

UC Merced

Proceedings of the Annual Meeting of the Cognitive Science Society

Title

Savvy software agents can encourage the use of second-order theory of mind by negotiators

Permalink

<https://escholarship.org/uc/item/2nx2w4ks>

Journal

Proceedings of the Annual Meeting of the Cognitive Science Society, 37(0)

Authors

de Weerd, Harmen

Broers, Eveline

Verbrugge, Rineke

Publication Date

2015

Peer reviewed

Savvy software agents can encourage the use of second-order theory of mind by negotiators

Harmen de Weerd, Eveline Broers, Rineke Verbrugge
Institute of Artificial Intelligence, University of Groningen

Abstract

In social settings, people often reason about unobservable mental content of other people, such as their beliefs, goals, or intentions. This ability helps them to understand and predict the behavior of others. People can even take this ability further, and use higher-order theory of mind to reason about the way others use theory of mind, for example in 'Alice believes that Bob does not know about the surprise'. However, empirical evidence suggests that people do not spontaneously use higher-order theory of mind in strategic games. In this paper, we let participants negotiate with computational theory of mind agents in the setting of Colored Trails. We find that even though participants are unaware of the level of sophistication of their trading partner, within a few rounds of play, participants offers are more indicative of second-order theory of mind reasoning when their trading partner was using second-order theory of mind as well.

Keywords: theory of mind; social cognition; negotiation; strategic games

Introduction

In social settings, people reason about unobservable mental content, such as beliefs, desires, and goals, to predict and interpret the behavior of others. This *theory of mind* (Premack & Woodruff, 1978) allows people to reason explicitly about the goals of others, such as deciding whether the behavior of others is accidental or intentional. Empirical evidence from second-order false belief tasks (Perner & Wimmer, 1985; Miller, 2009) reveals that people are also capable of reasoning about the theory of mind of others. People use *second-order theory of mind* when they reason about the beliefs others have about the beliefs of yet other people, and realize that such nested beliefs can be incorrect. Second-order theory of mind allows people to form nested beliefs such as "Alice believes that Bob does not know about the surprise party", and use these beliefs to interpret and predict Alice's behavior.

While participants readily use second-order theory of mind reasoning in the second-order false belief task, empirical evidence suggests that in strategic games, participants do not appear to make spontaneous use of higher-order (i.e. at least second-order) theory of mind (Hedden & Zhang, 2002; Camerer et al., 2004; Wright & Leyton-Brown, 2010; Goodie et al., 2012). Over a range of unrepeated single-shot games, Camerer et al. (2004) estimate the distribution of the level of sophistication used by human participants. They find that participant reasoning is typically limited to the use of zero-order or first-order theory of mind. Only few participants are found to be well-described as higher-order theory of mind reasoners (Wright & Leyton-Brown, 2010). When games are repeated, participants can successfully adjust their level of reasoning to accurately predict the behavior of other theory of mind reasoners (Hedden & Zhang, 2002; Zhang et al., 2012; Goodie

et al., 2012; Meijering et al., 2010, 2011, 2014; De Weerd et al., 2014; Devaine et al., 2014), although participants typically need many trials before their behavior is consistent with higher-order reasoning.

In this paper, we investigate human-agent interactions in the influential Colored Trails setting, introduced by Grosz, Kraus, and colleagues (Lin et al., 2008; Gal et al., 2010)¹, which provides a useful test-bed to study how different aspects of mixed-motive settings change interactions among agents and humans. In previous work in this negotiation setting, we presented a computational model for theory of mind agents to study the effectiveness of higher-order theory of mind reasoning (De Weerd et al., 2013). We found that the use of first-order and second-order theory of mind allows software agents to balance competitive and cooperative aspects of the game. This way, the use of theory of mind prevents negotiations from breaking down the way they do for agents without theory of mind. In the current paper, we use these agents to determine to what extent human participants reason at higher orders of theory of mind, by letting software agents interact directly with human participants.

Colored Trails

The game we study in this paper is a variation on the influential Colored Trails game. Colored Trails is a board game designed as a research test-bed for investigating decision-making of people and software agents (Lin et al., 2008; Gal et al., 2010). We consider a specific setting in which two negotiating agents alternate in making offers. We have previously used this setting to test the effectiveness of higher-order theory of mind in negotiations (De Weerd et al., 2013).

The game is played on a square board consisting of 25 patterned tiles, like the one depicted in Figure 1a. At the start of the game, each player receives a set of four patterned chips, selected at random from the same four possible patterns as those on the board. Each player is initially located on the center tile of the board, indicated with the letter *S* in Figure 1b. The goal of each player is to reach their personal goal location, which is drawn randomly from the board tiles that are at least three steps away from the initial location (gray tiles in Figure 1b). To move around on the board, players use their chips. A player can move to a tile adjacent to his current location by handing in a chip of the same pattern as the destination tile. Figure 1a shows an example of a Colored Trails board as well as a possible path across the board. A player following the path from location *A* to the blank tile marked *B*

¹Also see <http://coloredtrails.atlassian.net/wiki/display/coloredtrailshome/>.

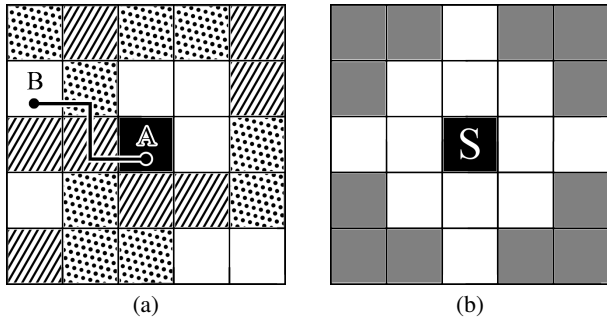


Figure 1: In Colored Trails, players spend chips to move around on a 5 by 5 board. (a) To follow the path from location A to location B, a player needs to hand in one blank, one striped, and one dotted chip. (b) Each player is initially located on the central tile S and is assigned a goal location drawn randomly from the gray tiles.

would have to hand in one striped chip, one dotted chip, and one blank chip.

Players are scored based on their success in reaching their goal location. Each player receives 50 points at the start of the game. If the player successfully reaches his goal, he receives an additional 50 points. However, if the goal is not reached, 10 points are deducted from the player’s score for each step needed to reach the goal location. Finally, any chip that has not been used to move around the board is worth an additional 5 points to its owner. For example, consider the situation in Figure 2, and suppose that player i has goal location G . With his initial set of chips, player i can obtain a score of 50 points. However, if player i would receive one of agent j ’s blank chips, he could obtain a score of 110 points.

To get closer to their goal location, players can trade chips with their co-player. To capture the dynamic aspect of negotiation, trading among players takes the form of a sequence of alternating offers. When a player makes an offer to redistribute the chips a certain way, his trading partner decides whether or not to accept this offer. If the offer is accepted, the proposed distribution of chips becomes final, the players move as close to their respective goal locations as possible, and the game ends. Alternatively, the trading partner may decide to withdraw from negotiations, which makes the initial distribution final. As a third option, the trading partner may decide to continue the game by rejecting the current offer and making his own offer for a redistribution of chips.

There are no restrictions on the offers that players can make. For example, a player is allowed to repeat an offer that has been previously rejected by his trading partner, or make an offer that he has previously rejected himself. However, the game ends if six offers have been rejected. In any game, each player can therefore make at most three offers. If a game ends because the maximum number of offers has been exceeded, the initial distribution of chips becomes final.

Although a player’s score is based only on how closely he approaches his own goal, Colored Trails is not a purely com-

petitive game. Since a player may need a different set of chips to achieve his goal than his trading partner, there may be an opportunity for a cooperative trade that allows both players to obtain a higher score. That is, although the score of a player does not depend on how closely his trading partner approaches his goal location, players can still benefit from taking into account the goal of their trading partner. Importantly, however, Colored Trails is a game of imperfect information: while players know what chips are in possession of their trading partner, they do not know the goal location of their trading partner at the start of the game.

In this paper, we investigate to what extent human participants reason using higher orders of theory of mind when playing Colored Trails with a software agent as trading partner, and to what extent participants adjust their level of theory of mind reasoning in response to the behavior of their trading partner. Since simulation experiments with agents have shown that second-order theory of mind can help agents to avoid negotiation failure and balance cooperative and competitive aspects of the game (De Weerd et al., 2013), the Colored Trails setting may facilitate theory of mind reasoning in human participants as well.

Theory of mind software agents

The theory of mind agents presented here as trading partners of human participants are adapted from De Weerd et al. (2013) to allow for games with a known finite horizon. That is, the computational agents know that the game cannot last more than six turns. In this section, we describe the way these make use of theory of mind. The mathematical details of these agents can be found in De Weerd et al. (2013).

Zero-order theory of mind

A zero-order theory of mind (ToM_0) agent is unable to reason about unobservable mental content of its trading partner, including its goal location. Instead, a ToM_0 agent models the behavior of its trading partner in terms of the offers that the trading partner is willing to accept. Based on previous experience in the Colored Trails game, a ToM_0 agent constructs beliefs about the likelihood that certain offers will be accepted by the trading partner. For example, over repeated games, a ToM_0 agent will learn that the trading partner rarely accepts offers that assign few chips to the trading partner, while offers that assign many chips to the trading partner are accepted with a high frequency.

Using these zero-order beliefs, a ToM_0 agent can calculate the expected gain of making a particular offer, and choose the action that the agent expects to yield it the highest gain. Based on the actions of the trading partner, the ToM_0 agent then updates its zero-order beliefs. This way, the ToM_0 agent can play Colored Trails without attributing any mental content to others. In particular, although the ToM_0 agent’s zero-order beliefs eventually reflect the desires of its trading partner, the ToM_0 agents cannot reason about such desires explicitly.

In terms of negotiation strategies, the ToM_0 agent engages purely in *positional bargaining* (Fisher & Ury, 1981), by rea-

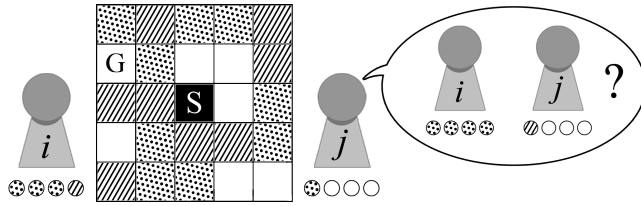


Figure 2: Example of a negotiation setting in Colored Trails. Agent j offers to trade the striped chip owned by agent i against the dotted chip owned by agent j . Since this trade would make it harder for agent i to reach his goal location (tile G), agent i will reject this offer. The goal location of agent j is not shown.

soning only about offers and the likelihood that these offers will be accepted by its trading partner.

First-order theory of mind

In addition to its zero-order beliefs, a first-order theory of mind ToM_1 agent can also determine what its own decision would have been if it had been in the position of its trading partner. This way, a ToM_1 agent can consider that its trading partner has beliefs and goals similar to its own that determine whether or not an offer will be accepted.

A ToM_1 agent believes that an offer will only be accepted if it increases the score of both the agent itself and its trading partner since the ToM_1 agent itself would only accept offers that increase its own score. In the same way, the ToM_1 agent realizes that its trading partner only makes offers that would increase its own score. This means that the offers made by the trading partner contain information about its goal location. For example, consider the situation depicted in Figure 2. In this example, agent j offers to trade its dotted chip for the striped chip owned by agent i . From this offer, a ToM_1 agent would conclude that the striped chip allows agent j to move closer to its goal location. Secondly, since the offered trade would leave agent j without any dotted chips, a ToM_1 agent would believe that agent j does not need any dotted chips to reach its goal location. This excludes several possible goal locations for agent j .

Importantly, a ToM_1 agent's first-order theory of mind is additional to its zero-order beliefs. Through repeated interactions, a ToM_1 agent may come to believe that first-order theory of mind fails to accurately model the behavior of its trading partner and that the use of zero-order theory of mind would result in a higher score. In this case, a ToM_1 agent may decide to play Colored Trails as if he were a ToM_0 agent.

Second-order theory of mind

Agents capable of second-order theory of mind can also consider the possibility that their trading partner is a ToM_1 agent. A second-order theory of mind (ToM_2) agent believes that its trading partner may be trying to interpret the offers made by the ToM_2 agent to determine the ToM_2 agent's goal. This allows the ToM_2 agent to reason about the way different offers

influence the beliefs of the trading partner about the agent's goal, and select the offer that provides its trading partner with as much information about its goal location as possible.

For example, suppose agent i in Figure 2 is a ToM_2 agent with goal location G . Using second-order theory of mind, agent i knows that making any offer in which the striped chip is assigned to agent j , agent j can conclude that agent i does not need a striped chip to reach its goal location. This allows agent j to exclude many possible goal locations for agent i , which can help agent j to make an offer that is acceptable to agent i . In this case, although agent i knows that agent j has a goal location, agent i remains unaware of what that goal location is. A second-order theory of mind agent can therefore engage in *interest-based negotiation* (Fisher & Ury, 1981), by choosing its offers in such a way that they communicate the agent's interests to its trading partner.

Similar to the ToM_1 agent, a ToM_2 agent does not know the extent of its trading partner's theory of mind abilities. Instead, a ToM_2 agent has zero-order, first-order, and second-order beliefs about the behavior of its trading partner. While negotiating in Colored Trails, the ToM_2 agent keeps updating its beliefs concerning which of these beliefs most accurately describes the actual behavior of its trading partner. This means that a ToM_2 agent may sometimes behave as if it were a ToM_0 agent, while behaving like a ToM_2 agent on other occasions.

Methods

Participants

Twenty-seven students (10 female) of the University of Groningen participated in this study. All participants were informed that after the conclusion of the study, the three participants with the highest score in the negotiation game received €15, €10, and €5, respectively. Each participant gave informed consent prior to admission into the study.

Materials

Twenty-four games were selected from a set of randomly generated games. To ensure that these games would allow us to distinguish between different orders of theory of mind reasoning of participants, they were selected so that:

- The participant's goal could be reached with the eight chips in the game;
- Simulations with computational agents predicted different outcomes for participants using zero-order, first-order, and second-order theory of mind; and
- Simulations with computational agents predicted that the game would last at least two turns and at most six turns.

These games were divided into three blocks of eight games each. The level of theory of mind reasoning of the software agent was fixed within each block, and varied between blocks.

Design and procedure

Before the start of the experiment, participants were tested on colorblindness. Participants were asked to distinguish patches of blue and orange, with four possible intensities of each color. All participants passed the colorblindness test. Next, participants played several Marble Drop games (Meijering et al., 2011).

The Colored Trails experiment consisted of a familiarization phase and an experimental phase. At the start of the familiarization phase, participants were asked to imagine themselves as an attorney for a major corporation. In this function, they would be involved in a number of negotiations with different clients. Participants were told that their trading partner was a computer player (Alex), which would always react on their offer as quickly as possible in a way it believed would maximize its own score. To ensure understanding of the Colored Trails game, participants answered a few questions about the rules, scoring, and movements on the game board.

In the experimental phase, participants played three blocks of eight games each. In each block, the participant either faced a ToM_0 , ToM_1 , or ToM_2 agent. The order of the blocks was randomized across participants. Participants were not informed that the level of reasoning of the trading partner would change over the course of the experiment, but participants were told that they would face different clients. At the start of the experiment, it was randomly decided whether the participant or the software agent would make the initial offer of the first game. In subsequent games, participant and agent alternated in the role of initiating player.

Participants were allowed 60 seconds to decide on their next action. During each round, the remaining decision time was presented to participants by means of a countdown timer. If a participant had not made a decision within 60 seconds, the game continued without an offer being made, and the software agent took its turn.

The zero-order beliefs of theory of mind agents were initialized by playing 200 randomly generated Colored Trails games against another agent. This allowed agents to learn what kind of offers were more likely to be acceptable to their trading partner. To conform to our cover story in which participants were told that they would face a number of different clients, the agent's beliefs were reset to this initial value at the start of each game. Additionally, theory of mind agents started every game reasoning at the highest order of theory of mind available to them. This means that although software agents learned from a participant's offers within a single game, and adjusted their behavior accordingly, agents did not exhibit any learning across games. This way, agents were prevented from adapting to specific participants, and every participant faced the same agent in every scenario.

After the Colored Trails games, participants answered a short questionnaire about the perceived difficulty of the task, the behavior of their trading partner, and the participant's reasoning strategies. In addition, participants took a test for their interpersonal reactivity index (Davis, 1983).

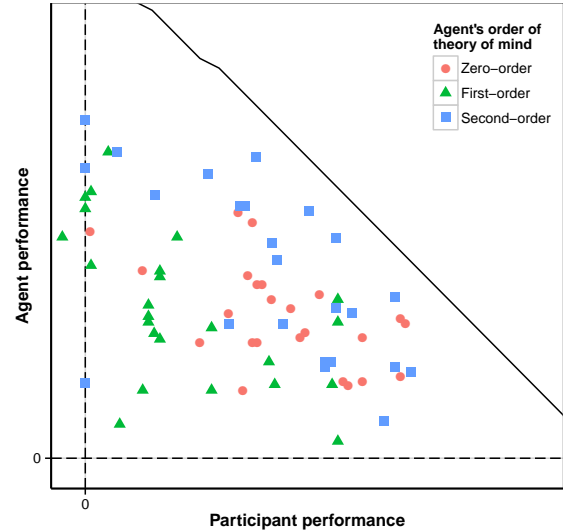


Figure 3: Outcomes of Colored Trails per block. The solid line indicates Pareto optimal outcomes. Dashed lines show the score participants and software agents would achieve if they were to withdraw from negotiation in every game.

Results

Figure 3 shows the outcomes of the Colored Trails game. The graph shows how the score of agents and participants changed as a result of negotiation for each participant and for each block. Dashed lines indicate the zero performance line, which is the score that players would have received if every game of the block had ended with withdrawal from negotiation. A score below the dashed line indicates that a player decreased its score through negotiation. As Figure 3 shows, only once a participant received a negative score in one of the blocks.

The solid line in Figure 3 shows the boundary of Pareto efficient outcomes. This boundary shows those outcomes for which neither the participant nor the software agent could have received a higher score without a decrease in the score of the other player. The Pareto boundary gives an impression of how well participants and software agents played Colored Trails. Participants and software agents generally negotiated mutually beneficial solutions, while neither player systematically exploited the other. Additionally, Figure 3 shows that when participants negotiate with ToM_2 agents, they tend to end up closer to the Pareto optimality, while negotiations between participants and ToM_1 agents typically end up far from the Pareto boundary.

Importantly, the use of theory of mind agents allows us to estimate to what extent participants make use of theory of mind while playing Colored Trails. We use a ToM_3 'spectator' agent that observes the offers of a participant and determines whether these offers are most consistent with zero-order, first-order, or second-order theory of mind reasoning. The software agent constructs a *confidence* for each order of theory of mind at which it can reason to decide which order

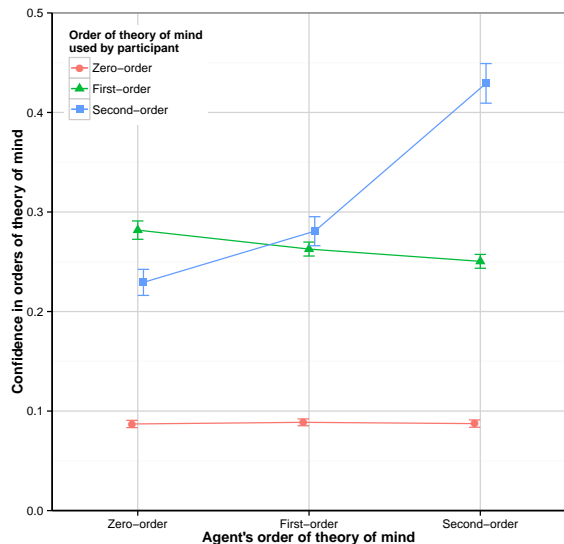


Figure 4: Estimated similarity of participant offers to the offers of ToM_0 (red circles), ToM_1 (green triangles), and ToM_2 (blue squares) agents in each of the three blocks. Brackets indicate one standard error.

of theory of mind would yield the best outcome (De Weerd et al., 2013). Each time the participant makes an offer O , the ToM_3 agent updates its confidence that this participant is using k th-order theory of mind by calculating the likelihood that a ToM_k agent would have made an offer similar to offer O .

For each of the three blocks, Figure 4 shows how similar participant offers were to offers of ToM_0 , ToM_1 , and ToM_2 agents, as judged by the ToM_3 agent. Red circles indicate the average similarity of a participant’s offers to zero-order theory of mind reasoning, green triangles indicate the similarity to first-order theory of mind reasoning, and blue squares show the similarity to second-order theory of mind reasoning. Interestingly, Figure 4 shows that participant offers are more similar to first-order and second-order theory of mind reasoning than they are to zero-order theory of mind reasoning.

Figure 4 also shows that the similarity ratings of participant offers vary depending on the order of theory of mind of the computer trading partner. Although similarity ratings for zero-order and first-order theory of mind reasoning show no variation across different levels of sophistication of the trading partner ($X^2_{(2)} = 0.52$, ns, and $X^2_{(2)} = 2.67$, ns, respectively), participant offers were significantly more similar to second-order theory of mind reasoning when they were facing a ToM_2 trading partner ($X^2_{(2)} = 24.89$, $p < 0.001$).

Previous studies into negotiations show that the opening bid of a negotiation can serve as an anchor for the entire negotiation process, making the first bid of a game especially influential in the negotiation process (Raiffa et al., 2002; Van Poucke & Buelens, 2002). In our experiment, the identity of the initiating player indeed influences negotiation outcomes. In general, both players ended up with an extra 15

points on average after negotiation when the software agent made the initial offer rather than when the participant was the first to propose a trade. The only exception to this rule was that participants negotiating with a ToM_2 agent ended up with a higher score when they made the initial offer themselves. This effect can be explained by the way agents of different orders of theory of mind construct their offers. Both ToM_0 and ToM_1 agents make offers that they believe will be accepted by their trading partner. In contrast, ToM_2 agents make offers that inform their trading partner about their own goals. As a result, initial offers made by ToM_0 and ToM_1 agents are typically more favorable to their trading partner than those made by ToM_2 agents. Similarly, when participants reasoned more like ToM_2 agents, their initial offers were more favorable to themselves than to their trading partner.

Discussion and conclusion

Experimental evidence suggests that participants do not make spontaneous use of higher-order theory of mind reasoning in unrepeatable games (Hedden & Zhang, 2002; Camerer et al., 2004; Wright & Leyton-Brown, 2010; Goodie et al., 2012), though participants can successfully adjust their level of reasoning to accurately predict the behavior of other theory of mind reasoners (Hedden & Zhang, 2002; Goodie et al., 2012; Meijering et al., 2010, 2011, 2014; Devaine et al., 2014). Surprisingly, in this paper, we find that participants show behavior consistent with second-order theory of mind reasoning in a negotiation game that lasts a few rounds only.

In our experiments, human participants negotiated with software agents that dynamically change their order of theory of mind reasoning in response to the behavior of their trading partner. We use these agent-based models to analyze participant behavior in a dynamic setting. Our model explicitly takes into account that participants may differ in the order of theory of mind at which they reason, and that a participant may change the order of theory of mind at which he reasons over the course of a single game. Based on this agent-based analysis, we find that participants make offers that are more consistent with second-order theory of mind reasoning when their trading partner is capable of second-order theory of mind as well. Interestingly, while participants knew that they would face different trading partners, they were unaware that these trading partners differed in their theory of mind abilities. That is, the behavior of higher-order theory of mind agents apparently encouraged participants to make use of higher-order theory of mind as well.

Experiments with adults typically show that individuals reason at low orders of theory of mind, and are slow to adjust to an opponent that reasons using theory of mind (Hedden & Zhang, 2002; Camerer et al., 2004; Wright & Leyton-Brown, 2010; Goodie et al., 2012). In our setting, however, participants exhibited second-order theory of mind within a few games. It is possible that the negotiation setting, which involves both cooperative and competitive goals, emphasized the social nature of the task. Such social framing has been

shown to encourage the use of theory of mind (Goodie et al., 2012; Devaine et al., 2014).

Our results show that mixed groups of human and software agents can successfully negotiate a mutually beneficial outcome. However, none of the negotiation outcomes in our experiment were Pareto efficient. That is, each participant could have received a higher score without reducing the score of their trading partner. This indicates that there is still room for significant improvement. One factor that may have limited the effectiveness of negotiations in Colored Trails is the time limit on the decisions of participants. Although participants failed to make a decision before time ran out in only four occasions, it is likely that participants selected suboptimal actions due to time constraints. Removing this time constraint in future experiments may increase negotiation performance.

Our results suggest that computational theory of mind agents can be used as a training tool for negotiation. When participants negotiated with a trading partner capable of second-order theory of mind, the outcome was generally closer to a Pareto optimal solution than when participants faced less sophisticated trading partners. In addition, our results show that agents could benefit more from making the opening bid than participants. Experience with theory of mind agents may therefore allow participants to learn to leverage the anchoring effect of the initial offer.

Acknowledgments

This work was supported by the Netherlands Organisation for Scientific Research (NWO) Vici grant NWO 277-80-001, awarded to Rineke Verbrugge for the project ‘Cognitive systems in interaction: Logical and computational models of higher-order social cognition’.

References

- Camerer, C., Ho, T., & Chong, J. (2004). A cognitive hierarchy model of games. *Quarterly Journal of Economics*, *119*(3), 861–898.
- Davis, M. H. (1983). Measuring individual differences in empathy: Evidence for a multidimensional approach. *Journal of Personality and Social Psychology*, *44*(1), 113–126.
- Devaine, M., Hollard, G., & Daunizeau, J. (2014). The social Bayesian brain: Does mentalizing make a difference when we learn? *PLoS Computational Biology*, *10*(12), e1003992.
- de Weerd, H., Verbrugge, R., & Verheij, B. (2013). Higher-order theory of mind in negotiations under incomplete information. In G. Boella, E. Elkind, B. T. R. Savarimuthu, F. Dignum, & M. K. Purvis (Eds.), *Prima 2013: Principles and practice of multi-agent systems, Dunedin, New Zealand, december 2013* (pp. 101–116).
- de Weerd, H., Verbrugge, R., & Verheij, B. (2014). Theory of mind in the Mod game: An agent-based model of strategic reasoning. In A. Herzig & E. Lorini (Eds.), *Proceedings of the european conference on social intelligence (ECSI-2014)* (pp. 129–136).
- Fisher, R., & Ury, W. (1981). *Getting to yes: Negotiating agreement without giving in*. Penguin Books.
- Gal, Y., Grosz, B., Kraus, S., Pfeffer, A., & Shieber, S. (2010). Agent decision-making in open mixed networks. *Artificial Intelligence*, *174*(18), 1460–1480.
- Goodie, A. S., Doshi, P., & Young, D. L. (2012). Levels of theory-of-mind reasoning in competitive games. *Journal of Behavioral Decision Making*, *25*(1), 95–108.
- Hedden, T., & Zhang, J. (2002). What do you think I think you think?: Strategic reasoning in matrix games. *Cognition*, *85*(1), 1–36.
- Lin, R., Kraus, S., Wilkenfeld, J., & Barry, J. (2008). Negotiating with bounded rational agents in environments with incomplete information using an automated agent. *Artificial Intelligence*, *172*(6), 823–851.
- Meijering, B., Taatgen, N. A., van Rijn, H., & Verbrugge, R. (2014). Modeling inference of mental states: As simple as possible, as complex as necessary. *Interaction Studies*, *15*, 455–477.
- Meijering, B., van Maanen, L., van Rijn, H., & Verbrugge, R. (2010). The facilitative effect of context on second-order social reasoning. In S. Ohlsson & R. Catrambone (Eds.), *Proceedings of the 32nd annual conference of the cognitive science society* (pp. 1423–1428).
- Meijering, B., van Rijn, H., Taatgen, N. A., & Verbrugge, R. (2011). I do know what you think I think: Second-order theory of mind in strategic games is not that difficult. In L. Carlson, C. Hoelscher, & T. F. Shipley (Eds.), *Proceedings of the 33rd annual conference of the cognitive science society* (pp. 2486–2491).
- Miller, S. A. (2009). Children’s understanding of second-order mental states. *Psychological Bulletin*, *135*(5), 749–773. doi: 10.1037/a0016854
- Perner, J., & Wimmer, H. (1985). “John thinks that Mary thinks that...”. Attribution of second-order beliefs by 5 to 10 year old children. *Journal of Experimental Child Psychology*, *39*(3), 437–71.
- Premack, D., & Woodruff, G. (1978). Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, *1*(4), 515–526.
- Raiffa, H., Richardson, J., & Metcalfe, D. (2002). *Negotiation analysis: The science and art of collaborative decision making*. Belknap Press.
- Van Poucke, D., & Buelens, M. (2002). Predicting the outcome of a two-party price negotiation: Contribution of reservation price, aspiration price and opening offer. *Journal of Economic Psychology*, *23*(1), 67–76.
- Wright, J. R., & Leyton-Brown, K. (2010). Beyond equilibrium: Predicting human behavior in normal-form games. In *Proceedings of the twenty-fourth conference on artificial intelligence* (pp. 901–907).
- Zhang, J., Hedden, T., & Chia, A. (2012). Perspective-taking and depth of theory-of-mind reasoning in sequential-move games. *Cognitive Science*, *36*(3), 560–573.