that was 1 higher than the participant’s previous choice. Due to this effect, we decided not to include the  $ToM_0$  agent in our experiment.

Each block consisted of twenty rounds of the Mod game, in which participants played against the same opponent for all twenty rounds. Between blocks, the opponent was changed to a  $ToM$  agent of a different order. The order in which participants faced the different opponents was randomly drawn from four possible sequences:  $[?,1,2,3]$ ;  $[3,?,1,2]$ ;  $[2,1,3,?]$ ; and  $[1,?,3,2]$ , where the question mark (?) represents a randomizing agent, whose order of theory of mind reasoning was randomized each round. The different sequences were chosen to rule out the effect of sequence on the performance. Participants were informed about the  $ToM$  order of the opponent they were playing against, except in the blocks in which they faced the randomizing agent. During the rounds against the randomizing agent, the order of the opponent was not shown to the participants, in order to see if the participants’ behavior also changed if the  $ToM$  order of the opponent was not known.

Each participant faced the four different opponents in two of the aforementioned sequences, so every participant played 2 repetitions each of 4 different opponents. In total, every participant played 40 rounds against every opponent. A certain opponent was played against for twenty rounds before the agent’s order changed. The number of twenty rounds per opponent was chosen because people typically need many trials before showing higher-order reasoning behavior [9].

### 5.3 Procedure and materials

Participants first read some short information about theory of mind and the different orders that were used in this experiment ( $ToM_1$ ,  $ToM_2$ , and  $ToM_3$ ). Note that while this procedure may have primed participants to make use of theory of mind strategies, evidence from Marble Drop experiments shows that even in this case, participants may have difficulties implementing higher-order theory of mind strategies [10,16]. Participants then read an explanation of the

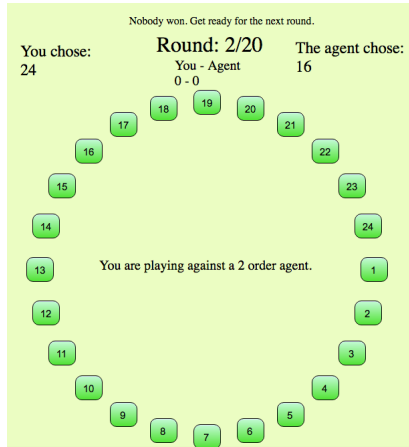


Fig. 2: Interface of the Mod24 Game experiment.

experiment itself, including the rules of the Mod game and an explanation of the interface. Before the experimental rounds started, the participants completed three test rounds to confirm they understood the interface. As the participants started with the experimental trials, they saw an interface with twenty-four buttons, numbered from 1 to 24, placed in a circle (see Figure 2). The placement of the buttons was constant throughout the whole experiment. The interface also showed what level of theory of mind the agent used (except during the randomizing agent blocks). The participants could also see how many rounds they had already played, how many rounds they would play against the same opponent, the current score of both players, and the chosen actions of both players in the previous round of play.

At the end of each block, participants were informed that they would continue playing against a new opponent. After four blocks, participants could take a break before continuing with the next four blocks. Once all eight blocks were finished, another pop-up was shown, informing the participants that the experiment was finished. Upon finishing the experiment, the participants were thanked for their cooperation and received payment. Each participant was equally compensated for their effort: the reward was not dependent on the points obtained during the experiment.

## 6 Results

In this section, we present the results of the experiment described in Section 5. To determine to what extent participants made use higher-order theory of mind reasoning, we applied RFX-BMS analysis (see Section 4) to the choices of both the participants as well as their computer opponents.

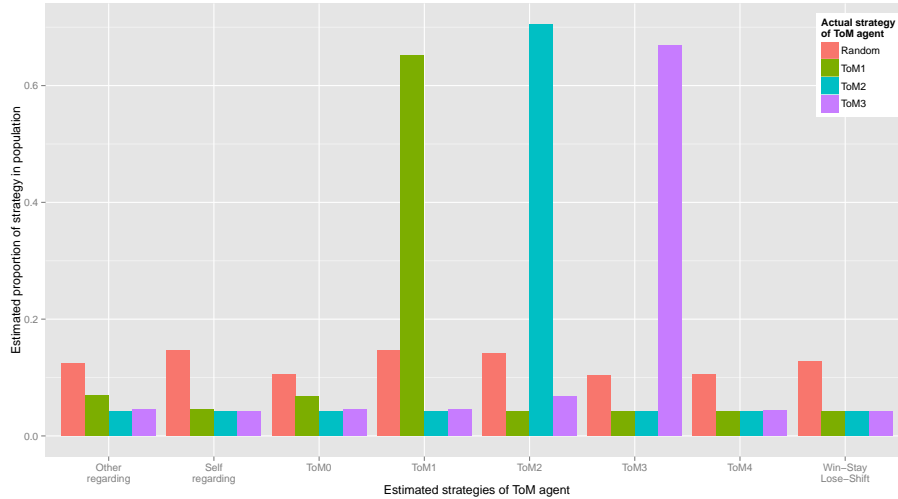


Fig. 3: Estimated strategies of the agents based on their choices throughout the experiment.

### 6.1 Agent behavior

Figure 3 shows the estimated strategies of all four agents in our experiment. By analyzing the choices of the computer opponents, we make sure that the behavior exhibited by the agents in interaction with human participants can be accurately distinguished as theory of mind strategies by our RFX-BMS analysis.

As Figure 3 shows, the RFX-BMS estimation correctly classifies agent behavior of the  $ToM_1$ ,  $ToM_2$ , and  $ToM_3$  agents as consistent with the corresponding order of theory of mind reasoning. Moreover, the randomizing agent is classified as using a strategy that is equally consistent with all strategies. This shows that the RFX-BMS estimation can accurately distinguish different order of theory of mind reasoning, and also considers the randomizing agent to be unpredictable.

### 6.2 Participant behavior

No significant influence of the sequence on the participant scores was found ( $F = 0.483, p = 0.487$ ). That is, there is no reason to believe that the sequence in which the agents appeared affected the performance of the participants in any way. In the remainder, we therefore present results that are aggregated across the different sequences.

Figure 3 shows the estimated strategies of the participants in our experiment. Since we expect that participants adjust their strategy to their opponent, we analyze participant behavior for each of the four agent opponents separately.

The RFX-BMS analysis for the blocks in which participants played against a  $ToM_1$  agents are indicated by green bars in Figure 4. The figure shows that a

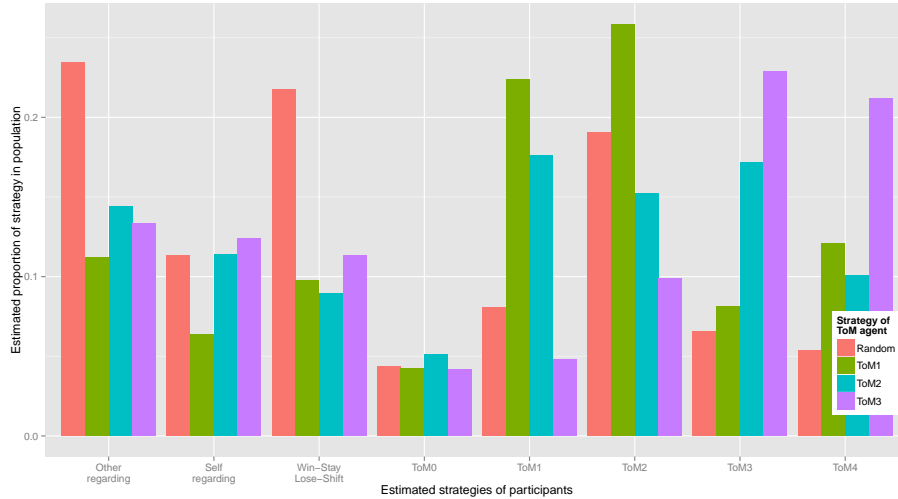


Fig. 4: Estimated strategy use of participants in the Mod game across the four different opponent types.

large proportion of participants are classified as making use of theory of mind. When playing against a  $ToM_1$  agent, an estimated 25.8% of the participants made use of a  $ToM_2$  strategy, while another estimated 22.4% of the participants used first-order theory of mind. It seems that, according to the participants, the  $ToM_2$  strategy was the best strategy. This also makes sense in relation to the theory, if you want to win a  $ToM$  game, it is most beneficial to think exactly one order of theory of mind higher than your opponent does. That is, when playing against a  $ToM_1$  opponent, it is best to make use of second-order theory of mind.

When playing against a  $ToM_2$  agent, participants mainly played according to first-order, second-order, or third-order theory of mind, as indicated by the blue bars in Figure 4. Note that it would be optimal for participants to reason using third-order theory of mind, while following first-order theory of mind would mean that the computer opponent consistently outsmarts the participant.

The purple bars of Figure 4 show the estimates strategies of participants while playing against a  $ToM_3$  agent. According to our RFX-BMS analysis, the strategies that best explained the largest percentages of the population are the  $ToM_3$  and  $ToM_4$  strategies, that represent an estimated 22.9% and 21.2% of the population respectively. Interestingly, the higher the order of theory of mind reasoning of the opponent, the longer participants took to react. In addition, those participants that were classified as using a higher order of theory of mind reasoning also obtained higher scores on average.

Note that for each of the  $ToM_1$ ,  $ToM_2$ , and  $ToM_3$  opponent agents, theory of mind strategies were estimated to explain the largest proportions of participant choices. However, this is not true when participants play against the randomizing

*ToM* agent, as shown by the red bars in Figure 4. While the behavior of a sizable proportion of participants is consistent with the *ToM*<sub>2</sub> strategy, participant behavior is generally better described as other-regarding or as a win-stay, lose-shift strategy. That is, unlike when participants play against the *ToM*<sub>0</sub>, *ToM*<sub>1</sub>, or *ToM*<sub>2</sub> agent, participant behavior against the random *ToM* opponent is better explained by behavior-based strategies than it is by theory of mind strategies.

## 7 General discussion and conclusion

In many social settings, people rely on the use of theory of mind to guide their actions. However, experimental evidence shows that adults more readily show signs of theory of mind reasoning in story comprehension tasks than in strategic games and that children learn to apply theory of mind in games much later than in story tasks [12,19,5,1,16].

In this paper, we let participants play a specific competitive game known as the Mod game [8,21] against artificial theory of mind agents, to determine to what extent participants use theory of mind reasoning. Our results show that, when playing against agents that use increasingly higher orders of theory of mind, a larger proportion of participants is estimated to make use of higher-order theory of mind. Participants that played against a first-order theory of mind agents relied mostly on first-order or second-order theory of mind themselves, while participants that played against a third-order theory of mind agent were better described as using third-order or fourth-order theory of mind. Moreover, when faced against a randomizing agent, who randomly selected to play as if it were a first-order, second-order, or third-order at the start of each round, participants were better described as relying on simpler behavior-based strategies.

In strategic games, participants are typically found to rely on low orders of theory of mind, and to be slow to adjust their level of theory of mind reasoning to more sophisticated opponents [10,2,25,9]. Earlier empirical research suggests that the use of first-order and second-order theory of mind in games can be facilitated by creating a believable story or insightful visual representation around an abstract problem [15,3], by creating a clear competition or negotiation setting [9,24], or by providing stepwise training from games that require zero-order *ToM* to second-order *ToM* games [14].

Impressively, our results in the Mod game suggest that participants even make use of an unprecedented *fourth-order* theory of mind reasoning when playing against a higher-order theory of mind opponent in the Mod game, even though they only faced each opponent for twenty consecutive rounds of play. Moreover, participant that were classified as using higher order of theory of mind tended to obtain higher scores. In future work, it would be interesting to determine to what extent participants exhibit the same behavior when facing more than one opponent (i.e., for  $n > 2$ ).

Interestingly, our results in the Mod game show higher levels of theory of mind reasoning than results of similar experiments with the matching pennies game [20,3]. One possible explanation is that adults in the matching pennies

experiment played 60 rounds in which they were given 1300 ms each to make a decision. In contrast, in our experiment, participants took an average 3700 ms to respond in the easiest blocks. That is, participants in the matching pennies experiment may not have had enough time to engage in higher-order theory of mind reasoning, and relied on simpler behavior-based strategies instead.

Alternatively, while the Mod game is arguably more difficult than matching pennies, it is easier to distinguish strategies from one another in the Mod game than in matching pennies. In matching pennies, players can only choose two possible actions. As a result, the choice of a particular action gives little information about the underlying strategy. In our Mod game experiment, however, there are 24 possible actions, which allows players to more easily interpret the actions of their opponent and draw conclusions about the underlying strategy.

The surprisingly high level of theory of mind reasoning of participants can also be partially explained by the representation of the game. In our specification of theory of mind agents, we assumed that the zero-order theory of mind agent reasons about opponent actions. However, because the actions are meaningfully arranged in a circle, it is also possible to define a zero-order theory that reasons about the change of actions (i.e. rate) of opponents rather than their actual choices. A first-order theory of mind agent in this former specification would have similar behavior to a zero-order theory of mind agent in the latter specification. That is, the representation of the zero-order theory of mind model is important in determining at which order of theory of mind participants are reasoning [15].

## Conclusion

Since the estimated level of theory of mind reasoning of participants increases with the increasing order of theory of mind reasoning of the artificial opponent, our results suggest that artificial agents can indeed encourage people to make use of higher order of theory of mind reasoning, thereby providing people with better results.

## References

1. Arslan, B., Verbrugge, R., Taatgen, N., Hollebrandse, B.: Teaching children to attribute second-order false beliefs: A training study with feedback. In: *CogSci*. pp. 108–113 (2015)
2. Camerer, C., Ho, T., Chong, J.: A cognitive hierarchy model of games. *Quarterly Journal of Economics* 119(3), 861–898 (2004)
3. Devaine, M., Hollard, G., Daunizeau, J.: The social Bayesian brain: Does mentalizing make a difference when we learn? *PLoS Computational Biology* 10(12), e1003992 (2014)
4. Devaine, M., Hollard, G., Daunizeau, J.: Theory of mind: Did evolution fool us? *PloS ONE* 9(2), e87619 (2014)
5. Flobbe, L., Verbrugge, R., Hendriks, P., Krämer, I.: Children’s application of theory of mind in reasoning and language. *Journal of Logic, Language and Information* 17(4), 417–442 (2008)

6. Franke, M., Galeazzi, P.: On the evolution of choice principles. In: Szymanik, J., Verbrugge, R. (eds.) *Proceedings of the Second Workshop Reasoning About Other Minds: Logical and Cognitive Perspectives*. CEUR Workshop Proceedings, vol. 1208, pp. 11–15 (2014)
7. Frey, S.: *Complex collective dynamics in human higher-level reasoning: A study over multiple methods*. Ph.D. thesis, Indiana University (2013)
8. Frey, S., Goldstone, R.L.: Cyclic game dynamics driven by iterated reasoning. *PLoS ONE* 8(2), e56416 (2013)
9. Goodie, A.S., Doshi, P., Young, D.L.: Levels of theory-of-mind reasoning in competitive games. *Journal of Behavioral Decision Making* 25(1), 95–108 (2012)
10. Hedden, T., Zhang, J.: What do you think I think you think?: Strategic reasoning in matrix games. *Cognition* 85(1), 1–36 (2002)
11. Imuta, K., Henry, J.D., Slaughter, V., Selcuk, B., Ruffman, T.: Theory of mind and prosocial behavior in childhood: A meta-analytic review. *Developmental Psychology* 52(8), 1192–1205 (2016)
12. Kinderman, P., Dunbar, R.I., Bentall, R.P.: Theory-of-mind deficits and causal attributions. *British Journal of Psychology* 89(2), 191–204 (1998)
13. Liddle, B., Nettle, D.: Higher-order theory of mind and social competence in school-age children. *Journal of Cultural and Evolutionary Psychology* 4(3-4), 231–244 (2006)
14. Meijering, B., van Rijn, H., Taatgen, N., Verbrugge, R.: I do know what you think I think: Second-order theory of mind in strategic games is not that difficult. In: *CogSci*. pp. 2486–2491 (2011)
15. Meijering, B., Van Maanen, L., Van Rijn, H., Verbrugge, R.: The facilitative effect of context on second-order social reasoning. In: *CogSci*. pp. 1423–1428 (2010)
16. Meijering, B., Taatgen, N.A., van Rijn, H., Verbrugge, R.: Modeling inference of mental states: As simple as possible, as complex as necessary. *Interaction Studies* 15(3), 455–477 (2014)
17. Premack, D., Woodruff, G.: Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences* 1(04), 515–526 (1978)
18. Stephan, K.E., Penny, W.D., Daunizeau, J., Moran, R.J., Friston, K.J.: Bayesian model selection for group studies. *Neuroimage* 46(4), 1004–1017 (2009)
19. Stiller, J., Dunbar, R.I.: Perspective-taking and memory capacity predict social network size. *Social Networks* 29(1), 93–104 (2007)
20. de Weerd, H., Diepgrond, D., Verbrugge, R.: Estimating the use of higher-order theory of mind using computational agents. *The B.E. Journal of Theoretical Economics* (2017)
21. de Weerd, H., Verbrugge, R., Verheij, B.: Theory of mind in the Mod game: An agent-based model of strategic reasoning. In: *Proceedings ECSI*. pp. 128–136 (2014)
22. de Weerd, H., Verheij, B.: The advantage of higher-order theory of mind in the game of limited bidding. In: *Proceedings Workshop ‘Reasoning about other Minds’, CEUR Workshop Proceedings*, vol. 751, pp. 149–164 (2011)
23. de Weerd, H., Verbrugge, R., Verheij, B.: How much does it help to know what she knows you know? An agent-based simulation study. *Artificial Intelligence* 199–200, 67–92 (2013)
24. de Weerd, H., Verbrugge, R., Verheij, B.: Negotiating with other minds: The role of recursive theory of mind in negotiation with incomplete information. *Autonomous Agents and Multi-Agent Systems*, 31, 250–287 (2017)
25. Wright, J.R., Leyton-Brown, K.: Beyond equilibrium: Predicting human behavior in normal-form games. In: *Proceedings of the Twenty-Fourth Conference on Artificial Intelligence*. pp. 901–907 (2010)