Using Our Theory of Mind for Inferences in Strategic Reasoning

Alisa R. Zoltowski

University of Michigan

A Thesis Submitted in Partial Fulfillment of the Requirements for the Degree of Bachelor of

Science with Honors in Biopsychology, Cognition, & Neuroscience from the University of

Michigan 2015

Advisor: Dr. Jun Zhang

Abstract

Prior investigations into theory of mind have used strategic games to examine how an opponent's known preferences are used to determine behavior. However, in the social setting, preferences are not always known. Mental state information often comes from observing another person's decision and inferring what motivated that decision. This study extended the strategic game framework to explore the levels of reasoning used when making these types of inferences. The participants were 47 undergraduates at the University of Michigan. Each participant received game cards with decisions coming from a "prior player" and was asked if particular choices for that player's missing payoff information could have led to those decisions. The responses to these questions indicated whether participants attributed a strategy to the prior player that anticipated the next game move (a "predictive" strategy) or did not (a "myopic" strategy). The initially assumed strategy was found to be neutral, with predictive inferences becoming more common in successive game sets. These results support the role of experience in engaging greater depths of reasoning when interpreting another person's decisions, even in the absence of feedback to influence this shift. Additionally, the task of being instructed to interpret another person's decisions may engage in-depth theory of mind reasoning more readily than is naturally assumed during an occurring interaction. Further experiments are necessary to determine if and how the findings from abstract games correspond to real-world theory of mind use.

*Keywords*: theory of mind, strategic games, inferences

**Using Our Theory of Mind for Inferences in Strategic Reasoning**

Why do we care what other people think? Humans have developed the remarkable capacity to have a "theory of mind," which is the ability to attribute distinct mental states such as thoughts, desires, and intentions to another person. This ability is important for social interactions because having an accurate representation of the mental states guiding another person helps determine what behavior would be appropriate. There are plenty of everyday situations that our theory of mind helps us navigate. One example familiar to those in academia would be the decision whether or not to apply to a particular graduate school. To make an informed decision, the individual would need to estimate how favorably the admissions committee would view his/her application in order to determine if applying to that school is worth the effort.

Early signs of the ability to take another person's perspective seem to emerge during the first year of life (Onishi & Baillargeon, 2005). In infancy, perspective taking is demonstrated by implicit expectations. This ability later develops into making explicit verbal predictions about another person's mental states (Perner, Leekam, & Wimmer, 1987). The classic test for having a theory of mind assesses whether an individual can explicitly attribute a belief to another person that the individual taking the test exclusively knows to be false. Children begin to reliably demonstrate this ability at age 4 (Perner et al., 1987). The most advanced theory of mind milestone that has been studied is the ability to use higher levels of mental state reasoning, e.g., reasoning about what one person believes about a second person's mental states. This capacity seems to emerge between ages 6 and 7 (Perner & Wimmer, 1985). Studies of how these skills are acquired provide evidence that they are gradually developed through practice, as opposed to sudden insights of understanding (Amsterlaw & Wellman, 2006). A more recent focus of

research has been to study how these skills mature beyond childhood and how they are typically

used in adult interactions.

Although some types of basic perspective taking appear to be automatic (Kovacs, Teglas,

& Endress, 2010; Samson, Apperly, Braithwaite, Andrews, & Scott, 2010), more advanced

theory of mind usage more closely resembles a reasoning process that can be spontaneously

activated in relevant situations (Back & Apperly, 2010; Lin, Keysar, & Epley, 2010). As with

other human reasoning, there are biases and variability in how we reason about others' mental

states (Birch & Bloom, 2007; Converse, Lin, Keysar, & Epley, 2008; Mitchell, Robinson, Isaacs,

& Nye, 1996; Nickerson, 1999). Nickerson's (1999) cognitive model of how we form and update

models of other people's minds neatly accounts for some of the biases that are encountered. His

account sets our own mind as an initial model, with alterations that can be made in one of three

ways. Knowledge can be subtracted based on information we have reason to believe is held by us

alone, added based on group associations about the other person, or updated based on

interactions with the other person. Biases can result from a mistake at any of these steps. The

most general type of bias is a failure to make the appropriate distinctions between someone

else's mind and our own, which is known as an egocentric bias.

Studies of these biases have provided valuable evidence about how theory of mind

matures into adulthood. Epley, Morewedge, & Keysar (2004) looked at egocentric biases in

children and adults, with results indicating that theory of mind in adults is a more practiced and

efficient version of the skills seen in children. The investigations about a more specific type of

egocentric bias called a "knowledge bias" give further information about how theory of mind

develops. Complementary with Nickerson's (1999) account, this bias occurs when someone fails

to subtract his/her own unique knowledge about the situation when forming mental state

assumptions. Mitchell et al. (1996) surprisingly found some knowledge biases that are more prevalent in adults than children, because the integration of different types of mental state information changes as we age. Children seem to primarily weight information based on what another person has seen, whereas adults' consideration of visual versus verbal information to the other person can be influenced by the observer's own knowledge. Although this study demonstrated limitations due to the complexity of situations, being able to integrate multiple sources of mental state information is clearly an advantageous development.

Although integrating mental state information is already challenging, the next step is using this information in interactions. One type of task that has been used to assess how theory of mind is demonstrated behaviorally is a competitive two-player game with interactions between each player's strategies (Goodie, Doshi, & Young, 2010; Hedden & Zhang, 2002; Meijering, van Rijn, Taatgen, & Verbrugge, 2012). Game theory analysis can then be used to analyze the implications of players' choices. These research projects have typically looked at the process of participants choosing their own strategy, which if done optimally, necessarily involves predicting the move of an opponent with well-defined motives. Making predictions requires theory of mind; one must take the opponent's perspective in order to predict his/her strategy. These studies have shown clear distinctions between optimal strategies and how participants have behaved (Hedden & Zhang, 2002; Meijering et al., 2012). The optimal strategy in these games is to use backwards induction, by which a player selects the ideal final outcome and then compares prior decision points to determine which choices lead to that outcome. However, there is evidence that players do not necessarily use this strategy. Meijering et al. (2012) have tracked participants' eye movements and found that players seem to use the sub-optimal strategy of initially looking forward through the game scenario, similar to causal reasoning, then checking decisions by using

backwards comparisons. Although these studies are limited in their ability to extend to the social environment, the use of causal reasoning instead of the optimal game theory solution suggests that the approaches taken to solving these games may be a general mechanism used in other scenarios. Thus, it remains promising that the findings may extend to how we approach daily social interactions.

Several studies have used what is known as a "Stackelberg game" in the game theory literature, which is defined by sequential moves and a finite number of stages (Osborne & Rubinstein, 1994). Specifically, these studies have used two-player games with at most three stages (Goodie et al., 2010; Hedden & Zhang, 2002; Meijering et al., 2012). The decision at each of the three stages is to either end the game and receive the current outcome or continue the game to the next stage. The participant has control over the first and third stages, but an opponent controls the second stage. This format allows for theory of mind analysis because the decision whether or not to continue the game to the next stage depends on what the player expects the opponent to decide at that stage. The limitation to three stages allows for distinct strategies reflecting different types of reasoning.

A focus of these analyses has been the levels of theory of mind reasoning that can be used. As discussed before, higher levels of reasoning typically become available during childhood (Perner & Wimmer, 1985). In these specific games, the participant, Player I, can consider that the opponent, Player II, may be anticipating the participant's strategy at the third decision point. As in Hedden and Zhang (2002), a Player II who does indeed anticipate Player I can be termed "predictive." The alternative, a Player II who only compares his/her own outcomes without predicting Player I, is known as "myopic." More complicated options therefore exist in Player I's perspective of the game. A Player I who assumes a predictive Player

II is using second-order reasoning. A Player I who assumes a myopic Player II is using first-order reasoning. Although not found in these studies, Player I could also be myopic and not consider Player II at all. Hedden and Zhang (2002) found that first-order theory of mind reasoning, the less complex option, is most prevalent initially. However, second-order reasoning becomes increasingly more common throughout experience with a Player II who does indeed act predictively. This experiment was conducted using mixed-motive games, in which the players' payoff structures vary independently of each other. Another option is to study strictly competitive games, in which one player's loss corresponds to the other player's gain. Goodie et al. (2010) found that when using strictly competitive games rather than mixed-motive games, second-order reasoning, which is more complex, becomes the default. These mixed results indicate that these games contain a wealth of information to be processed and the resulting reasoning may reflect the difficulty of integrating information in each type of game.

Zhang et al. (2012) investigated this idea by attempting to dissociate information processing constraints from theory of mind reasoning limitations in the game. They found that even when participants were given the same amount of information, assigning them a different perspective influenced the likelihood of adopting the predictive game strategy. This finding provides evidence that advanced theory of mind reasoning requires cognitive resources beyond those solely used to process information. However, it does not rule out the possibility that decreasing executive demands may open up resources to be used for more complex perspective taking. Different types of tasks with different demands may provide tradeoffs in their conduciveness to theory of mind depth. Flobbe, Verbrugge, Hendriks, & Krämer (2008) found support that, in children, theory of mind development varies between linguistic and strategic tasks. Therefore, a different task such as explicitly asking individuals to reason about an

opponent's decision may lead to different levels of perspective-taking and motivates the present investigation.

The aim of this project is to use the strategic game framework to delve further into mental state reasoning. Previous literature has analyzed how known preferences were used to determine behavior. Less has been studied about the reverse process: inferring preferences from observed behavior. Much of our information about other people comes from observing their actions and decisions; less often is an interaction completely specified or preferences explicitly known. It is up to the observer to infer what motivations underlie the given decisions. Since assumptions about motivation provide valuable mental state information that influences our daily interactions, they are worth questioning. How do we make those inferences? The answers are likely more complex given the full context of an interaction involving comparable preferences, so this project started with the simplified game model to provide a starting point for addressing the more complicated, real-world picture.

In particular, the current project addresses two key questions: what level of reasoning we attribute to another person and how that level of reasoning changes with experience. To answer these questions, this project examined the scenario of providing participants with decisions supposedly from previously played games, but not providing all of the relevant payoff values. To study how motivations were inferred, we analyzed whether participants would endorse particular values as plausibly leading to that decision. The investigation into the level of reasoning was exploratory, since there was not enough evidence to inform how participants would react to this new task. However, we suspected that not all participants would reason predictively initially and hypothesized that participants would show increasing use of predictive reasoning as they gained experience with the task.

**Method**

**Participants**

Participants were 47 undergraduate students (17% males, 83% females) at the University of Michigan, who were recruited from the introductory psychology subject pool and received class credit for their participation. All participants gave informed consent before starting the experiment and were debriefed about the origin of the provided "Player II" decisions at the conclusion.

**Design**

**Game design.** The games in this experiment are two-player games that follow the same format as those used in Hedden and Zhang (2002). The four possible outcomes are the four cells "A," "B," "C," and "D." These cells each contain respective payoff values for Player I and Player II. Each player's goal is to get his/her highest possible value, which are independently ranked from "1" as the worst outcome to "4" as the best. The game starts in cell A and ends either at cell D or when a player chooses to "Stop." When the game ends, each player gets his/her respective payoff value from the ending cell. Players have opportunities to make decisions in three sequential stages. At these stages, the player in control can choose to either "Stop" and end the game in the current cell or "Go" to the next cell. Player I has control over the first and third stages; Player II has control over the second stage. A diagram of a sample game and the decision points are shown in Figure 1.

**Training games.** The 12 games in the training session were chosen to be trivial, assigning Player II's values so that there is one clear decision that satisfies both myopic and predictive reasoning, regardless of what Player I decides at the third stage. The two situations that precipitate a "Go" decision are strictly ascending Player II values (1-2-3-4) or a "1" in cell B

(e.g., 3-1-4-2). In these cases, both the cell C and cell D values are better outcomes than the cell B value, so a decision to move is always optimal. The two situations that precipitate a "Stop" decision are strictly descending Player II values (4-3-2-1) or a "4" in cell B (e.g., 2-4-1-3). Conversely, in both of these cases, Player II is guaranteed a better outcome at cell B than at either cell C or cell D, so stopping is always optimal. Player I's values were chosen in conjunction with these strategies to provide equal numbers of games ending in each of the possible cells if players choose the optimal strategies. Information about the training games is shown in Table 1.

Experimental games. Several design concerns were used to determine the following components for the 48 experimental games.

Player I values. Player I's payoff sets were chosen from the 2 x 6 orderings of {1, 2, 3, 4} that start with either a 1 or a 2. These lower start values were selected to provide plausibility that Player I had moved at the first stage and thus Player II was able to complete a turn.

Player II values. In the experimental games, two of Player II's potential payoff values were missing. Since the game action in cell A is only influenced by Player I's decision, this value is not relevant for interpreting Player II's decision in the game. Therefore, in all cases, the cell A value was missing so that there could be two missing values but only one critical cell for analysis. Cell D was not chosen as a critical cell because it may bias participants towards predictive reasoning, since looking ahead to cell D is the mark of a predictive player. The critical cells were chosen to be cell B and cell C because these cells provide crucial information for either a predictive or myopic Player II. The two types of games that had either the cell B value missing or the cell C value missing are considered the "B-Missing" type and "C-Missing" type for the analyses. The sets of provided values were chosen to be {2, 4} or {1, 3}, with the missing

values assumed to complete the set of {1, 2, 3, 4}. There were four patterns containing each

value set for each game type.

     *Diagnostic and non-diagnostic games.* Player I's values determined whether each game

was "diagnostic," meaning that different levels of theory of mind reasoning can be determined

from the participants' responses. The games in which Player I's cell D value is greater than the

cell C value are diagnostic because in these games, Player I would presumably decide to move at

the third stage. A predictive Player II would anticipate this move, but a myopic player would not,

so the difference between these two strategies is observable. The remaining games, in which

Player I has a greater value in cell C than cell D, are "non-diagnostic." In these games, Player I

would presumably not move at the third stage; therefore, predictively anticipating this decision

leads to the same Player II decision as myopically comparing the cell B and cell C values. Non-

diagnostic games were used to balance expected responses and to ensure the participants'

understanding of the game.

     *Decisions.* Although the diagnostic games provide a framework for observing the

different strategies, Player II's values also contribute to whether this difference will be observed.

Within the diagnostic games, each of Player II's payoff patterns has a trivial missing value and a

discriminating missing value. Substituting the trivial value into the critical cell provides similar

patterns to those used for the training games, for which the optimal decision does not depend on

Player I's third stage decision. The "prior Player II" decisions written on the experimental cards

were chosen to correspond to this optimal decision in order to be plausible for either strategy. On

the other hand, substituting the discriminating missing value into the critical cell leads to a

different decision for a myopic Player II than for a predictive one. The experimental questions

were therefore only concerned with the discriminating values. Although the pattern of expected

responses differed by strategy, for both strategies there were equal numbers of games where substituting the discriminating value would have led to the given decision and games where it would have led to the opposite decision.

  *Game categories.* Within the 48 experimental games, there were 8 different game categories: a diagnostic and a non-diagnostic group for each of the 4 patterns of Player II values. The games were presented to the participants in 3 sets of 16 to allow the participant sufficient time to process each game card and question. Each set consisted of 2 games from each of the 8 different categories. The game order was pre-determined using random selection among and within these categories. Participants received these games in a fixed order that appeared entirely random.

  A complete description of the experimental games is shown in Table 2.

## Materials

  The experimental materials consisted of 60 game cards. The cards had borders of different colors to distinguish the training games, the B-Missing games, and the C-Missing games. Each game card was a square containing the 4 cells, A-D. Each cell listed two values in different fonts to distinguish the respective payoff values for Player I and Player II. Twelve of these game cards were used in a training session played with the experimenter to gain familiarity with the game format. The remaining 48 were used as testing materials. The 12 training cards contained full information about each player's payoff values. The 48 experimental cards, on the other hand, had black tape covering two of Player II's possible outcomes. In addition, a decision attributed to a prior Player II was written on the front of these experimental cards. The participant was given a question sheet accompanying the experimental cards that asked one question per game card.

**Procedure**

      **Training session.** To begin the training session, the experimenter instructed the participant on the rules of the game. Following these instructions were a series of games in which the participant played as Player I and the experimenter played as Player II. The training games were played interactively by using a coin to track the progression of the game as determined by each player's decisions. Once an endpoint was reached, the experimenter acknowledged the end of the game and recorded each player's point values on a whiteboard. No feedback about particular strategies was given except if the participant failed to move on the first training game. The first game contained a 1 in cell A, so ending the game at that point was the worst possible outcome for the participant. In these cases, the experimenter would question the decision in order to check understanding about the goal and/or rules. Additionally, the experimenter would provide answers if the participants had further questions.

      **Experimental session.** In the experimental session, the participant was given a set of game cards and told that the decisions written on the cards were made by prior participants who played as Player II at the second stage in these games. These decisions were supposedly made under full information, but some values had been covered by tape for the purpose of this new experiment. The participants were then given an accompanying sheet of questions. For these questions, they were asked to circle yes or no to their belief about whether a particular value could be in that game's critical cell, based on Player II's decision. After giving the instructions, the experimenter remained in a divided section of the room and the participant was told to contact the experimenter upon finishing a set in order to receive the next one. Once the participant had completed all three sets, he/she was asked to fill out an exit questionnaire to

provide further information about his/her reasoning in the games. When the questionnaire was completed, the participant was thanked and debriefed.

**Scoring**

**Predictive scores.** As mentioned in the design, for the diagnostic games, the value in the corresponding question would lead to a different decision by a myopic Player II than for a predictive one. Thus, the participant's theory of mind level can be discerned by their answer to the plausibility of that value. A theory of mind reasoning score was calculated within each set by taking the proportion of diagnostic games answered in agreement with the expected predictive response. Answers opposite from this response can be assumed to be myopic, since that is the logical alternative. Scores were calculated separating the B-Missing type and the C-Missing type of games. Each set contained four diagnostic games of each type, so these scores ranged from 0 to 1 in increments of 0.25.

**Accuracy scores.** Since the myopic and predictive strategies yielded the same response for the non-diagnostic games, there was no logical alternative to the expected answer in these cases. Therefore, responses on these questions were simply scored as "Correct" if the response agreed with the joint predictive/myopic response or "Incorrect" if the response disagreed. Accuracy was calculated as each individual's percentage of correct responses for all 24 non-diagnostic games in the experiment.

<div align="center">

**Results**

</div>

**Data Inclusion**

Among the 47 participants, the decision to remove specific cases was at the discretion of the researcher. Eight of the 47 cases (17%) had a positive response to the item on the exit questionnaire regarding any suspicions that the decisions on the cards did not come from prior

participants. The most common reason for suspecting planned decisions was the valid doubt that a psychology study would provide uncontrolled materials. Although these cases were considered for removal, they were ultimately kept because all participants' exit questionnaire responses described dynamic reasoning with a focus on the opponent's intentions.

Evidence about participants' understanding of the task was taken from the performance on the non-diagnostic games. The distribution of accuracy scores used to determine outliers is shown in Figure 2. Five cases that were outliers on the accuracy distribution were removed due to questionable understanding of the task. This left a total of 42 cases for the remaining analyses. Overall understanding seemed generally strong among the remaining cases, with an average of 92% correct responding.

**Predictive Scores by Game Type**

The percentage of participants with each predictive score, the average predictive scores, and their standard errors are shown for the three sets within the B-Missing type (see Figure 3) and the C-Missing type (see Figure 4). Scores of 0.75 or 1 indicate clear use of the predictive strategy. For both types, the distribution charts indicate a greater proportion of respondents using the predictive strategy in each successive set. Additionally, for both types, the average predictive score increases by set. These results provide preliminary support that the level of reasoning shows an increasing trend among game sets.

In order to investigate whether participants were responding similarly for the two different game types, individuals' overall predictive scores were calculated for each type as the average of the three set scores. These values were plotted as seen in Figure 5. A linear model was fit with the B-Missing score as an explanatory variable and the C-Missing score as a response variable. The scores were found to be highly correlated ($R^2 = 0.85$, $F(1, 40) = 225.1$, $p < .001$).

Additionally, the line of best fit for the model indicates a nearly one-to-one relationship between average predictive scores on the B-Missing type and on the C-Missing type ($\beta = 1.04$, $t(40) = 15.00$, $p < .001$). Based on this model, participants seem to perform similarly on the two different types of games.

   To test the statistical significance of the set and type variables, game set and game type were both entered as within-subjects variables in a repeated-measures ANOVA on the mean prediction score. The main effect of game set was found to be highly significant ($F(2,82) = 13.92$, $p < 0.001$). However, the main effect of game type was not found to be significant ($F(1,41) = 1.84$, $p = 0.18$). These results confirm the exploratory analysis that participants' predictive scores significantly differ among the three sets but not within the two game types. Indicated by the data, this difference in predictive scores is explained by increasingly predictive reasoning in each set. These findings are therefore consistent with the initial hypothesis of predictive reasoning becoming more common as participants gain experience with the task.

**Predictive Scores by Time**

   Since both the linear model and ANOVA results indicate that predictive scores were similar in the two game types, these types were collapsed for the remaining analyses. Predictive scores were recalculated into six time points that divided each set into halves. From the fixed design, each half-set had two diagnostic games of the B-Missing type and two diagnostic games of the C-Missing type to calculate a similar predictive score ranging from 0 to 1 in increments of 0.25. For the six time points, the proportion of respondents with each score, the average predictive scores, and their standard error are shown in Figure 6. The average predictive score at the first time point is 0.51, which indicates that on average, the starting strategy is neither clearly myopic nor clearly predictive. The plots show that the average predictive score substantially

increased between the first and second time points and gradually increased at all remaining time points. These results indicate that the early phase in which participants are first learning the task may be the most important time for establishing whether participants will learn to reason predictively.

## Discussion

Participants were trained as Player I similarly to the Hedden and Zhang (2002) design. However, since participants did not have to make their own decisions in the games, it cannot necessarily be assumed that they took Player I's perspective when analyzing Player II. Nonetheless, participants often described their strategies as transferring their Player I perspective from their training experience, using "I" pronouns to relate to the strategy of the hypothetical past Player I. However, some participants chose to shift to the perspective of Player II, describing analyzing Player II as comparable to "playing as Player II." It was not possible to determine the exact prevalence of these perspectives, since not all participants gave a clear indication of the perspective with which they primarily identified. Predictive reasoning in either of these cases can be described as second-order reasoning, since participants in both cases were reasoning about a prior player anticipating an opponent. The perspective taking demands were therefore comparable to prior studies.

Unlike Hedden and Zhang (2002), this study found a mixed strategy at the outset that tended towards predictive reasoning by the second time point. This is likely in part due to the task demands that specifically instructed participants to analyze a decision by a prior Player II, which may engage theory of mind reasoning more readily than anticipating a player in real time. Although Zhang et al. (2012) found that limits in perspective taking were not solely explained by working memory constraints, this task provided a situation where executive demands were

reduced. Participants' only demands were to answer a question based on a prior player's decision and did not need to decide any moves themselves. This may have freed cognitive resources to devote to higher order reasoning.

Similarly to prior findings, the reasoning level showed a tendency to increase with additional experience. Based on both trends in scores and from participants' descriptions, this often occurred because of a shift from myopic to predictive reasoning. Interestingly, the observed increase in reasoning level occurred even in the absence of feedback. Participants did not get any indication of the accuracy of their choices, unlike real-time players, who would receive lower point scores if they were incorrectly anticipating the opponent. The design contributed to this lack of feedback because decisions were chosen to be plausible for either reasoning strategy. Alternatively, there are certain scenarios where decisions would only be plausible for one of the strategies, regardless of which of the two missing values were substituted. These scenarios were not used because our goal was to explore which strategy was naturally assumed. Since participants did not receive evidence to influence strategy choice, it seems that experience with merely the task of analyzing another person's decisions is sufficient to boost reasoning abilities.

One divergence from prior studies is the presence of the experimenter as the opponent in the training session, as opposed to a confederate posing as a peer to the participant. Although this should not have affected participants' ability to learn the games, it may have provided an early bias when participants were learning to anticipate Player II's strategy. Participants may have more naturally expected the experimenter to anticipate their decisions at the third stage than they would for a peer opponent, since the presumed role of the experimenter is to analyze their performance. This expectation may have transferred to the "prior subjects" when the participants

began the task. Hedden and Zhang (2002) did not find an impact of perceived intelligence of the opponent on performance in their study, but the participants in that case did not have continual face-to-face contact with the opponent, as they did with the experimenter in this study. Therefore, this distinction may also explain the trend of earlier predictive reasoning than seen in some prior studies. This observation leads to the interesting question of whether baseline theory of mind reasoning levels vary between peers and "authority figures" that are assumed to have some degree of meta-knowledge about the task. Since many real-world strategic interactions are indeed with authority figures, this question could have practical implications.

**Limitations**

Initially, we intended to implement this project with a computer opponent giving real-time feedback. Since coding this design was beyond the current project's resources, game cards listing "prior decisions" were used instead. As mentioned, this design limited task constraints so they were not directly comparable with prior studies. Another concern is a disproportionate number of female participants, which occurred by chance based on the introductory psychology students who chose to sign up for the study. This limits the generalizability of the study to both genders, although whether gender affects performance on this task is unclear. A final limitation is that the timing of the project allowed for a smaller sample size than desired. Possible future studies may attempt to replicate these findings with more subjects.

**Conclusions and Future Directions**

The results of this study combined with prior studies provide an insightful picture of our theory of mind capacities in two-player games. We tend to view others as highly rational when directly asked to do so or when put in their position (Zhang et al., 2012). However, we seem to have a harder time attributing full rationality to others in real-time interactions, although this is

more readily done in direct competition (Goodie et al., 2012; Hedden & Zhang, 2002). Importantly, this study supported that predictive reasoning becomes more likely with experience, regardless of feedback on the task. This could potentially extend usefully to social interactions, because our ability to fully understand others' motivations may simply increase with the practice of consciously questioning them.

There are several possible future research questions stemming from this design that could supplement findings on this topic. One option is to use decisions that provide evidence about the player's strategy in order to investigate how participants adapt to this evidence. Another option is to readapt this task to a real-time computer format similar to prior studies, to explore how inferences from real-time decisions differ from analyzing decisions from a more distant past. Strategic gaming paradigms have provided a useful framework for studying theory of mind, but the ecological validity has yet to be determined. It will be very interesting to see how studies that more fully capture the real-world environment relate to these findings.

References

Amsterlaw, J. & Wellman, H. M. (2006). Theories of mind in transition: A microgenetic study of the development of false belief understanding. *Journal of Cognition and Development, 7*(2), 139-172. doi:10.1207/s15327647jcd0702_1

Back, E. & Apperly, I. A. (2010). Two sources of evidence on the non-automaticity of true and false belief ascription. *Cognition, 115*(1), 54-70. doi:10.1016/j.cognition.2009.11.008

Birch, S. A. J. & Bloom, P. (2007). The curse of knowledge in reasoning about false beliefs. *Psychological Science, 18*(5), 382-386. doi:10.1111/j.1467-9280.2007.01909.x

Carbone, E., & Hey, J. D. (2001). A test of the principle of optimality. *Theory and Decision, 50*(3), 263-281. doi:10.1023/A:1010342908638

Colman, A. M. (2003). Cooperation, psychological game theory, and limitations of rationality in social interaction. *Behavioral and Brain Sciences, 26*(2), 139-153. doi:10.1017/S0140525X03000050

Colman, A. M. (2003). Depth of strategic reasoning in games. *Trends in Cognitive Sciences, 7*(1), 2-4. doi:10.1016/S1364-6613(02)00006-2

Converse, B. A., Lin, S., Keysar, B., & Epley, N. (2008). In the mood to get over yourself: Mood affects theory-of-mind use. *Emotion, 8*(5), 725-730. doi:10.1037/a0013283

Epley, N., Morewedge, C. K., & Keysar, B. (2004). Perspective taking in children and adults: Equivalent egocentrism but differential correction. *Journal of Experimental Social Psychology, 40*(6), 760-768. doi:10.1016/j.jesp.2004.02.002

Flobbe, L., Verbrugge, R., Hendriks, P., & Krämer, I. (2008). Children's application of theory of mind in reasoning and language. *Journal of Logic, Language and Information, 17*(4), 417-442. doi:10.1007/s10849-008-9064-7

German, T. P., & Hehman, J. A. (2006). Representational and executive selection resources in

'theory of mind': Evidence from compromised belief-desire reasoning in old age. *Cognition,*

*101*(1), 129-152. doi:10.1016/j.cognition.2005.05.007

Gigerenzer, G. (2001). The adaptive toolbox In G. Gigerenzer & R. Selten (Eds.), *Bounded*

*rationality* (pp. 37-50). Cambridge, MA: MIT Press.

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.

doi:10.1002/bdm.717

Happé, F. G. E., Winner, E., & Brownell, H. (1998). The getting of wisdom: Theory of mind in

old age. *Developmental Psychology, 34*(2), 358-362. doi:10.1037/0012-1649.34.2.358

Hedden, T., & Zhang, J. (2002). What do you think I think you think? strategic reasoning in

matrix games. *Cognition, 85*(1), 1-36. doi:10.1016/S0010-0277(02)00054-9

Kovacs, A. M., Teglas, E. & Endress, A. D. (2010). The social sense: Susceptibility to others'

beliefs in human infants and adults. *Science, 330*(6012), 1830-1830.

doi:10.1126/science.1190792

Li, J., Liu, X., & Zhu, L. (2011). Flexibility of the theory of mind in a matrix game when

partner's level is different. *Psychological Reports, 109*(2), 675-685.

doi:10.2466/04.10.22.PR0.109.5.675-685

Lin, S., Keysar, B., & Epley, N. (2010). Reflexively mindblind: Using theory of mind to interpret

behavior requires effortful attention. *Journal of Experimental Social Psychology, 46*(3),

551-556. doi:10.1016/j.jesp.2009.12.019

Meijering B, van Rijn H, Taatgen NA, Verbrugge R (2012) What Eye Movements Can Tell

about Theory of Mind in a Strategic Game. PLoS ONE 7(9): e45961.

doi:10.1371/journal.pone.0045961

Mitchell, P., Robinson, E. J., Isaacs, J. E., & Nye, R. M. (1996). Contamination in reasoning

about false belief: An instance of realist bias in adults but not children. *Cognition, 59*(1), 1-

21. doi:10.1016/0010-0277(95)00683-4

Nickerson, R. S. (1999). How we know—and sometimes misjudge—what others know:

Imputing one's own knowledge to others. *Psychological Bulletin, 125*(6), 737-759.

doi:10.1037/0033-2909.125.6.737

Onishi, K. H., & Baillargeon, R. (2005). Do 15-Month-Old Infants Understand False Beliefs?

*Science*, *308*(5719), 255–258. doi:10.1126/science.1107621

Osborne, M. J. & Rubinstein, A. (1994). *A course in game theory*. Cambridge, MA: MIT Press.

Perner, J., Leekam, S. R., & Wimmer, H. (1987). Three-year-olds' difficulty with false belief:

The case for a conceptual deficit. *British Journal of Developmental Psychology, 5*(2), 125-

137. doi:10.1111/j.2044-835X.1987.tb01048.x

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), 437-471. doi:10.1016/0022-0965(85)90051-7

Samson, D., Apperly, I. A., Braithwaite, J. J., Andrews, B. J. & Scott, S. E. B. (2010). Seeing it

their way: Evidence for rapid and involuntary computation of what other people see.

*Journal of Experimental Psychology.Human Perception and Performance, 36*(5), 1255-

1266. doi:10.1037/a0018729

Saxe, R, Schulz, L. E., & Jiang, Y. V. (2006). Reading minds versus following rules:

Dissociating theory of mind and executive control in the brain. *Social Neuroscience, 1*(3-4),

284-298. doi:10.1080/17470910601000446

Zhang, J., Hedden, T and Chai, A. (2012). Perspective-taking and depth of theory of mind

    reasoning in sequential-move games. *Cognitive Science*, *36*(3), 560-573.

    doi:10.1111/j.1551-6709.2012.01238.x
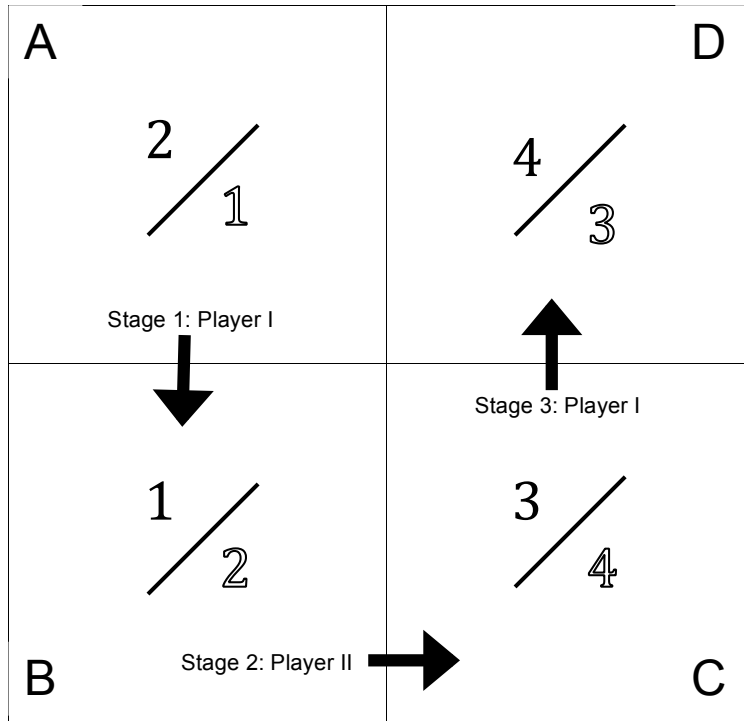
Table 1

*List of training block games*

| Player I | Player II | Player II Decision | End Point |
|----------|-----------|--------------------|-----------|
| 1324 | 1234 | Go | D |
| 4321 | 2413 | Stop | A |
| 3412 | 4321 | Stop | B |
| 2431 | 3142 | Go | C |
| 3241 | 3421 | Stop | A |
| 3421 | 2413 | Stop | B |
| 1432 | 4123 | Go | C |
| 2134 | 1243 | Go | D |
| 3124 | 4312 | Stop | A |
| 1234 | 4321 | Stop | B |
| 2143 | 1234 | Go | C |
| 3214 | 3142 | Go | D |

*Note.* This table lists each game's payoff sets by cell (in the order of ABCD) for Player I and

Player II, the consistent decision made by the experimenter, and the expected end point if both
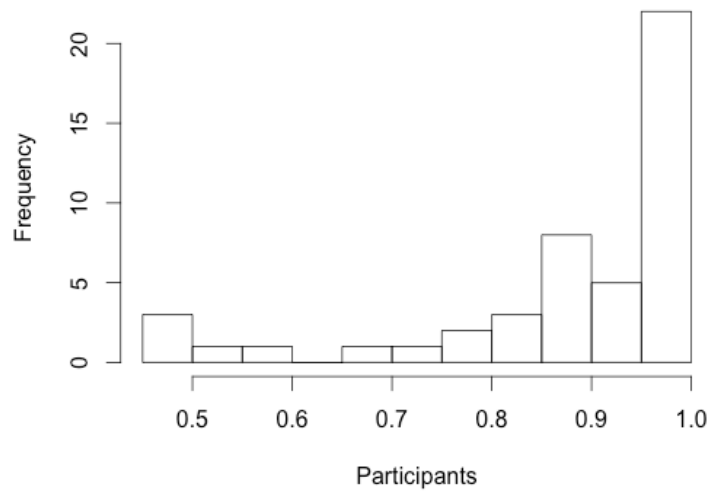
players decide optimally.

Table 2
*List of experimental games*

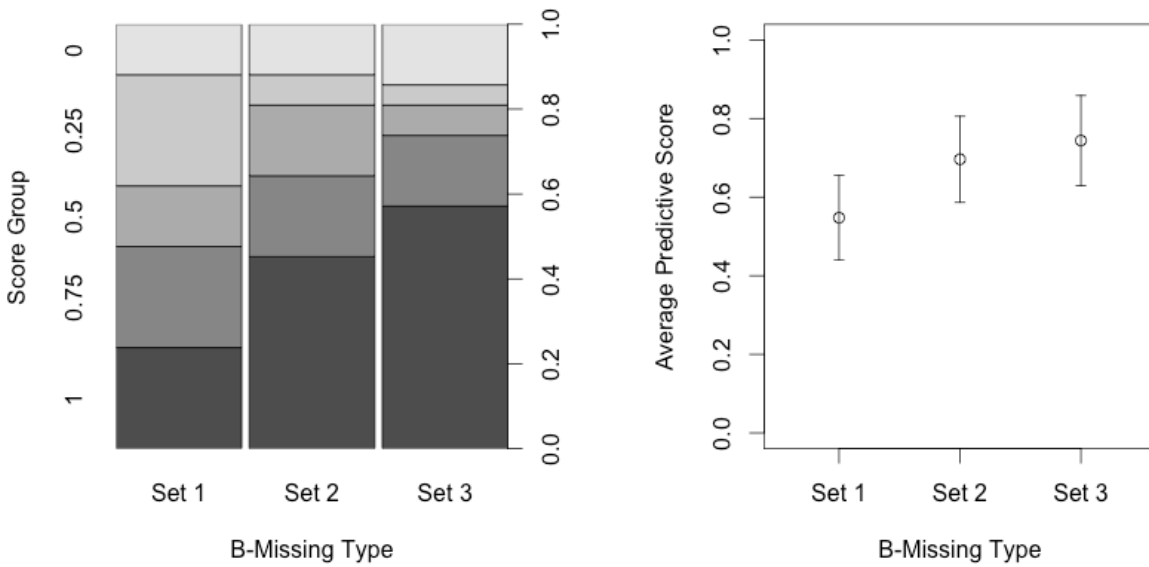| | B-Missing Type | | | | | |
|---|---|---|---|---|---|---|
| | Pattern 1 | | | Pattern 2 | | |
| | ?-?- 1- 3 | | | ?-?- 4- 2 | | |
| | Question: Could there be a 2 in cell B? | | | Question: Could there be a 3 in cell B? | | |
| Diagnostic | ID | Decision | Predictive Response | ID | Decision | Predictive Response |
| 1234 | 10 | STOP | No | 20 | GO | No |
| 1324 | 11 | STOP | No | 21 | GO | No |
| 1423 | 12 | STOP | No | 22 | GO | No |
| 2314 | 13 | STOP | No | 23 | GO | No |
| 2413 | 14 | STOP | No | 24 | GO | No |
| 2134 | 15 | STOP | No | 25 | GO | No |
| Non-Diagnostic | | | Accurate Response | | | Accurate Response |
| 1243 | 16 | STOP | Yes | 26 | GO | Yes |
| 1342 | 17 | STOP | Yes | 27 | GO | Yes |
| 2341 | 18 | STOP | Yes | 28 | GO | Yes |
| 2431 | 19 | STOP | Yes | 29 | GO | Yes |
| 1432 | 50 | STOP | Yes | 52 | GO | Yes |
| 2134 | 51 | STOP | Yes | 53 | GO | Yes |
| | C-Missing Type | | | | | |
| | Pattern 3 | | | Pattern 4 | | |
| | ?- 2- ?- 4 | | | ?- 3- ?- 1 | | |
| | Question: Could there be a 1 in cell C? | | | Question: Could there be a 4 in cell C? | | |
| Diagnostic | ID | Decision | Predictive Response | ID | Decision | Predictive Response |
| 1234 | 30 | GO | Yes | 40 | STOP | Yes |
| 1324 | 31 | GO | Yes | 41 | STOP | Yes |
| 1423 | 32 | GO | Yes | 42 | STOP | Yes |
| 2314 | 33 | GO | Yes | 43 | STOP | Yes |
| 2413 | 34 | GO | Yes | 44 | STOP | Yes |
| 2134 | 35 | GO | Yes | 45 | STOP | Yes |
| Non-Diagnostic | | | Accurate Response | | | Accurate Response |
| 1243 | 36 | GO | No | 46 | STOP | No |
| 1342 | 37 | GO | No | 47 | STOP | No |
| 2341 | 38 | GO | No | 48 | STOP | No |
| 2431 | 39 | GO | No | 49 | STOP | No |
| 1432 | 54 | GO | No | 56 | STOP | No |
| 2134 | 55 | GO | No | 56 | STOP | No |

*Note.* This table lists each game's payoff sets by cell (in the order of ABCD) for Player I and Player II ("?" indicates a covered value). The decisions in this table were provided on the game card. In diagnostic games, the myopic response is the opposite response from the predictive. In non-diagnostic games, the accurate response applies to both strategies. The same question was used within each pattern, focusing on the discriminatory missing value and the critical cell.
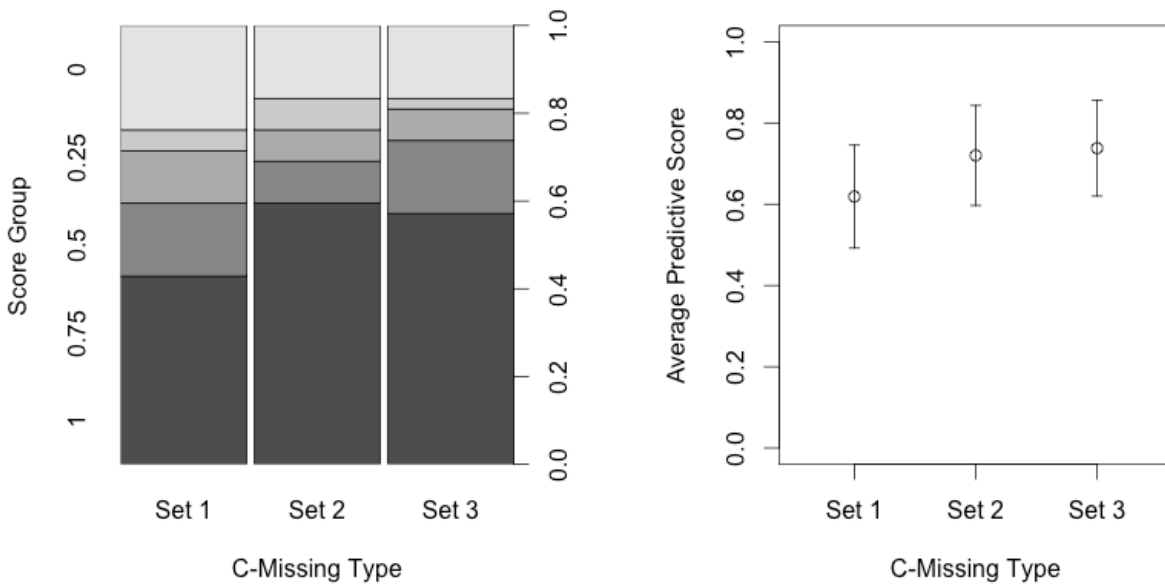
*Figure 1.* Diagram of sample game with stages. The top values correspond to Player I's payoff values and the bottom values correspond to Player II's payoff values. The arrows indicate the potential decision at each of the game stages.
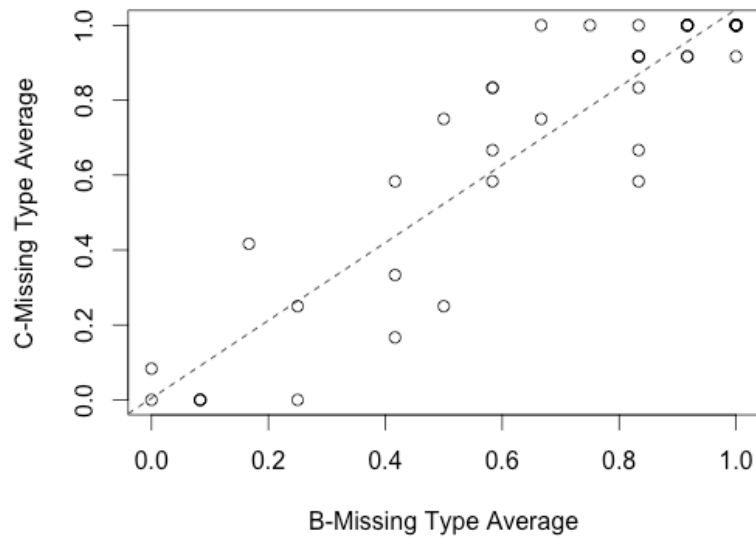
*Figure 2.* Accuracy distribution of all participants. Most participants had a total accuracy score above 65% correct, but there were a few outliers with scores lower than this percentage. The outliers have been removed from further analysis.
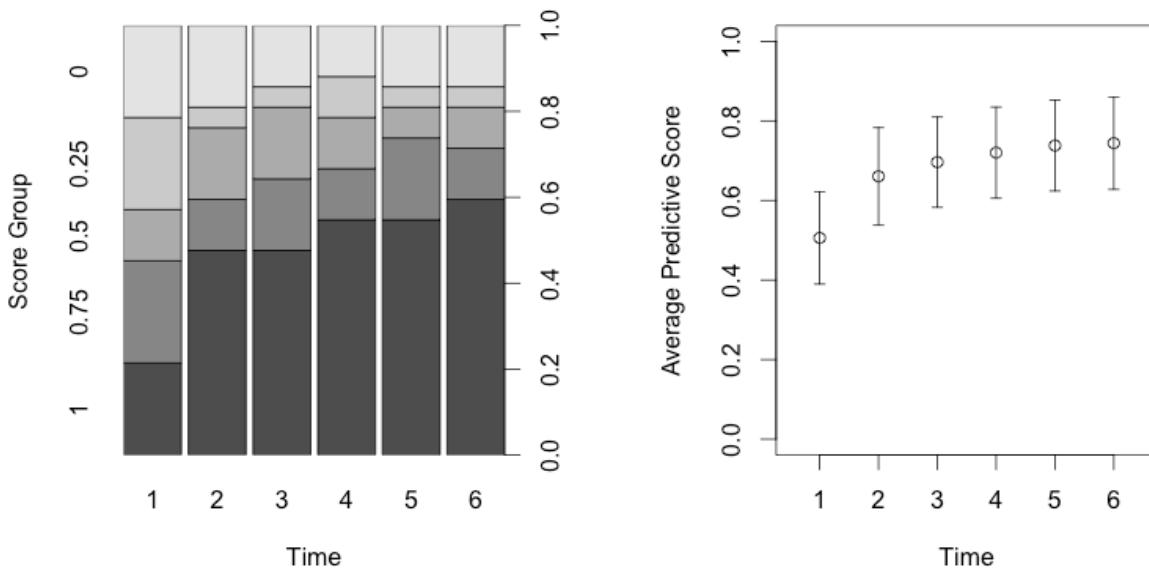
*Figure 3.* Predictive scores by set for the B-Missing type of game. The first chart shows the distribution of players at each possible score for the three sets. The second chart shows the average predictive score for each set, with bars indicating the standard error. The average predictive score increases in each successive set, primarily explained by the increasing proportions of respondents with a perfect "1" predictive score.

*Figure 4.* Predictive scores by set for the C-Missing type of game. The first chart shows the

distribution of players at each possible score for the three sets. The second chart shows the

average prediction score for each set, with bars indicating the standard error. The average

predictive score increases in each successive set, although the increase in the third set is less

notable than for the B-Missing type. The increase in the second set can be explained by a greater

proportion of respondents with a perfect "1" predictive score and the smaller increase in the third

set can be explained by a greater proportion of respondents with the second highest "0.75" score.

*Figure 5*. Plot of each individual's overall average predictive scores on the B-Missing and C-Missing types. The fitted model was C-Missing Type Average = 0.00 + 1.04 B-Missing Type Average, $R^2$ = .85, $p$ < .001. Individuals seem to perform similarly on each game type.

*Figure 6.* Predictive scores by time point. Each time point was half of a set and included two games from each type. The first chart shows the distribution of players at each possible score for each time point. The second chart shows the mean prediction score for each time point, with bars indicating the standard error. Although all time points show an increase in reasoning score, the most substantial increase seems to be between the first and second time point. Proportion of respondents with perfect "1" predictive scores shows a similar trend.