

# **Reasoning about self and others**

**Ben Meijering**

Printed by Gildeprint Drukkerijen,  
Enschede, the Netherlands

ISBN printed version: 978-90-367-7062-0  
ISBN digital version: 978-90-367-7061-3

© Ben Meijering, Groningen, The Netherlands, 2014



university of  
 groningen

# Reasoning about self and others

## PhD thesis

to obtain the degree of PhD at the  
 University of Groningen  
 on the authority of the  
 Rector Magnificus Prof. E. Sterken  
 and in accordance with  
 the decision by the College of Deans.

This thesis will be defended in public on

Friday 6 June 2014 at 14.30 hours

by

**Ben Meijering**

born on 5 February 1982  
 in Zevenaar

**Supervisors**

Prof. L.C. Verbrugge

Prof. N.A. Taatgen

**Co-supervisor**

Dr. H. van Rijn

**Assessment committee**

Prof. J. Perner

Prof. M.E.J. Raijmakers

Prof. P. Hendriks

# Contents

<i>Chapter 1</i>	
<i>Introduction</i>	6
<i>Chapter 2</i>	
<i>Integrating recursive application of theory of mind in decision making in sequential games</i>	10
<i>Chapter 3</i>	
<i>Reasoning about self versus others: Changing perspective is hard</i>	26
<i>Chapter 4</i>	
<i>Reasoning about diamonds, physics, and mental states: The cognitive costs of theory of mind</i>	38
<i>Chapter 5</i>	
<i>What eye movements can tell about theory of mind in a strategic game</i>	48
<i>Chapter 6</i>	
<i>Modeling inference of mental states: As simple as possible, as complex as necessary</i>	62
<i>Chapter 7</i>	
<i>Summary and discussion</i>	76
<i>Samenvatting</i>	82
<i>Dankwoord</i>	85
<i>References</i>	86
<i>Publication list</i>	92

# **Chapter 1**

## **Introduction**

## Theory of mind: Understanding others

We live in a social world in which we frequently interact with others. In our jobs, for example, we may collaborate with colleagues and negotiate with superiors, and in our leisure time we may compete with friends when playing a board game. These different scenarios have in common that in each of them we are trying to understand one another. For example, to collaborate with a colleague we need to know what her goals are to be able to find commonalities, as a basis for working together. To compete with a friend in a board game we need to know what her goals are to be able to anticipate her actions. In either case, we are reasoning about goals, which happen to be intangible, hidden to the eye. In spite of this characteristic of goals, human beings are quite proficient at inferring goals and other so-called mental states such as beliefs, desires, and intentions. In fact, the ability to infer mental states of others seems to be unique to human intelligence (Call & Tomasello, 2008). This ability has often been referred to as theory of mind (Onishi & Baillargeon, 2005; Premack & Woodruff, 1978; Wimmer & Perner, 1983), and because the ability to infer mental states of others is pivotal to human nature, many studies investigate the nature and development of theory of mind (Wellman, Cross, & Watson, 2001).

### Studying theory of mind

The most influential paradigm to investigate theory of mind is the false-belief task, also known as the Sally-Anne task (Baron-Cohen, Leslie, & Frith, 1985; Wimmer & Perner, 1983). In this task, children listen to a story with two characters: Sally and Anne. Sally is playing with a marble, and before she leaves the room she places the marble in her basket. In Sally's absence, Anne moves the marble and hides it in another location, a box. Then, Sally returns and children are asked where she will look for her marble. To pass this task, children need to understand that Sally will look for the marble in the basket, whereas in fact the marble is now in the box.

Some developmental studies have shown that infants as young as 15 months are already susceptible to the mental states of others. Onishi and Baillargeon (2005), for example, have shown that 15-month-old infants are surprised when their expectation that an agent will act according to her beliefs is violated. In Onishi and Baillargeon's false belief task, children saw that an agent, who believed that a slice of watermelon was stored in a box, was unaware of the watermelon being moved to a second box. When the agent reached for the watermelon, after it had been moved, the children looked longer at the scene if the agent reached for the second box instead of the first. The children looked longer because the agent did not act according to her false belief that the watermelon was hidden in the first box. Importantly, Onishi and Baillargeon administered an implicit test of theory of mind, using the violation of expectation paradigm, which did not require the infants to explicitly reason about the agent's mental states. As the violation of expectation paradigm does not require explicit reasoning, the infants' understanding of mental states would not be obscured by cognitive functions that are yet to be developed.

By the age of 4, children have sufficiently developed cognitive functions to pass the original false belief task (Wellman et al., 2001; Wimmer & Perner, 1983). This task requires a deliberate

response, as children are explicitly asked where the agent will look for the toy. A deliberate response is an additional step that involves effortful processing, which is clearly illustrated by a dissociation between two distinct measures of theory of mind within the same sample (Clements & Perner, 1994; Ruffman, Garnham, Import, & Connolly, 2001). In Ruffman et al.'s study, for example, children looked at the correct location, where the agent falsely believed the toy should be, but answered incorrectly that the agent would look at the toy's actual location. To this day, this remarkable finding still does not have one succinct explanation, and there are many studies investigating which distinct processes constitute the developmental trajectory of theory of mind.

Roughly speaking, there are two clearly opposing camps that have different ideas about which processes are involved in the development of theory of mind. One camp argues that children's understanding of mental states undergoes a conceptual change (Gopnik & Wellman, 1992). Not until children understand that, for example, beliefs do not need to reflect reality, do they understand that others can have beliefs that differ from their own. The other camp argues that children first need to develop other important cognitive functions that, once sufficiently developed, will enable inference of mental states (e.g., Carlson, Moses, & Breton, 2002; Sabbagh, Xu, Carlson, Moses, & Kang, 2006). These cognitive functions are required to support the costly computational processes that are involved in mental state inference. Of course, children still need to be able to discern mental states, for example, by means of a so-called theory of mind mechanism (e.g., Leslie, Friedman, & German, 2004), or a minimal understanding of concepts in general. Accordingly, there seems to be merit in hybrid theories as well, combining the theories of the two camps delineated above, as application of theory of mind seems to be an effortful process that requires sufficiently complex mental state representations.

Given that application of theory of mind is a costly computational process, it may not be surprising that it takes another two years before children, of 5 to 6 years old, start understanding that others, too, apply theory of mind (e.g., Miller, 2009). They learn to understand recursive structures such as "John thinks that Mary thinks that..." (Perner & Wimmer, 1985). Comprehension of such structures requires application of second-order theory of mind, which is a challenge for both children and adults (Flobbe, Verbrugge, Hendriks, & Krämer, 2008; Hedden & Zhang, 2002; Perner & Wimmer, 1985; Qureshi, Apperly, & Samson, 2010; Raijmakers, Mandell, Van Es, & Counihan, 2013). Application of higher orders of theory of mind would exceed the cognitive resources of most people. The dynamics of declarative and procedural memory, for example, would probably not support such complex reasoning. Fortunately, in many circumstances second-order theory of mind seems to be the highest level of theory of mind that is still advantageous for people to use (De Weerd, Verbrugge, & Verheij, 2013).

## **Contribution of this dissertation**

This thesis describes a detailed investigation of theory of mind in adults, who have sufficient declarative and procedural knowledge of mental states to interpret the behavior of self and others, in contrast to infants and children. The thesis therefore contributes to the field in at least one obvious way: As adult understanding of mental states has undergone the earlier



proposed conceptual changes (Gopnik & Wellman, 1992), the findings reported here can be interpreted as a measure of the cognitive functions and resources required to apply theory of mind. In other words, the findings in this thesis help shed light on the procedural building blocks that constitute inference of mental states. A better understanding of the procedural building blocks is important, because it may be of help in both the clinical and practical domains.

Another contribution of this thesis is that it shows how two-player games may be a good alternative to test for theory of mind. In contrast to false-belief stories, games can be presented many times and in many different configurations to the same individuals. The assumption that these games are interpreted in terms of mental states remains implicit in some chapters (Chapters 2, 5, and 6). However, Chapters 3 and 4 explicitly test whether two-player games require a theory of mind. It turns out that these games are especially successful in varying demands on mental state reasoning.

The thesis starts with a broad question in Chapter 2. The question is whether theory of mind is a fixed skill, or an ability that is susceptible to improvement by means of supportive measures that help structure inference of mental states. Chapters 3 and 4 detail an investigation into the computational costs associated with taking either one's own or someone else's perspective. Most definitions of theory of mind state that it is an understanding of mental states of self and others. However, few studies seem to investigate the difference between taking one's own and someone else's perspective. Chapters 5 and 6 provide a detailed analysis of particular strategies that people use when they are reasoning about others' mental states. These studies show that people rather use simple strategies that may not be optimal but yield correct responses in most cases. Only when they really have to, do people invest their resources into forming complex mental state representations. In other words, people may have the ability to infer complex mental states, but still fail or choose not to use it.

## Chapter 2

# Integrating recursive application of theory of mind in decision making in sequential games

### Abstract

In collaborative, competitive, and negotiation situations we need to reason about each other's goals, intentions, beliefs, and desires, which requires a so-called Theory of Mind (ToM). This study investigates decision making in which ToM has to be applied recursively: "I think that you think that I think...". Participants were presented with sequential games in which the payoff for one player depended on the decisions of another player. Previous studies typically found suboptimal decisions in these games. One possible explanation is an overall inability to apply recursive ToM. Instead, we argue that suboptimal decisions are caused by unsuccessful integration of recursive ToM in the decision making process. This hypothesis is tested by means of three experimental manipulations that each should facilitate this integration process. First, each player's decision options are introduced and explained in a stepwise fashion during the training phase. Second, participants are prompted to predict the other player's decision. Third, participants are presented with a concrete and realistic task representation that visually cues the recursive structure of the decision making problem. The results show that performance was better in those conditions that specifically targeted the integration of recursive mental states in the decision making process.

Parts of this chapter were previously published in the *proceedings of the 32nd and 33rd Annual Conference of the Cognitive Science Society* (2010, 2011).

## Introduction

Making decisions can be a complicated process, especially when the actions of others have to be factored in. To understand, or even predict, the actions of others we need to infer their goals, beliefs, desires, et cetera. It is difficult to infer these mental states because we cannot directly observe them. Still, children already learn to reason about the mental states of others from the age of 4. They develop a so-called theory of mind (ToM), an understanding that others may have beliefs that differ from their own and do not need to reflect reality (e.g., Wellman, Cross, & Watson, 2001; Wimmer & Perner, 1983). However, making decisions based on mental state inferences remains challenging, even for adults (e.g., Flobbe, Verbrugge, Hendriks, & Krämer, 2008; Hedden & Zhang, 2002; Zhang, Hedden, & Chia, 2012).

Developmental studies have shown that children typically make optimal decisions if they are required to infer relatively simple mental states (Flobbe et al., 2008; Raijmakers, Mandell, Van Es, & Counihan, 2013). In Raijmaker et al.'s (2013) sequential games, for example, children performed well if they had to infer mental states such as “The other player *intends* to stop the game when it is his turn”. In contrast, few children were able to make optimal decisions if they were required to infer more complex mental states such as “The other player *believes* that I *intend* to continue when it is my turn”. Most children could not incorporate such second-order, or recursive, beliefs into their decision making process (for similar findings see Flobbe et al., 2008).

Some studies have shown that adults, too, find it difficult to make decisions based on second-order mental states (Flobbe et al., 2008; Hedden & Zhang, 2002; Johnson, Camerer, Sen, & Rymon, 2002; McKelvey & Palfrey, 1992; Verbrugge & Mol, 2008; Zhang et al., 2012). In Hedden and Zhang's (2002) sequential games, for example, adult participants initially did not seem to realize that the other player was reasoning about them. As a consequence they made suboptimal decisions that were based on inaccurate mental state representations.

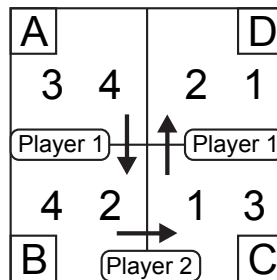


Figure 2.1: Example of a matrix game; adapted from Hedden and Zhang (2002); labels (A – D, Player 1 / 2) and arrows are added for illustrative purposes and were not depicted during the experiment.

Figure 2.1 depicts an example of Hedden and Zhang's sequential games, which are so-called matrix games. Each cell of a matrix game contains a pair of rewards, or so-called payoffs, that range from 1 to 4. The left payoff of a pair is Player 1's payoff; the right payoff is Player 2's. Both players alternately decide whether to stop the game in the current cell, or to continue it

to the next. Each player's goal is that the game stops in the cell that contains his or her highest possible payoff, irrespective of the payoff for the other player. Based on the assumption that each player is rational, both players should know that their outcome depends on the decisions of the other player: Either player can decide to stop or continue the game when it is their turn. Crucially, as each player's decision is based on beliefs about the next decision, recursive beliefs have to be incorporated into the decision making process.

In the game in Figure 2.1, for example, Player 1 should decide to continue the game from A to B, as she could have inferred that a rational Player 2 will decide to stop the game in that cell. Player 2 would not continue to C, as he would have inferred that a rational Player 1 would continue to D, which contains a smaller payoff for him than B. As B contains a higher payoff for Player 1 than A, Player 1 should decide to continue from A to B. In sum, Player 1 needs to reason about Player 2's belief about Player 1's intention, thus applying second-order ToM by attributing first-order ToM to Player 2.

Most of the earlier mentioned studies concluded that adults and children do not have sufficient ability or cognitive capacity to apply second-order ToM (Flobbe et al., 2008; Hedden & Zhang, 2002; Raijmakers et al., 2013; Verbrugge & Mol, 2008; Zhang et al., 2012). Lacking the ability to apply ToM recursively, these studies claim, one makes suboptimal decisions that are based on inaccurate mental state representations. However, as children typically perform well in false-belief tasks that require them to keep track of stories with multiple agents that each have their own set of recursive beliefs (e.g., Apperly, Back, Samson, & France, 2008; Apperly et al., 2010; Hollebrandse, Hobbs, De Villiers, & Roeper, 2008; Perner & Wimmer, 1985; Sullivan, Zaitchik, & Tager-Flusberg, 1994; Wimmer & Perner, 1983), it seems that in some tasks sufficient ability or capacity is available to apply second-order ToM. We therefore hypothesize that suboptimal decisions are due to unsuccessful integration of recursive mental states in the decision making process. In many sequential games, participants have to attribute recursive mental states to the other player and, on top of that, they have to combine these attributions to choose their own best possible course of action. Suboptimal decisions arise if the integration of recursive mental states breaks down, even if the second-order representations are accurate.

To test whether suboptimal decisions in sequential games are due to unsuccessful integration of second-order mental states, we devised three experimental manipulations that each should scaffold the integration of mental states: (1) To make participants aware of the dependency between all decision points, each subsequent decision is introduced and explained in a stepwise fashion during the training phase; (2) To train integration of the other player's mental states, participants are prompted to predict the other player's decision; (3) To make the recursive structure of the decision making problem more salient, participants are presented with a new visual task representation. All three manipulations should clarify the "who, what, where" while making decisions in sequential games.

## Method

### Participants

Ninety-three first-year psychology students (63 female) with a mean age of 21 years (ranging between 18 and 31 years) participated in exchange for course credit. All participants had normal or corrected-to-normal visual acuity. Informed consent as approved by the Ethical Committee Psychology of the University of Groningen was obtained before testing.

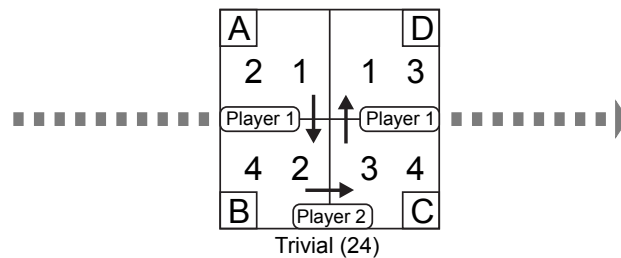
### Design

The experimental design comprised three factors: training, prompting predictions, and task representation. All factors were administered between participants. The experiment consisted of two phases: a training phase, followed by a test phase.

#### Training

The training phase was included to familiarize participants with the rules of sequential games. Participants were randomly assigned to one of two training procedures. In one training procedure participants were presented with Hedden and Zhang's (2002) 24 original training

#### Undifferentiated Training



#### Stepwise Training

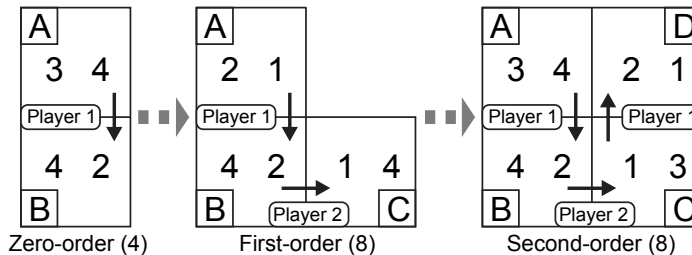


Figure 2.2: Schematic overview of the Undifferentiated and Stepwise training procedures. Undifferentiated training consists of 24 so-called trivial games (see text for explanation). Stepwise training consists of 4 zero-order games, 8 first-order games, and 8 second-order games. The actual training items all had different payoff distributions.

games (see Figure 2.2; top panel). These training games are considered easier to play than truly second-order games such as in Figure 2.1, because Player 2 does not have to reason about Player 1's last possible decision: Player 2's payoff in B is either lower or higher than both his payoffs in C and D. Consequently, Player 1 does not have to attribute ToM to Player 2. These games are therefore referred to as trivial games by Hedden and Zhang. This training procedure will henceforth be referred to as Undifferentiated training, as all games have three decision points.

In the other training procedure participants were presented with three blocks of games that are simple at first and become increasingly more complex with each subsequent block (see Figure 2.2; bottom panel). This procedure will henceforth be referred to as Stepwise training. The first block consisted of 4 games with just one decision point. These games are so-called zero-order games, as they do not require application of ToM. The second block consisted of 8 games with two decision points. These games require application of first-order ToM, as the participant is required to reason about the other player. The third block consisted of 8 games with three decision points that require application of second-order ToM, as the participant has to reason about the other player, and take into account that the other player is reasoning about them.

We hypothesize that the Stepwise training procedure provides scaffolding to support representation of increasingly more complex mental states. Stepwise introduction, explanation, and practice of each additional decision point helps participants integrate mental states of increasing complexity into their decision making process.

### *Prompting predictions*

The second factor, prompting participants for predictions, was manipulated in the test phase. Hedden and Zhang (2002) prompted their participants to predict Player 2's decision (in B), before making a decision themselves. By prompting participants for predictions, participants were explicitly asked to take the other player's perspective, and we hypothesize that these prompts helped participants to integrate the other player's perspective in their decision making process.

We tested this hypothesis by means of two test blocks. In the first, we asked half of the participants, assigned to the Prompt group, to predict Player 2's move before making their own decision. Participants assigned to the No-Prompt group, in contrast, were not explicitly asked to predict Player 2's move. The second test block was added to test whether prompting had long-lasting effects on performance. No participant was asked to make predictions, and performance differences between the two groups would indicate lasting effects of prompting.

### *Task representation*

The third factor is task representation. Before the training phase started, participants were assigned to one of two task representations, which did not change anymore during the remainder of the experiment. The Matrix Game was one of these task representations, and we devised a second. The new task representation was devised to clarify the recursive structure of the decision making problem. Depicted in Figure 2.3, the new task representation shows a more intuitive display of who decides where and what the consequences of each decision are. We argue that this new representation, henceforth referred to as Marble Drop, provides scaffolding that supports the integration of decisions and underlying mental states.

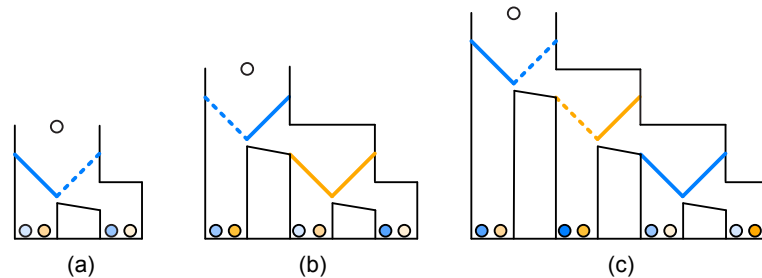


Figure 2.3: Examples of zero-order (a), first-order (b), and second-order (c) Marble Drop games. The blue player (i.e., Player 1) has to obtain the darkest possible blue marble, the orange player (i.e., Player 2) the darkest possible orange marble. The dashed lines are added for illustrative purposes and represent the trapdoors that proficient ToM players should remove to obtain their darkest possible marble. See text for additional explanation.

Importantly, this representation is isomorphic to matrix games and thus requires the same reasoning.

Figure 2.3 depicts examples of zero-order, first-order, and second-order Marble Drop games. Participants are told that a white marble is about to drop, and that its path can be manipulated by removing trapdoors. Their goal is to let the white marble drop into the bin containing the darkest possible marble of their target color, blue in these example games, by controlling only the blue trapdoors. Player 2's goal is to obtain the darkest possible orange marble, but Player 2 can only control the orange trapdoors. The marbles are ranked from light to dark, with darker marbles preferred over lighter marbles, yielding payoff structures isomorphic to those in matrix games.

## Stimuli

### *Payoffs*

The payoffs in matrix games are numerical, ranging from 1 to 4, whereas the payoffs in Marble Drop games are color-graded marbles that have a one-to-one mapping to the numerical values in the matrix games. The colors of the marbles are four shades of orange and blue, taken from the HSV (i.e., hue, saturation and value) space. A sequential color palette is computed by varying saturation, for a given hue and value. This results in four shades (with saturation from .2 to 1) for each of the colors orange (hue = .1, value = 1) and blue (hue = .6, value = 1).

### *Payoff structures*

The payoff structures are selected so that the order of ToM reasoning mastered by the participants can be derived from their first decision<sup>1</sup>. The total set of payoff structures, balanced for the number of decisions to continue or stop a game, is limited to 16 items. These items are listed in Appendix A. Detailed discussion of the rationale behind the exclusion

<sup>1</sup> We look at the entire set of first decisions, as an individual decision cannot discriminate between multiple strategies. For example, in one particular game both guessing and applying second-order ToM might yield a correct decision. However, by looking at the entire set of decisions, we can discriminate guessing from applying second-order ToM, as guessing would only yield a correct decision in 50% of the games.

criteria is given in Appendix B.

## Procedure

To familiarize participants with the rules of sequential games, they were first presented with a training block that either consisted of Undifferentiated training or Stepwise training. The instructions, which appeared on screen, explained how to play sequential games and what the goal of each player was. The instructions also mentioned that participants were playing against a computer-simulated player, as Hedden and Zhang (2002) have shown that inclusion of a cover story did not affect ToM performance. Each training game was played until either the participant or the computer-simulated player decided to stop, or until the last possible decision was made. After each training game participants were presented with accuracy feedback indicating whether the highest attainable payoff was obtained. In case of an incorrect decision, an arrow pointed at the cell / bin that contained the highest attainable payoff. As the feedback never referred to the other player's mental states, participants had to infer these themselves.

The two test blocks consisted of second-order games. As mentioned above, the procedure for participants in the Prompt and No-Prompt groups differed in the first test block. Participants in the Prompt group were first asked to enter a prediction of Player 2's decision before they were asked to make a stop-or-continue decision at their own decision point. Participants in the No-Prompt group, in contrast, were not asked to make predictions. They were only asked to make decisions. Accuracy feedback was presented both after entering a prediction and after entering a decision, but the arrow was not shown anymore in the test blocks. This block consisted of 32 trials; each of the 16 payoff structures was presented twice, but not consecutively as the items were presented randomly.

The second test block was the same for all participants. They were asked to make decisions only.

## Results and discussion

To account for multiple sources of random variation (i.e., participants and payoff structures were both sampled from a larger population), the data were analyzed by means of linear mixed-effects (LME) models (Baayen, 2008; Baayen, Davidson, & Bates, 2008). We included random intercepts to allow the intercepts of the regression models to vary across participants and items. Random slopes were included to allow the effects (i.e., slopes) of training, prompting, and task representation to vary across items (Barr, Levy, Scheepers, & Tily, 2013). The correctness of the decisions was analyzed by means of *logistic* LMEs, as correctness of decisions is a binary variable. These models are provided by the *lme4* package (version 0.999375-42; Pinheiro & Bates, 2000) in R ([www.r-project.org](http://www.r-project.org), version 3.0.1). Two separate models were fit to the data because the factors training and prompting were manipulated in different blocks. All figures depict means and standard errors, which are represented by error bars.



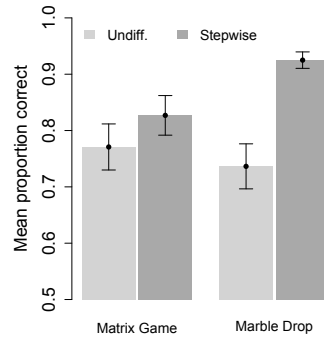


Figure 2.4. Mean proportions correct decisions; depicted separately for Undifferentiated training (light grey) and Stepwise training (dark grey) for both Matrix Games and Marble Drop games.

## Effect of training and representation

Type of training was manipulated before either test block was administered. Its effect on the correctness of decisions was analyzed in the first test block, together with the effect of task representation. Figure 2.4 depicts the mean proportions of correct decisions.

As can be seen in Figure 2.4, the effect size of Stepwise training (in contrast to Undifferentiated training) significantly varied with task representation,  $\chi^2(1) = 6.4$ ,  $p = 0.011$ . Figure 2.4 shows that stepwise training had a significant positive effect on the correctness of decisions,  $\chi^2(1) = 12.1$ ,  $p < 0.001$ , but this effect was driven by participants assigned to Marble Drop games,  $\beta = 2.063$  (SE = 0.491),  $z = 4.203$ ,  $p < 0.001$ . In fact, the effect of stepwise training was not significant for participants assigned to Matrix Games,  $\beta = .427$  (SE = .432), ns.

Even though there was no significant main effect of task representation, participants assigned to stepwise training did have a significantly higher probability of making a correct decision if they were assigned to Marble Drop games instead of Matrix Games,  $\beta = 1.307$  (SE = 0.483),  $z = 2.707$ ,  $p = 0.007$ . Participants assigned to undifferentiated training did not have an advantage if they were assigned to Marble Drop games,  $\beta = -.326$  (SE = .439), ns. In other words, the effectiveness of the Marble Drop task representation, positive in any case, varied with the type of training received by the participants.

In sum, both the Stepwise training procedure and the Marble Drop task representation positively affected the probability of making a correct decision, the latter factor primarily by means of an interaction.

## Effect of prompting

As prompting was manipulated in the first test block, we analyzed its short-term effect during the first test block and its longer-term effect in the second test block. Figure 2.5 depicts the accuracy, or proportion correct, of the decisions and predictions.

Log-likelihood ratio comparisons indicated that the factor task representation did not make the models fit the decisions better. Therefore, we report the statistics of a full factorial model that includes terms for *prompting*, *test block*, and an interaction between the two.

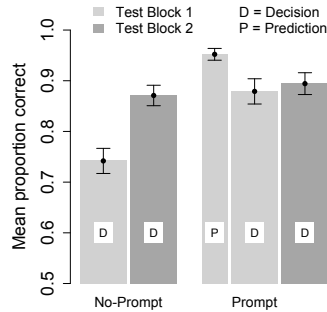


Figure 2.5. Mean proportion of correct predictions and decisions, depicted separately for the No-Prompt and Prompt groups, the first and second test block.

As can be seen in Figure 2.5, the extent to which the proportion of correct decisions increased from the first test block to the second depended on whether participants were prompted to predict the other player's decision,  $\chi^2(1) = 12.2$ ,  $p < 0.001$ . There was a significant main effect of *test block*,  $\chi^2(1) = 78.8$ ,  $p < 0.001$ , but it was mainly driven by the No-Prompt group,  $\beta = 1.072$  (SE = 0.12),  $z = 9.96$ ,  $p < 0.001$ . The probability of making a correct decision increased just slightly in the Prompt group,  $\beta = .351$  (SE = .156),  $z = 2.256$ ,  $p = .024$ .

There was also a main effect of prompting participants for predictions,  $\chi^2(1) = 13$ ,  $p < .001$ . However it was mostly present in the first test block, as the difference between the prompting conditions in the second test block was small and just significant,  $\beta = .639$  (SE = .314),  $z = 2.032$ ,  $p = 0.04$ . Thus, during Test Block 2, participants that were not prompted for predictions in Test Block 1 performed almost as well as participants that *were* prompted. Nevertheless, prompting did have a stronger positive effect in the short term during Test Block 1,  $\beta = 1.36$  (SE = .304),  $z = 4.467$ ,  $p < 0.001$ .

Figure 2.5 also shows that participants performed close to ceiling with respect to their predictions, which also require second-order theory of mind. However, their proportion of correct decisions was significantly lower,  $t(46) = -3.1827$ ,  $p < 0.01$ . Thus, a correct prediction did not always yield a correct decision, which implies that it is not trivial to incorporate a second-order inference in the decision making process. This is remarkable, as the application of theory of mind is required when making a prediction, and not anymore when a decision has to be made afterwards. Nevertheless, this finding supports our hypothesis that integration of mental states, when making decisions, is the crux of suboptimal performance.

## General conclusions

In this study we investigated decision making in the context of sequential games. Crucially, participants were asked to make decisions that required them to infer second-order, or recursive, mental states. Previous studies have found suboptimal performance in such games and concluded that incorrect decisions were caused by insufficient ability to apply second-order ToM (Flobbe et al., 2008; Hedden & Zhang, 2002; McKelvey & Palfrey, 1992; Raijmakers

et al., 2013; Zhang et al., 2012). Other studies, however, have found that participants have few difficulties applying second-order ToM if they are not required to make decisions (Flobbe et al., 2008). We therefore argue that the crux of suboptimal performance in sequential games is the integration of second-order ToM in the decision making process. Our results confirm this idea and show that decision making in sequential games can be improved by scaffolding the integration of second-order ToM.

Our hypothesis is best supported by the finding that correct predictions do not always result in correct decisions. Interestingly, a prediction of what the other player would do requires the application of second-order ToM, as participants have to reason about the other player, who in turn reasons about their last possible decision. Participants performed almost at ceiling when asked to predict the other player's decision. In other words, the participants had sufficient ability to apply second-order ToM. Nevertheless, they performed significantly worse when asked to make a decision based on their predictions. This finding implies that participants found it difficult to integrate second-order ToM in the decision making process, not to apply second-order ToM.

The other findings, that decision making improves by providing participants with Stepwise training, prompts for predictions, and the Marble Drop task representation, support our hypothesis as well. As mentioned previously, these conditions should scaffold the integration of recursive mental states. In Stepwise training, participants practice how to make decisions based on mental states of increasing complexity. When prompted to predict the other player's decision, participants are asked to take the other player's perspective, which is an essential step that should precede and be the basis of their own decision making. The Marble Drop task representation provides clear visual cues as to how the games are structured. In fact, a previous eye tracking study (Meijering, Van Rijn, Taatgen, & Verbrugge, 2012) has shown that participants use these cues while reasoning. For example, the participants in Meijering et al.'s (2012) study kept fixating the trapdoors during the entire experiment, even though the trapdoors remained a constant factor that did not change across the training and test phases. Meijering et al. argued that the trapdoors were used as placeholders for mental states, thereby providing scaffolding for the decision making process. In other words, the Marble Drop task representation, together with the other factors, facilitated the integration of second-order ToM in the decision making process.

The effects were most prominent in the first test block. Participants who were not assigned to stepwise training, who were not prompted to take the other player's perspective, and who played abstract Matrix Games during the entire experiment, did learn in the long run how to incorporate complex mental states into their decisions. In the second-test block the differences in performance between the conditions failed to reach significance. Thus, it seems that participants benefitted most from supporting structure early on in the task when they had not yet developed a strategy for which their cognitive resources were sufficient, as we will explain below.

Before we continue with a cognitive explanation, we would first like to address the concern that people were playing against a computer-simulated player instead of a real player. Based on Hedden and Zhang's (2002; 2012) findings we did not construct a cover story to make participants believe they were playing against a real player. In Hedden and Zhang's study it did not matter whether participants thought to be playing against a real or a computer-simulated player. Furthermore, the level of intelligence participants attributed to the real player did not

affect their default mental model of the other player. In other words, the participants were aware of the nature of the other player, but that awareness did not affect what beliefs, goals, and intentions they attributed to that player. Moreover, other studies have shown, too, that people do attribute (human) mental states to computer-simulated players (Flobbe et al., 2008; Meijering et al., 2012; Meijering, Van Rijn, Taatgen, & Verbrugge, 2013).

Given that sequential games evoke the application of (recursive) ToM, we argue that our findings are generalizable to other ToM settings, including everyday social interactions. Still, generalizability depends on the extent to which supportive measures are effective in reducing demands on cognitive resources. A slightly different notion of generalizability is the extent to which experience on our task will improve application of ToM in other domains. Here, we are more conservative because Flobbe et al.'s (2008) study has shown that performance on one particular ToM task does not necessarily correlate with performance on another ToM task. As Flobbe et al. have shown that application of ToM may not be a unitary skill, it is conceivable that training in one ToM task does not necessarily generalize to other ToM tasks (also see Thoermer, Sodian, Vuori, Perst, & Kristen, 2012).

A cognitive explanation for the facilitative effects of training, prompting, and task representation may be in terms of reducing demands on executive functions. Examples of executive functions are planning, resistance to interference, set-shifting, and working memory, which all help to combine alternate perspectives and find the best possible decision (Apperly & Butterfill, 2009; Bull, Phillips, & Conway, 2008; Dumontheil, Apperly, & Blakemore, 2010). Our experimental manipulations may have reduced demands on executive processes in three ways. First, in the Stepwise training condition, participants received naturally delineated chunks of instruction and training of each successive ToM-order. A clear outline of the task at hand helped planning and reduced demands on working memory. Second, by prompting for predictions, we may have structured participants' reasoning by providing them with an efficient method of 'solving' games. A structured method may not only help planning and thereby reduce demands on working memory; it may also help set-shifting, that is, switching between goals, beliefs, and intentions. Third, the Marble Drop task representation visually cued possible actions and consequences, and provided placeholders for mental states. These cues and placeholders allowed participants to preserve working memory capacity and help them planning their own actions.

Support for these explanations comes from studies that have shown that application of ToM is an effortful process (e.g., Apperly et al., 2010; Lin, Keysar, & Epley, 2010). In Lin et al.'s (2010) study, for example, working memory capacity was positively correlated with efficiency in applying ToM. Participants with low working memory capacity were less efficient in applying ToM than participants with high working memory capacity. Moreover, a second experiment in Lin et al.'s study demonstrated that participants' ability to apply ToM was significantly reduced by a secondary task. These findings show that even though people may have the capacity to apply ToM, they may fail in correctly using ToM, which was already noted by Keysar et al. (2003). Lin et al.'s findings imply that, the other way around, efficiency may improve if demands on working memory, and other executive functions, are reduced. Our results seem to corroborate this notion.

A similar explanation has been proposed for suboptimal behavior in the multitasking setting. Borst et al. (2010a; 2010b) have shown that multitasking in itself is not difficult, but that instead working memory constraints cause the often-claimed difficulties associated with

multitasking. If both tasks cause high working memory load, performance breaks down. However, performance does not deteriorate if at most one of the tasks causes high working memory load. By reducing working memory load in our task, participants had more freely available cognitive resources to incorporate mental states into their decision making process (see also Van Rij, Van Rij, & Hendriks, 2010; 2013).

To conclude, decision making can be a complex task, depending on the variables involved. In this study, decisions were dependent on mental states. Participants had to apply second-order ToM and use the outcome as a basis for their decision making. This proved especially difficult: Correct predictions, which required inference of second-order mental states, did not always result in correct decisions. We argue that participants failed to make optimal decisions because they had difficulties integrating the complex mental states. Our experimental manipulations specifically and successfully targeted this integration process, showing that decision making on the basis of mental states appears to be flexible in the sense that it is susceptible to improvement. Generalizing our findings to everyday life, one could argue that the complex decision making that we engage in in many social settings is a skill that we can develop. We can improve it by taking measures that help us incorporate the mental states of others.

## Appendix A

The table below lists the payoff structures that were used in the experiment. The prediction and decision columns list whether the correct responses were either to stop (i.e., 0) or to move (i.e., 1) at the corresponding decision points.

Payoffs Player 1				Payoffs Player 2				Prediction	Decision
A	B	C	D	A	B	C	D		
Zero-order payoff structures									
1	3			3	2				1
3	1			2	1				0
4	2			2	4				0
3	4			2	4				1
First-order payoff structures									
2	1	3		1	2	3		1	1
2	1	3		3	2	1		0	0
2	3	1		1	2	3		1	0
2	3	1		3	2	1		0	1
3	2	4		2	3	4		1	1
3	2	4		4	3	2		0	0
3	4	2		2	3	4		1	0
3	4	2		4	3	2		0	1
Second-order payoff structures									
3	1	2	4	2	3	4	1	0	0
3	1	2	4	4	2	3	1	0	0
3	2	1	4	1	3	4	2	0	0
3	2	1	4	4	2	3	1	0	0
3	4	1	2	1	3	2	4	1	0
3	4	1	2	2	3	1	4	1	0
3	4	1	2	3	2	1	4	1	0
3	4	1	2	4	2	1	3	1	0
3	4	1	2	1	3	4	2	0	1
3	4	1	2	2	3	4	1	0	1
3	4	1	2	3	2	4	1	0	1
3	4	1	2	4	2	3	1	0	1
3	1	2	4	2	3	1	4	1	1
3	1	2	4	4	2	1	3	1	1
3	2	1	4	1	3	2	4	1	1
3	2	1	4	4	2	1	3	1	1
Trivial payoff structures									
2	4	3	1	1	2	4	3	1	1
3	4	2	1	1	2	3	4	1	0
2	1	4	3	4	3	1	2	0	0
3	4	1	2	4	3	1	2	0	1
2	3	4	1	1	2	3	4	1	1
3	4	1	2	4	3	2	1	0	1

Payoffs Player 1				Payoffs Player 2				Prediction	Decision
A	B	C	D	A	B	C	D		
2	1	3	4	4	3	1	2	0	0
3	4	2	1	1	2	4	3	1	0
3	4	2	1	4	3	2	1	0	1
2	1	4	3	4	3	2	1	0	0
3	4	1	2	1	2	4	3	1	0
2	4	3	1	1	2	3	4	1	1
3	4	2	1	4	3	1	2	0	1
3	4	1	2	1	2	3	4	1	0
2	3	4	1	1	2	4	3	1	1
2	1	3	4	4	3	2	1	0	0
3	2	4	1	4	3	1	2	0	0
2	4	1	3	4	3	2	1	0	1
2	1	4	3	1	2	3	4	1	1
3	4	1	2	4	1	2	3	1	0
2	1	3	4	1	2	4	3	1	1
2	3	1	4	4	3	2	1	0	1
3	4	2	1	4	1	2	3	1	0
3	1	4	2	4	3	2	1	0	0

## Appendix B

Payoff structures are excluded if Player 1's payoff in A is either a 1 or a 4, because Player 1 would not need to reason about the Player 2's decision. It is obvious that Player 1 should continue the game if his payoff in A is a 1 and stop if his payoff in A is a 4. The game in Figure B.1b is an example of a game in which Player 1 should decide to stop in A. In line with Hedden and Zhang (2002), we focused on so-called 2- and 3-starting games, associated with payoff structures in which Player 1's first payoff was a 2 or a 3, respectively.

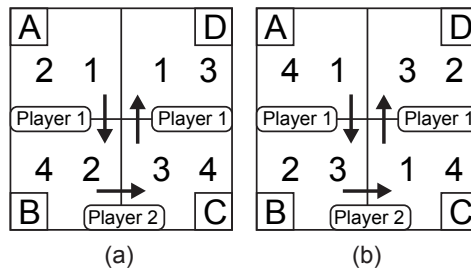


Figure B.1: Two example matrix games. The game in *a* is a so-called trivial game. See text for explanation. The game in *b* does not require any ToM reasoning at all, because Player 1's maximum payoff is already available in A.

We also excluded payoff structures in which Player 2's payoff in B was either a 1 or a 4, because Player 2 would not need to reason about Player 1's last possible decision between C and D. Accordingly, first-order reasoning on the part of Player 1 would suffice.

Of the remaining payoff structures, we excluded the so-called trivial ones in which Player 2's payoff in B was either lower or higher than both his payoffs in C and D. Figure B.1a depicts an example of such a game: Player 2 does not need to reason about Player 1's last possible decision, as his payoffs in C and D are both more preferable than his payoff in B.

The next two exclusion criteria are based on the type of Player 2 that a participant could be reasoning about. In line with Hedden and Zhang, we distinguished between a hypothesized zero-order Player 2 and a hypothesized first-order Player 2. Hedden and Zhang consider these Player 2 types to be either myopic or predictive, respectively. A participant, always assigned to the role of Player 1, might be reasoning about a zero-order Player 2, which would *not* reason about the participant's last possible decision. Hedden and Zhang consider such a Player 2 to be myopic, as it only considers Player 2's own payoffs in B and C. In contrast, a participant might be reasoning about a hypothesized first-order Player 2, which *does* reason about the participant's last possible decision.

As Player 1's decision depends on Player 2's decision, payoff structures that yield the same answer for zero-order and first-order Player 2s cannot inform us on the level of ToM reasoning on the part of Player 1. These payoff structures are considered non-diagnostic, and are therefore excluded from the final set of stimuli. In contrast, Hedden and Zhang included these payoff structures as long as the decision (and thus prediction) of Player 2's move at the second decision point was opposite for imagined zero- and first-order Player 2s. As half of the participants were not prompted for predictions in the first test block and none of them in the second test block, we had to exclude these items.



We selected a final set of stimuli, which we were able to (double-)balance for both the number of *stop* and *continue* decisions of Player 1 and the number of *stop* and *continue* decisions of Player 2. As this was only possible for 3-starting games, we excluded the 2-starting games. This left us with 16 unique payoff structures, all 3-starting games.

Using the same criteria as mentioned above, we selected 4 zero-order and 8 first-order payoff structures for the Stepwise Training condition. These training payoff structures, as well as the 16 second-order payoff structures, are listed in Appendix A. The trivial games used in the Undifferentiated Training condition, which were described earlier, are also listed in Appendix A.

## **Chapter 3**

# **Reasoning about self versus others: Changing perspective is hard**

### **Abstract**

To understand others, we need to infer their mental states, such as beliefs, desires, and intentions. Some developmental studies have suggested that general reasoning ability plays a crucial role in inference of mental states. In contrast, we show that the most important factor for successful inference of mental states is the ability to change perspective. In our experiment, participants either had to make a decision, or predict how another person would make the exact same decision. Crucially, the required steps to solve both problems were the same. Nevertheless, participants made more mistakes and required more time while making predictions instead of decisions. This finding implies that perspective taking, while making predictions, employs computational processes that are unique to the mental aspects of a problem.

This chapter was submitted to a journal where it is currently under review.

We are living in a social world in which many of our daily activities involve interactions with others. These interactions can take many forms: negotiating a higher salary, gossiping with a friendly neighbor, bluffing in a card game, or having a conversation via Facebook or Twitter. Irrespective of form, a social interaction requires the ability to infer another's knowledge, beliefs, desires, and intentions. For example, if we are bluffing in a game of poker we are reasoning about what the other player may think our intentions are if we raise the bet. Reasoning about beliefs and other mental states requires a so-called Theory of Mind (ToM), which develops around the age of 4 (Wellman, Cross, & Watson, 2001; Wimmer & Perner, 1983). That is relatively late considering that infants as young as 7 months already seem to be susceptible to the beliefs of others (Kovács, Téglás, & Endress, 2010; also see O'Neill, 1996; Onishi & Baillargeon, 2005). Furthermore, the process of putting oneself in another's shoes seems to be almost effortless and automatic (Cohen & German, 2010; Kovács et al., 2010; Ramsey, Hansen, Apperly, & Samson, 2013; Samson, Apperly, Braithwaite, Andrews, & Bodley Scott, 2010). Given these findings, it is surprising that it takes so long before children learn to infer the beliefs of others and that adults still frequently fail to use ToM (Apperly et al., 2010; Keysar, Lin, & Barr, 2003; Lin, Keysar, & Epley, 2010).

In recent years, some developmental studies have investigated whether preschoolers' difficulties with ToM are specific to the domain of mental states (Leekam, Perner, Healey, & Sewell, 2006; Perner & Leekam, 2008; Sabbagh, Xu, Carlson, Moses, & Kang, 2006). They set out to investigate whether problems in ToM tasks arise because preschoolers have to put themselves into another's shoes, or because the underlying logical problems in these tasks require cognitive functions that have not yet matured. The preschoolers were presented with two tasks. In one, the false-belief task, the preschoolers were presented with a story in which a person, call her Anne, stored a toy in a box before leaving the room. While Anne was away, a second person, Bob, moved the toy to another location without Anne being aware of it. After Anne returned, the preschoolers were asked where she would look for the toy. To answer correctly they had to reason about Anne's false belief that the toy was still in the box. Performance in this task was compared with performance on the so-called false-sign task, which has the same logical structure but does not involve mental states. A toy is stored in one location, and a sign is pointing at that location. Next, the toy is moved to another location, but the sign is still pointing at the original location. After the story had been told, the preschoolers were asked to indicate where the toy should be according to the sign. In this task, children did not have to reason about mental states (i.e., a false belief), but they did have to refrain from indicating the actual location of the toy, similar to the false-belief task. Because performance in this task correlated with performance in the false-belief task, it seemed that preschoolers' difficulties in understanding false beliefs were not solely confined to mental states (Leekam et al., 2006; Sabbagh et al., 2006). In fact, performance on both tasks correlated with measures of general reasoning ability (Sabbagh et al., 2006), which suggests that inference of mental states suffered from limited capacity of general cognitive functions, not ToM proficiency (also see Bloom & German, 2000).

These findings may seem compelling, but why then do adults, whose general cognitive functions have matured, still frequently fail to apply ToM to interpret the behavior of others? Adults have sufficiently developed general cognitive functions to be able to comprehend many logical problems. Yet, it seems that they reflexively reason about their own beliefs when interpreting the behavior of others (Apperly et al., 2010; Keysar et al., 2003; Lin et al., 2010).

We therefore argue, in contrast to the developmental studies, that people *do* find it difficult to put themselves in another’s shoes. The preschoolers in the developmental studies may have had such difficulties too, but they did not have sufficiently matured general reasoning ability to understand the logical structure of false beliefs and false signs, to begin with. In contrast, adults do have mature general reasoning ability and therefore are the appropriate population to investigate the role of perspective taking during inference of mental states: If adults understand the logical structure of a given ToM task but find it difficult to put themselves in another person’s shoes, unsuccessful application of ToM would reflect a deficiency of a specialized function, not a deficiency of general reasoning ability.

To test specifically whether adults find it difficult to put themselves in another’s shoes, we devised a task in which they either had to make a decision themselves, or had to predict how another would make the very same decision (Appendix A). The task consisted of two-player games in which the participant and a computer-simulated player alternately made decisions. There were two conditions and the only difference between the conditions was the required level of perspective taking; all other task aspects were the same. Importantly, this was a direct and specific test of perspective taking, because participants were not asked to reason about distinct types of representations, which was the case in the false-belief and false-sign tasks. Crucially, the steps to ‘solve’ the games were the same in both conditions, irrespective of the required level of perspective taking. The only difference for a participant was the instruction, or prompt, given at the start of each game. In one condition the prompt was “Decide” what

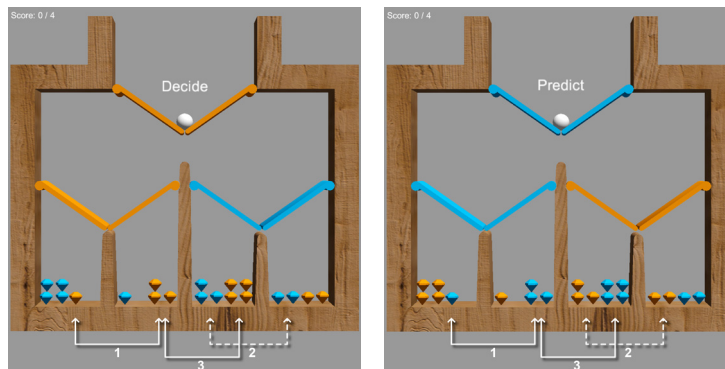


Figure 3.1: Screenshots of isomorphic Decision (left panel) and Prediction (right panel) games. The arrows were added for illustrative purposes and show that the games require the same comparisons to obtain the best possible outcome. In these particular games, the participant is assigned the target color orange and has to obtain as many orange diamonds as possible. The game in the left panel prompts the participant to make a decision: “Decide”. To make a decision, the participant needs to switch perspective once, as his outcome depends on the decision of the other player, who decides at the bottom-right trapdoors and whose goal is to obtain as many blue diamonds as possible. The game in the right panel prompts the participant to predict the other player’s decision: “Predict”. The participant needs to switch perspectives twice: The first time to reason about the other player’s intention at the topmost trapdoors, and a second time to switch back to his own perspective, because the other player’s decision depends on the participants decision at the bottom-right trapdoors.

you would do in this game, and in the other condition it was “Predict” what the other player would do in this game (see Figure 3.1). Again, the steps to arrive at a decision or a prediction were the same, but it mattered whether participants were asked to execute these steps from their own perspective or from the other’s, as we will show.

At the start of the experiment, each player was assigned their own target color (orange or blue) and the goal was to obtain as many target-colored diamonds as possible (Appendix A). In each game a white marble dropped into a contraption (see Figure 3.1) and both players could influence its path by opening the trapdoors depicted in their target color. Each player obtained the target-colored diamonds located in the bin into which the marble dropped. Both players had to reason about one another, because the other player’s decisions affected their outcomes. In games such as the one in the left panel of Figure 3.1, for example, participants were prompted to make a decision at the topmost trapdoors. They could infer that their best possible decision is to open the right-side trapdoor, because the other player’s intention is to open the left-side trapdoor, as his goal is to obtain as many blue diamonds as possible: He would obtain 3 blue diamonds, and the participant would obtain 4 orange diamonds. In these so-called Decision games participants had to switch perspective just once while they were reasoning about the other player. In so-called Prediction games, however, participants had to switch perspective twice (Figure 3.1; right panel): To predict the other player’s decision at the topmost trapdoors, they had to switch from their own perspective to that of the other player. Furthermore, as the topmost decision was based on their own decision at the bottom-right trapdoors they had to switch back again from the other player’s perspective to their own. In other words, participants had to reason about recursive mental states: “*the other player thinks that I intend to open the left-side trapdoor, as he knows that my goal is to obtain the highest possible number of orange diamonds*”. Importantly, the Decision and Prediction games were structurally equivalent: The order of the decisions was the same in both types of games, as was the distribution of diamonds. Thus, participants had to make the same comparisons between the same distributions of diamonds, as can be seen in Figure 3.1, irrespective of the prompt given at the start. The only difference between the games was the perspective from which to make the comparisons.

In this study, we are comparing two hypotheses. Based on the earlier mentioned developmental studies, one would expect the response patterns to be the same in Decision and Prediction games, because these games require the same steps to ‘solve’ them (see Figure 3.1). As a consequence, Decision and Prediction games induce the same demands on general reasoning ability. The accuracy and response times are therefore not expected to differ, which makes this the null-hypothesis. The other hypothesis is that adults will perform differently in Decision and Prediction games, because these types of games differ with respect to the required level of perspective taking. A previous study has suggested that people do switch between perspectives each time they fixate a new set of trapdoors (Meijering, Van Rijn, Taatgen, & Verbrugge, 2012). Based on that study, we expect the accuracy and response times to differ, because the number of switches between perspectives differs between Decision and Prediction games.

As can be seen in Figure 3.2, both the accuracy (i.e., the proportion of correct responses) and the response times indicate that it is more difficult to reason about someone else’s decision-making than making the same decisions oneself (Appendix B). In the Prediction games, in comparison to the Decision games, accuracy is significantly lower,  $\chi^2(1) = 44.77$ ,  $p$

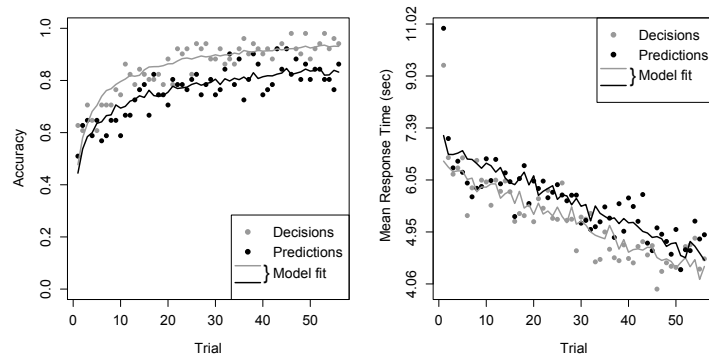


Figure 3.2: The left panel depicts the mean accuracy of decisions (grey) and predictions (black) across participants. The right panel depicts the actual response times on a logarithmic scale, also averaged across participants. The solid lines depict fits of linear mixed-effects regression models.

$< 0.001$ , and the response times are significantly longer,  $\chi^2(1) = 4.97$ ,  $p < 0.05$ . These findings imply that the reasoning processes were not symmetrical in the two types of games. We argue that the asymmetrical response patterns are due to differential demands on perspective taking, because the games are isomorphic otherwise. In prediction games participants had to switch perspectives twice, and as a consequence they made more mistakes and needed more time to produce a response. We accounted for the fact that Decision and Prediction games were presented in an intermixed fashion, which may have caused asymmetric reorientation costs while switching back and forth between Decision and Prediction games. However, this reorientation factor was not significant, nor was any interaction with it. Thus, the demand on perspective taking was the sole factor determining accuracy in Decision and Prediction games.

There is a ‘smart’ strategy to play these games, and that is to pretend to switch between colors, reasoning as if one is the ‘orange’ player in some games (e.g., Decision games) and the ‘blue’ player in the others (e.g., Prediction games). By using this strategy the participants would have reduced the required level of perspective taking of all games to level one. Strikingly, the participants did not use this strategy, as the accuracy and response times differ between Decision and Prediction games. Our findings therefore show that the participants were strongly committed to their own target color and reasoned from their own perspective when prompted to “Decide” and from the other player’s perspective when prompted to “Predict”. In fact, the difference in accuracy and the ratio of associated response times did not become any smaller during the experiment (Appendix B). This is remarkable as participants were presented with 112 games in total, which is ample opportunity to adopt a smart strategy. The fact that the differences did not become smaller is a strong indication that the participants committed to their own target color and engaged in perspective taking, which caused differential demands on mental state reasoning in the Decision and Prediction games.

There are several possible explanations for differential response patterns in Decision and Prediction games. For example, the process of modeling the other player’s mental states may require cognitive functions that are not as well developed as the cognitive functions to model

one's own mental states (Birch & Bloom, 2004). The process of modeling the other player's goals and intentions could therefore be prone to errors, yielding incorrect predictions. A related explanation is that people fall prey to the egocentricity bias and consequently interpret the behavior of the other player according to their own goals and intentions (Birch & Bloom, 2007; Keysar et al., 2003). Overcoming this bias may be difficult and cause high cognitive demands, as one has to inhibit one's own mental states (German & Hehman, 2006; Samson, Apperly, Kathirgamanathan, & Humphreys, 2005). Lastly, a representation of the other player's mental states could intrinsically be more complex (Gopnik & Wellman, 1992), because the mental states have to be labeled as belonging to that player. One's own mental states, in turn, do not have to be labeled. In any case, each of these possible factors is a consequence of perspective taking, as perspective taking alone distinguishes between Decision and Prediction games.

To conclude, the results show that reasoning about someone else's decision-making is more difficult than making the same decisions oneself, even if the conditions are equivalent. The steps to arrive at either a decision or a prediction were the same in this study, but apparently participants engaged in distinct processes. Critically, predictions required more switches between perspectives than did decisions, and as a consequence participants produced fewer optimal responses and required longer response times. This study shows that the bottleneck of mental state reasoning is perspective taking: Reasoning about self is not as difficult as reasoning about others.

## **Appendix A: Methods**

### **Ethical statement**

The study was approved by the ethical committee of the Psychology department of the University of Groningen.

### **Participants**

In this study, 51 first-year psychology students (34 female) participated in exchange for course credit. Their mean age was 21, ranging from 18 to 34 years. None of the participants were excluded from the data analyses. All participants had normal or corrected-to-normal visual acuity. Written informed consent was obtained from all participants.

### **Stimuli**

Each game had a unique distribution of payoffs (i.e., diamonds). Of all possible payoff distributions we excluded those that did not require inference of mental states. There were three exclusion criteria in total. First, a payoff distribution was excluded if the player deciding at the topmost trapdoors had her two highest payoffs on one side and her two lowest payoffs on the other side. In this case, the player would not have to consider the other player's decision. Second, a payoff distribution was excluded if the maximum payoff of the player deciding at the topmost trapdoors was behind her own two sequential trapdoors. In case of such a payoff distribution, she does not need to reason about the other player's goals and intentions. Third, a payoff distribution was excluded if both players had the same number of diamonds in each bin. Such a payoff distribution does not require inference of mental states, as there would not be any conflict between the two players' goals, beliefs, and intentions.

From the remaining 192 payoff distributions, 56 were randomly selected to be included as the set of Decision games. Another 56 items, also randomly selected from the 192 payoff distributions, comprised the set of Prediction games.

### **Procedure**

The participants were seated in front of a 24" monitor on which the games were played. Before the task started, instructions were given on paper. The instructions explained: the goals of the participant and computer-simulated player; who decided where; and what response participants should give in case of the "Decide" and "Predict" prompts. After reading the instructions, participants could ask for clarification if they had any questions. The experimenter answered these questions, but was careful not to give any information on the strategy of the computer-simulated player.

The decision and predictions games were presented in random order, and the participants played these games from start to end, until the marble dropped into one of the bins. The games were fully animated, and at the end of each game there was feedback, which indicated the participant's score in that particular game. For example, "+3" if the marble dropped into a bin



that contained three diamonds in the participant's target color. The total score was depicted in the top left corner of the screen. During the first twelve games, participants were provided with additional feedback after each game, which indicated whether the obtained score was the highest possible score they could have obtained in that particular game.

## Appendix B: Results

The individual responses (decisions and predictions) were analyzed by means of *logistic* linear mixed-effect models, as the data contained at least two sources of random variation: The participants and payoff structures were both sampled from larger populations. We constructed a full factorial model that comprised fixed effects of ToM *order* (level 1 / level 2), *switching* between ToM-orders (switch / no-switch), and the covariate *trial*. Trial was log-transformed to account for the non-linear increase in the proportion of correct decisions and predictions (see Figure B.1). The values of the *Akaike information criterion* (AIC) and the *Bayesian information criterion* (BIC) of the full factorial model were compared with those of simplified models from which interaction and main effects were removed. The AIC and BIC values indicate whether a better fit of a more complex model is justifiable given its extra parameters.

### Accuracy

The model with the most favorable (i.e., smallest) AIC and BIC values contained main effects of *order* and *log-trial*, and an interaction between the two. The absence of main effects and interaction effects on accuracy of switching between Decision and Prediction games, suggests that reorientation costs, if any, were negligible and did not differ between switching to Decision and switching to Prediction games.

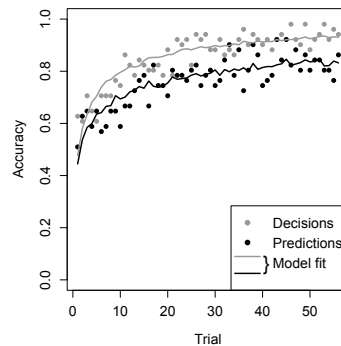


Figure B.1. Mean accuracy of decisions (orange) and predictions (blue), averaged across participants. The solid lines represent the fits of the linear mixed-effects models.

Figure B.1 shows that the accuracy of predictions is lower than the accuracy of decisions. This difference is significant,  $\chi^2(1) = 44.77$ ,  $p < .001$ . At first, the probability of making correct decisions and predictions does not significantly differ,  $\beta_{\text{prediction-decision}} = -.09$  ( $SE = .25$ ),  $z = -.37$ , ns. However, the difference becomes larger, as can be seen in Figure B.1. At the end of the experiment, the probability of making a correct prediction is significantly lower than the probability of making a correct decision,  $\beta_{\text{prediction-decision}} = -1.06$  ( $SE = .15$ ),  $z = -7.10$ ,  $p < 0.001$ .

Figure B.1 also shows that the accuracy of decisions and predictions increase with each game (i.e., trial) played, which is a significant main effect,  $\chi^2(1) = 289.65$ ,  $p < 0.001$ . This trend shows that the participants became better as they played more games. However, the trend was

smaller in prediction games than in decision games,  $\beta_{\log\text{-trial}} = 0.58$  (SE = 0.05),  $z = 10.78$ , and  $p < 0.001$ , and  $\beta_{\log\text{-trial}} = 0.82$  (SE = 0.06),  $z = 13.75$ , and  $p < 0.001$ , respectively. This interaction is significant,  $\chi^2(1) = 9.11$ ,  $p < 0.005$ . Thus, performance was susceptible to improvement, but more so in Decision games than in Prediction games.

## Response times

Participants' response (decision / prediction) times were also analyzed by means of linear mixed-effects models. The reaction times were first log-transformed to reduce skew in the distribution of response times. Performing the procedure of model comparison described above, we found a best fitting model that contained main effects, only, of ToM *order*, *switching*, and *trial*.

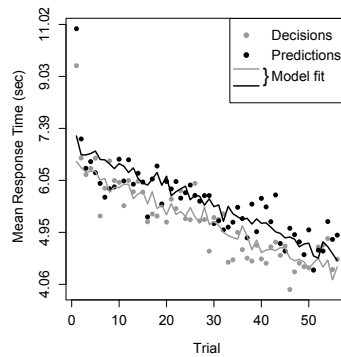


Figure B.2: The response times were averaged across participants; plotted separately for Decision games (orange) and Prediction games (blue) on a logarithmic scale.

The participants did not only make more mistakes in Predictions games than in Decision games (Figure B.1), they also required more time to predict the other player's decision than to make the same decision themselves (Figure B.2). This main effect is significant,  $\chi^2(1) = 4.97$ ,  $p < .05$ .

Figure B.2 shows that the log-transformed response times (log-RTs) decreased linearly during the experiment, in both the Decision and Prediction games. This effect of *trial* on the log-RTs is significant,  $\chi^2(1) = 354.01$ ,  $p < .001$ . The lack of an interaction between type of game (Decision / Prediction) and trial implies that the ratio of the actual decision and prediction times did not change over trial.

Interestingly, whereas switching between Decision and Prediction games did not have an effect on the accuracy of the participants' responses, switching did cause a significant time cost,  $\chi^2(1) = 49.13$ ,  $p < .001$ . This finding implies that switching back and forth between one's own and the other player's perspective did have an associated time cost, in addition to the time cost associated with playing Prediction games instead of Decision games. Importantly, there were no interaction effects that included this switching factor, which means that the time cost of switching between Decision and Prediction games, and vice versa, was symmetrical. Therefore, switching between Decision and Prediction games could not have been a confounding factor

in explaining, for example, shorter RTs in Decision games in comparison to Prediction games. That difference is solely to be attributed to the required order of theory of mind.



## Chapter 4

# Reasoning about diamonds, physics, and mental states: The cognitive costs of theory of mind

### Abstract

Theory of mind (ToM) is required when reasoning about mental states such as knowledge, beliefs, desires, and intentions. Many complex reasoning tasks require domain-general cognitive resources such as planning, resistance to interference, and working memory. In this paper we present a study of the additional cognitive costs of reasoning about mental states. We presented participants with sequential games in which they have to reason about another player. In the so-called *player* condition, the other player is reasoning about the participant, whereas in the so-called *balance* condition, the other player is reasoning about a balance scale. Both types of games require the same comparisons, but only differ in the required depth of ToM reasoning. Games in the player condition require one additional switch between perspectives. The results show that participants make different types of mistakes in the player condition as compared to the balance condition. This finding implies a different reasoning process when reasoning about mental states. The results also show faster decreasing reaction times in the balance condition than in the player condition. Based on these findings, we argue that reasoning about mental states requires unique cognitive resources.

Part of this chapter was previously published in the *proceedings of the 35th Annual Conference of the Cognitive Science Society* (2013).

## Introduction

In many social interactions we reason about one another. If, for example, our decisions or outcomes depend on someone else's actions, we try to predict what the other will do. Predicting the other's actions requires an understanding of how behaviors are caused by mental states such as beliefs, desires, goals, et cetera. Such an understanding is often referred to as *theory of mind* (Baron-Cohen, Leslie, & Frith, 1985; Premack & Woodruff, 1978; Wellman, Cross, & Watson, 2001; Wimmer & Perner, 1983).

A theory of mind, or ToM, is starting to develop around the age of three to four years (Wellman et al., 2001; Wimmer & Perner, 1983). However, younger infants already are susceptible to others' mental states (Kovács, Téglás, & Endress, 2010; O'Neill, 1996; Onishi & Baillargeon, 2005). One possible explanation is that they are able to read others' behavior, but cannot yet explicitly reason about the underlying mental states. Only after many interactions, reading many distinct behaviors, do children start to develop a theory of how behaviors generally correspond with beliefs, desires, intentions, et cetera.

So far, we have introduced ToM as being a *theory* (Gopnik & Wellman, 1992; Wellman et al., 2001). However, we do not want to exclude another definition of ToM that considers it to be an *ability* or skill to reason about mental states of oneself and others (Apperly, 2011; Leslie, Friedman, & German, 2004; Van Rij, Van Rij, & Hendriks, 2010; Wimmer & Perner, 1983). In fact, a theory alone would not suffice when reasoning about others' mental states. Such reasoning is an entire process of generating many possible mental state interpretations (Baker, Saxe, & Tenenbaum, 2009), and ToM reasoning might be qualitatively different from other kinds of reasoning.

Some studies have shown similar but uncorrelated developmental trends in ToM tasks and non-mental tasks that require similar representations (Arslan, Hohenberger, & Verbrugge, 2012; De Villiers, 2007; De Villiers & Pyers, 2002; Flobbe, Verbrugge, Hendriks, & Krämer, 2008; Hale & Tager-Flusberg, 2003). For example, a relative clause in the sentence "The goat that pushes the cat" requires a similar representation as the complement clause in "Alice knows that Bob is writing", but only the complement clause requires a mental state representation. As children become older, they get better at understanding both types of sentences. However, their performance does not correlate when the factor age is controlled for. These findings show that ToM tasks might consume unique cognitive resources. It is important to note, however, that these tasks might have differed with respect to other factors, besides the aspect of mental representations.

Some studies show similar performance in ToM tasks, on the one hand, and equivalent but non-mental control tasks, on the other. In the false-belief or Sally-Anne task, for example, children have to attribute a false belief about an object's current location to Sally (Baron-Cohen et al., 1985; Wellman et al., 2001; Wimmer & Perner, 1983). Sally stores an object at location A, but the object is moved from location A to location B while Sally is away. Therefore, Sally still thinks that the object is at location A. To pass this task, children should acknowledge that Sally falsely believes that the object is still at location A. The false-sign task is a similar but non-mental counterpart of the false-belief task. An object is first stored at location A, indicated by an arrow. Next, the object is moved from location A to location B, but the arrow still points at location A. The false sign in this task is the arrow pointing at location A, which is similar to Sally's false belief. Children's accuracy in both tasks is similar, and their

performance correlates, even after correcting for age (Leekam, Perner, Healey, & Sewell, 2006; Perner & Leekam, 2008; Sabbagh, Xu, Carlson, Moses, & Kang, 2006). This finding implies that mental state reasoning might not qualitatively differ from other kinds of reasoning.

Similar accuracy of responses in ToM tasks and their non-mental counterparts, however, does not necessarily imply a similar reasoning process. Moreover, differences might manifest themselves elsewhere, for example, in the reaction times. If, for example, both tasks require overlapping cognitive functions but ToM tasks require additional cognitive processing, the response patterns might not differ as much as the associated response times. Moreover, differences in accuracy might not manifest themselves until the tasks become more complex and exhaust cognitive resources.

Given these mixed findings, the question remains whether reasoning about mental states requires additional cognitive resources. Complex reasoning tasks consume cognitive resources, because oftentimes they require integration of information in the overall reasoning process. Integration of information and reasoning both require executive functions such as

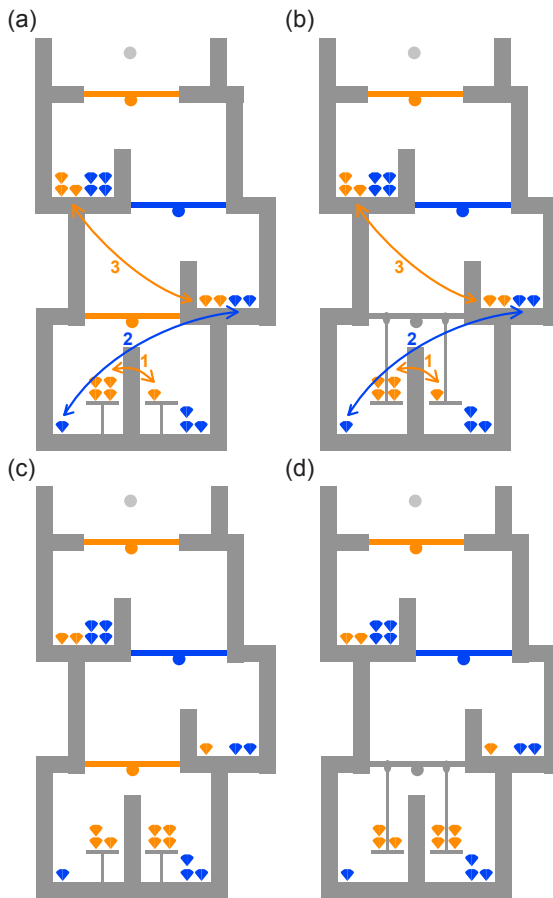


Figure 4.1. Examples of two-player Marble Drop games. A white marble is about to drop, and its path can be manipulated by turning the orange and blue trapdoors. In these example games, participants have to obtain as many orange diamonds as possible and they control the orange trapdoors. The other player has to obtain as many blue diamonds as possible and controls the blue trapdoor. The arrows in *a* and *b* indicate what comparisons should be made to make an optimal decision at the topmost trapdoor. In game *a*, the optimal decision for a participant is to let the white marble drop into the topmost bin, thereby obtaining 3 orange marbles. The 4 orange diamonds in the bottom-left bin are not obtainable for the participant, as the other (blue) player's optimal decision is to let the white marble drop into the middle bin: The other player knows that the optimal (orange) decision at the bottom trapdoors is to go left, yielding a suboptimal outcome of 1 blue diamond for Player 2. Games *a* and *c* are second-order games, because participants (as Player 1) have to reason about the other player (i.e., Player 2) who in turn has to reason about Player 1. The games in *b* and *d* are first-order counterparts of the games in *a* and *c*, respectively. They require the same comparisons, as the outcome of the balance scale

is congruent with Player 1's last correct / rational decision: Both depend only on Player 1's diamonds in the bottom two bins. However, the games with the balance require one fewer switch between Player 1 and Player 2 perspectives.



planning, set shifting, resistance to interference, and working memory. It is not yet obvious why these executive functions alone would not suffice to reason about mental states.

In this study we investigate whether reasoning about mental states consumes unique cognitive resources. Participants are presented so-called Marble Drop games (Figure 4.1) in which they have to reason about another player. Marble Drop games have a recursive structure because the best possible, or optimal, decision at the first trapdoor depends on the other player's decision at the second trapdoor, which in turn depends on the outcome at the third trapdoor (Meijering, Van Rijn, Taatgen, & Verbrugge, 2011). The crucial factor in this experiment is whether the outcome at the third trapdoor is determined by Player 1's decision (*player* condition) or by the physics of a balance scale (*balance* condition). Both conditions require the same comparisons, but games in the player condition require one additional switch between player perspectives: Player 1 has to reason about what Player 2 thinks that Player 1 will do at the final trapdoor. If reasoning about mental states requires additional cognitive resources, games in the player condition would be more difficult than games in the balance condition.

## Method

Participants are always assigned to the role of Player 1, and in both conditions they need to take the perspective of Player 2 to predict the outcome at the second trapdoor. This perspective taking requires ToM. As explained previously, the decision at the second trapdoor depends on the outcome at the third trapdoor. If the participants (i.e., Player 1) control that trapdoor, they need to switch perspective again. They need to re-take their own perspective from within Player 2's perspective. This requires second-order ToM. In the balance scale condition, participants do not have to switch perspective again, and thus need first-order ToM at most. They still need to make the same comparisons, as the outcome of the balance scale depends on Player 1's payoffs and this outcome is congruent with Player 1's goal to maximize his or her payoffs.

If ToM requires unique cognitive resources, we expect that participants respond faster in the balance condition than in the player condition, because the balance condition requires one switch less between Player 1 and Player 2 perspectives than the player condition. We also expect better performance in the balance condition, because Marble Drop games in which Player 1 controls the third trapdoor might appear to be less deterministic. The hypothesis, here, is that it is easier to attribute knowledge of physics to Player 2 than to attribute to Player 2 epistemic reasoning about Player 1, as epistemic reasoning involves testing of multiple possible Player 2 perspectives.

## Participants

Forty-two first-year Psychology students (30 female) participated in exchange for course credit. The average age was 21 years, ranging from 18 to 25. Each participant reported normal or corrected-to-normal visual acuity.

## Stimuli

Of all possible payoff structures, only those that are diagnostic of second-order ToM reasoning were included in the experiment. A game is diagnostic of second-order ToM reasoning if it requires a participant to reason about each decision point to arrive at the optimal decision. An example of a non-diagnostic payoff structure is one in which Player 1's first payoff, in the topmost bin, is the maximum payoff in that game. In that case, Player 1 would not need to reason about the second and third decision points. The payoff structures are listed in a table, which can be found at [http://www.ai.rug.nl/~meijering/marble\\_drop.html](http://www.ai.rug.nl/~meijering/marble_drop.html).

## Design

The experimental design consists of two between-subjects conditions: balance condition versus player condition. In the player condition, participants are presented with the original second-order ToM games (Meijering, Van Rijn, Taatgen, & Verbrugge, 2012). In the balance condition, participants play the games with the same payoff structures, but the third decision point is replaced by a balance scale. Importantly, the games in both conditions are equivalent, as they require the same comparisons between payoffs. In each game, the outcome of the balance is the same as the last correct / rational decision of Player 1, because both only depend on the number of Player 1 diamonds in the bottom two bins (see Figure 4.1).

## Procedure

After giving informed consent, participants were seated in front of a 24-inch iMac. They were randomly assigned to the balance scale condition or the player condition. The participants were instructed that their goal was to obtain as many diamonds as possible of their target color, either blue or orange, which was counterbalanced between participants. They were also instructed that Player 2's (i.e., the computer's) goal was to obtain as many marbles as possible of the other color.

The experimental procedure is the same in both the player and the balance conditions. Participants are presented 62 unique games. At the start of each game, participants have to decide whether to stop the game, by letting the white marble drop into the top bin, or to continue the game, by letting the white marble drop onto Player 2's trapdoor. The game stops if Player 2 decides to let the white marble drop into the middle bin. If Player 2 decides to let the white marble drop onto the third trapdoor, participants in the player condition have to decide whether to stop the game in the bottom-left or bottom-right bin. In the balance condition, the physics of the balance scale determine whether the marble drops into the bottom-left bin or the bottom-right bin. Importantly, the balance scale is set in motion as soon as the white marble drops onto it. Otherwise, Player 2 would not have to reason about the balance scale. Each game is fully animated. See Figure 4.1 for some example games.

After each game, participants receive feedback that mentions Player 1's outcome. If, for example, the marble drops into a bin that contains two diamonds for Player 1, the feedback mentions: "You get 2".

To familiarize participants with the rules of Marble Drop games, participants are presented additional feedback during the first 12 games. Feedback explicitly mentions whether the

outcome is the highest attainable Player 1 payoff. In case a participant obtains 3 diamonds and could not have obtained more, feedback is: “*Correct. You get 3. The highest possible payoff!*”. In case a participant obtains 3 diamonds, but could have obtained 4, feedback is: “*Incorrect. You get 3. You could have obtained 4*”.

## Results and discussion

The data consist of 62 unique Marble Drop games (i.e., payoff structures) for each participant. In the statistical analyses, the games are blocked to accommodate non-linear and differential learning rates: The first 12 ‘training’ games comprise the first block, and the remaining 50 games are split into 5 subsequent blocks of 10 games each. The graphs show means and standard errors, which are represented by error bars.

The data are analyzed by means of linear mixed-effects models (Baayen, 2008; Baayen, Davidson, & Bates, 2008; Gelman & Hill, 2007) to accommodate random sources of variation due to sampling of participants and items (i.e., payoff structures). Specifically, each model allows for by-participant and by-item adjustments of the intercept. For each analysis that we report below, we first constructed a full factorial model with all main and interaction effects. Based on likelihood ratio comparisons, we removed main and interaction effects for as long as the corresponding parameters were not justified. If a comparison preferred a simplified model, we report the log-likelihood statistics. The correctness of responses is analyzed by means of *logistic* linear mixed-effects models, as correctness of responses is a binary variable (incorrect vs. correct).

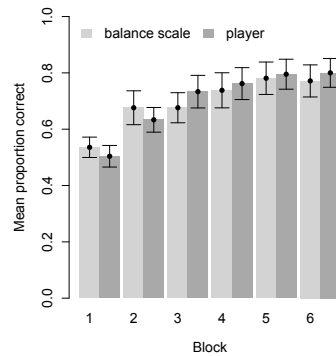


Figure 4.2: Mean proportion of correct responses per block; depicted separately for participants in the balance condition (light gray) and the player condition (dark gray).

### Mean proportion correct

The proportion of correct responses in each block is averaged across participants and depicted in Figure 4.2. The figure does not show great differences between performance in the balance scale and player conditions.

A full-factorial model with main effects and an interaction effect of *Condition* and *Block* did not fit the data better than an additive model,  $\chi^2(5) = 6.08$ , ns. The parameters of the additive model are discussed below.

There is a significant effect of *Block*,  $\beta = 1.37$ ,  $z = 10.89$ ,  $p < .001$ . As can be seen in Figure 4.2, performance increases over the course of playing many Marble Drop games.

There is no effect of *Condition*, as can be seen in Figure 4.2. In contrast to our hypothesis, the probability of making a correct decision does not differ between the balance scale and player conditions. An analysis of the types of errors (next section: Types of Errors), however, shows differential errors between the balance scale and player conditions.

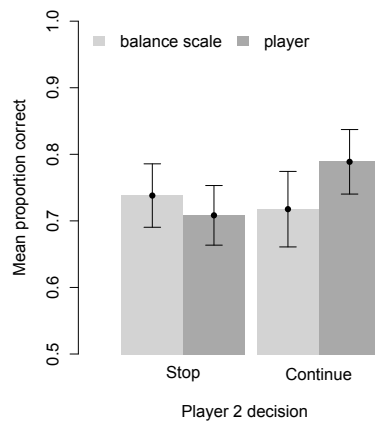


Figure 4.3: Mean proportion of correct responses across participants, depicted separately for the balance scale and player conditions, and Player 2's decision.

## Types of errors

The errors that participants made were categorized according to game type, as an overall analysis might not be sensitive enough to differentiate between the balance scale and player conditions. Two types of games were distinguished on the basis of Player 2's (programmed) decision, which is either stop the game or continue.

There is no main effect of *Player 2 decision*,  $\beta = -.08$ ,  $z = -.575$ , ns, which means that the difficulty of a game does not depend on Player 2's decision. This finding implies that there is no reason to believe that there are particular subsets of hard(er) payoff structures among the selected payoff structures.

There is a significant interaction effect between the factors *Condition* and *Player 2 response* (see Figure 4.3),  $\beta = .65$ ,  $z = 3.349$ , and  $p < 0.001$ . In the balance scale condition, the probability of making a correct decision does not differ between games in which Player 2's decision is to stop, on the one hand, and games in which Player 2's decision is to continue, on the other hand. In the player condition, in contrast, there is a difference. One possible explanation is that participants in the player condition expect Player 2 to continue in most games, and this expectation pays off in games in which Player 2 actually decides to continue. In each game, Player 2 has a greater payoff in one of the last two end states than in the earlier end state, and

participants might assign too great a probability to Player 2 going for that payoff. Participants in the balance condition, in contrast, might estimate those probabilities more accurately (i.e., *lower*), because games with a balance scale can be considered more deterministic.

## Reaction times

There are differences in the types of errors between participants in the balance scale and player conditions, but what about the reaction time data? RTs are analyzed to find out whether a switch between perspectives comes with a time-cost. The RTs are log-transformed as reaction

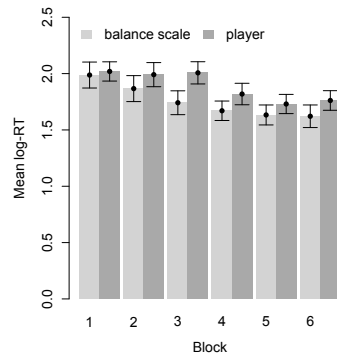


Figure 4.4: Average log-RT across participants plotted against block, separately for the balance scale and player conditions.

times are skewed to the right. Figure 4.4 shows the average log-RT across participants.

Figure 4.4 shows differential learning rates between participants in the balance scale and player conditions, especially in the first half of the experiment, in blocks 1 to 3. In the second half, blocks 4 to 6, the learning rates do not seem to differ that much. To specifically accommodate for differential learning rates, the factor *Block* was re-parameterized as a new factor *Half*, with levels 1 and 2, and a new factor *Block* with levels 1, 2, and 3 within each level of *Half*. The results of the full factorial LME with main and interaction effects of *Condition*, *Half*, and *Block* are discussed below.

The main effects of *Half* and *Block* (with linear contrast) are significant,  $\beta = -.22$ ,  $t = -7.82$ ,  $p < .001$ , and  $\beta = -.18$ ,  $t = -5.37$ ,  $p < .001$ , respectively. From the first to the second half of the experiment, and within each half, the RTs decrease linearly. The interaction between *Half* and *Block* is also significant,  $\beta = .15$ ,  $t = 3.19$ ,  $p = .0015$ . The decrease in RTs is stronger in the first half of the experiment than in the second half.

The interaction between *Condition* and *Block* is significant,  $\beta = .17$ ,  $t = -3.43$ ,  $p < .001$ . The decrease in RTs in the first half of the experiment is less strong in player condition than in balance scale condition. This finding is partly congruent with the hypothesis that RTs are shortest in the balance scale condition because it requires fewer switches between perspectives than the player condition. There is, however, no main effect of *Condition*,  $\beta = .14$ ,  $t = 1.2$ , ns. Thus, on average, the RTs do not differ between the balance scale condition and the player condition. However, participants in the balance scale condition *do* become faster towards the

end of the first half of the experiment, whereas participants in the player condition do not become faster. A possible explanation is that participants in the balance scale condition are quicker over the course of playing multiple games to attribute an understanding of gravity to Player 2. In contrast, participants in the player condition need to play more games and test multiple Player 2 perspectives.

The interaction between *Condition*, *Half*, and *Block* is also significant,  $\beta = -.14$ ,  $t = -2.83$ ,  $p < .005$ . As can be seen in Figure 4.4, the differential learning rates in the first half of the experiment disappear in the second half of the experiment, where the RT trends do not differ that much between the balance scale condition and the player condition.

In sum, there is an interaction effect of *Condition* and *Block* on the RTs, and this effect is mainly present in the first half of the experiment. There, the RTs decrease more in the balance scale condition than in the player condition. This interaction effect, between *Condition* and *Block*, seems to disappear in the second half of the experiment. A possible explanation for the latter finding is that, initially, participants in the balance scale condition settle more quickly on the correct Player 2 perspective than participants in the player condition, who test multiple Player 2 perspectives across multiple games.

## General conclusions

In this study we investigated whether ToM requires additional cognitive resources. We presented two types of games that required the same comparisons but differed with respect to the required depth of ToM reasoning: Games in the *player* condition required second-order ToM, as participants had to reason about a Player 2 that, in turn, reasoned about them; Games in the *balance scale* condition required first-order ToM, as participants had to reason about a Player 2 that reasoned about a balance scale. Our results show different errors between these conditions, which implies that the reasoning was not the same in the balance scale and player conditions. Moreover, the reaction time trends differed. The learning rate was faster for participants in the balance scale condition than for participants in the player condition. A faster learning rate in the balance condition is congruent with our hypothesis that it is easier to play against a Player 2 that reasons about gravity than playing against a Player 2 that reasons about mental states.

We assumed that games with a balance scale are easier to play because they appear to be more deterministic than games in which Player 1 has the last decision. This assumption is congruent with the RT data: Longer RTs in the player condition could be the cause of participants' testing of multiple possible Player 2 perspectives. Games in the balance scale condition, in contrast, require testing of fewer possible Player 2 perspectives, yielding shorter decision times.

Besides a faster learning rate in the balance condition, we expected a greater proportion of correct decisions. However, the probability of making a correct decision does not differ between the balance scale (i.e., first-order ToM) condition and the player (i.e., second-order ToM) condition. One possible explanation is that knowledge about gravity is not automatically attributed to Player 2. We expected that participants in the balance condition would automatically 'see' how Player 2's decision depends on the outcome of the balance, as young children have already mastered many balance scale configurations (Jansen & Van

der Maas, 2002; Van Rijn, Van Someren, & Van der Maas, 2003). However, attributing an understanding of gravity to Player 2 might be less of an automatic process than reasoning about gravity oneself.

Based on our findings, we conclude that participants do need theory of mind in Marble Drop games. Sequential games such as Marble Drop can be critiqued for not requiring ToM: If Player 2's strategy is known, the optimal (Player 1) decision can be determined without reasoning about Player 2's reasoning about Player 1's last possible decision. Applying backward induction, an algorithm based on sequential payoff comparisons, would yield the optimal decision. However, Meijering et al.'s (2012) eye tracking study (see also Chapter 5 in this dissertation) shows that participants use more complicated and diverse reasoning strategies, not only backward induction. Moreover, backward induction would not be able to account for different types of mistakes and differential reaction times in the two conditions, as backward induction always works the same, irrespective of condition. Therefore, our findings provide support for the idea that sequential games are not just a decision-making problem but also evoke reasoning about mental states and thus require ToM.

In fact, it seems that sequential games are a particularly good paradigm to test reasoning about mental states, as they require *active* application of ToM. If Player 2's strategy is not yet known, participants need to actively find the correct Player 2 perspective. In any given game, multiple Player 2 perspectives might apply, but only that of a rational Player 2 is consistent with Player 2's actual decisions across all games. Active application of ToM is required to test multiple perspectives and find that of a rational Player 2.

To conclude, our findings are congruent with findings from fMRI studies showing that mental state reasoning employs brain regions that differ from the regions involved in cognitive control (e.g., Apperly, 2011; Ramsey, Hansen, Apperly, & Samson, 2013; Saxe, Schulz, & Jiang, 2006). Our findings suggest that perspective taking requires additional cognitive resources, as opposed to just greater cognitive control, as one additional switch between perspectives induces not only longer reaction times but also qualitatively different decisions.

## Chapter 5

# What eye movements can tell about theory of mind in a strategic game

### Abstract

This study investigates strategies in reasoning about mental states of others, a process that requires theory of mind. It is a first step in studying the cognitive basis of such reasoning, as strategies affect tradeoffs between cognitive resources. Participants were presented with a two-player game that required reasoning about the mental states of the opponent. Game theory literature discerns two candidate strategies that participants could use in this game: either *forward reasoning* or *backward reasoning*. Forward reasoning proceeds from the first decision point to the last, whereas backward reasoning proceeds in the opposite direction. Backward reasoning is the only optimal strategy, because the optimal outcome is known at each decision point. Nevertheless, we argue that participants prefer forward reasoning because it is similar to causal reasoning. Causal reasoning, in turn, is prevalent in human reasoning (Gopnik et al., 2004).

Eye movements were measured to discern between forward and backward progressions of fixations. The observed fixation sequences corresponded best with forward reasoning. Early in games, the probability of observing a forward progression of fixations is higher than the probability of observing a backward progression. Later in games, the probabilities of forward and backward progressions are similar, which seems to imply that participants were either applying backward reasoning or jumping back to previous decision points while applying forward reasoning. Thus, the game-theoretical favorite strategy, backward reasoning, does seem to exist in human reasoning. However, participants preferred the more familiar, practiced, and prevalent strategy: forward reasoning.

This chapter was previously published *PLoS ONE* (2012).



## Introduction

Having a theory of mind (ToM) allows us to reason about other people's mental states, their knowledge, beliefs, desires, and intentions. This ability is helpful in social interactions, especially when our outcomes depend on the actions of others, and vice versa. Many studies have focused on the age at which ToM develops (Baillargeon, Scott, & He, 2010; Flobbe, Verbrugge, Hendriks, & Krämer, 2008; Onishi & Baillargeon, 2005; Perner & Wimmer, 1985; Wellman, Cross, & Watson, 2001; Wimmer & Perner, 1983), the proficiency of humans and nonhumans in ToM tasks (Call & Tomasello, 2008; Goodie, Doshi, & Young, 2012; Hedden & Zhang, 2002; McKelvey & Palfrey, 1992; Meijering, Van Maanen, Van Rijn, & Verbrugge, 2010; Meijering, Van Rijn, Taatgen, & Verbrugge, 2011; Premack & Woodruff, 1978; Zhang, Hedden, & Chia, 2012), and the brain regions associated with ToM (Gallagher & Frith, 2003; Saxe, 2006; Saxe, Schulz, & Jiang, 2006). In contrast, few studies have focused on the cognitive basis of ToM (Apperly, 2011; Apperly & Butterfill, 2009). Consequently, little is known about how inferences about mental states are achieved.

As findings from cognitive neuroscience have shown that participants in ToM tasks employ many brain regions rather than one single "ToM module" (Apperly, 2011; Gallagher & Frith, 2003; Saxe, 2006; Saxe et al., 2006), ToM reasoning probably consists of multiple serial and concurrent cognitive processes. Cost-benefit tradeoffs between these various resources will most likely have cascading effects on cognitive load (Borst, Taatgen, & Van Rijn, 2010) and thus ToM reasoning. Both task setting and strategies, in turn, have been shown to affect cost-benefit tradeoffs between cognitive resources (Fu & Gray, 2004; Gray, Sims, Fu, & Schoelles, 2006; Todd & Gigerenzer, 2000). Therefore, the study of strategies and task setting might be an appropriate first step in the study of the cognitive basis of ToM reasoning (Ghosh & Meijering, 2011; Ghosh, Meijering, & Verbrugge, 2010).

In this study, we investigate the ongoing process of ToM reasoning in a two-player game, referred to as Marble Drop (Meijering et al., 2010; 2011), see Figure 5.1. In this game, a white marble is about to drop, and each player's goal is that the white marble drops into the bin that contains the darkest possible marble of his or her allocated color. This is commonly known

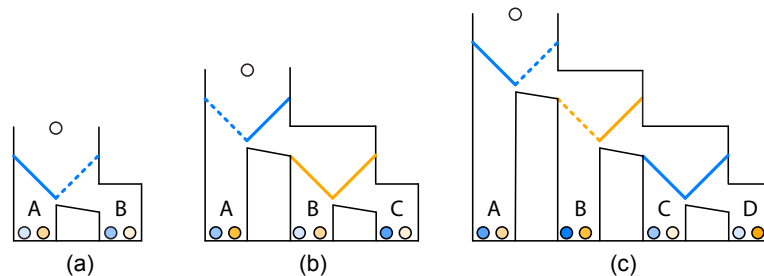


Figure 5.1: Examples of zero-order (a), first-order (b), and second-order (c) Marble Drop games. Each bin contains a pair of marbles, labeled A to D. For each player, the goal is that the white marble drops into the bin that contains the darkest possible marble of his or her allocated color. In this example, Player 1's marbles are blue, and Player 2's marbles are orange. Player 1 controls the blue trapdoors and Player 2 controls the orange trapdoors. The dashed diagonal lines represent the trapdoors that the players should decide to remove to obtain their maximum payoffs in these particular games.

among the players. Both players can remove trapdoors to control the path of the white marble. Marble Drop requires ToM because each player's outcomes depend on the decisions of the other player.

The example games in Figure 5.1 are of varying difficulty. With each additional decision point (i.e., set of trapdoors), the required reasoning becomes more complex. The game in Figure 5.1c is the most difficult, and requires second-order ToM. Below, we provide a possible reasoning scenario to explain how second-order ToM comes into play in this particular game.

By looking at payoff-pairs A to D in the game in Figure 5.1c, Player 1 will find out that B contains the darkest marble of his allocated color, blue. Player 1 has to ask himself whether that marble is attainable. In other words, Player 1 has to reason about whether Player 2 would remove the left orange trapdoor. Therefore, Player 1 has to look at the orange marbles in B to D to find out that D contains Player 2's darkest orange marble. ToM reasoning continues with Player 1 asking himself whether Player 2 thinks her orange marble in D is attainable. In other words, Player 1 has to reason about whether Player 2 thinks that he, Player 1, would remove the right blue trapdoor of the rightmost set of trapdoors. Player 1 knows that he would not remove that trapdoor, but that he would remove the left one instead. He also knows that Player 2 is aware of this, as both players are aware of each other's goals. Therefore, Player 1 knows that Player 2 knows that her darkest orange marble in D is unattainable. Therefore, Player 1 has to go back to the second decision point (i.e., the orange trapdoors). There, Player 2 would compare the orange marbles in B and C and decide to remove the left orange trapdoor, because the orange marble in B is the darkest orange marble that she can still attain. To conclude, Player 1 knows that his darkest blue marble in B is attainable, and will thus remove the right blue trapdoor of the leftmost set of trapdoors.

According to game theory literature there is just one strategy that undoubtedly yields the optimal outcome: reasoning by backward induction. We will refer to this strategy simply as *backward reasoning*. Backward reasoning proceeds from the last decision to be made back to the original problem or situation (Osborne & Rubinstein, 1994). The last decision in the game in Figure 5.1c is Player 1's decision between the blue marbles in payoff-pairs C and D. Player 1 would decide to remove the left trapdoor because C contains the darker blue marble. Backward reasoning would then proceed with the second-to-last decision, which is Player 2's decision between the orange marbles in payoff-pairs B and C. Player 2 would decide to remove the left orange trapdoor, because B contains the darker orange marble. Backward reasoning stops at the third-to-last decision, which is Player 1's decision between the blue marbles in payoff-pairs A and B. Player 1 would remove the right blue trapdoor, because B contains the darker blue marble. This scenario shows that backward reasoning is very efficient, because the optimal outcome is known at each decision point. Accordingly, few reasoning steps need to be retained, and working memory load would be small.

Game theory literature discerns another possible strategy, *forward reasoning*, but this strategy is not guaranteed to yield the optimal outcome (Ghosh & Meijering, 2011). Opposite to backward reasoning, the forward reasoning strategy starts at the first decision point in a game and blindly proceeds to the next for as long as higher outcomes are expected to be available at future decision points. A drawback of this strategy is that a player might not recognize the highest attainable outcome and continues the game to future decision points with lower outcomes. However, occasionally forward reasoning yields a quick solution, for example, if the maximum outcome is available at the first decision point.

Even though backward reasoning is the optimal strategy in games such as Marble Drop, it does not seem to be ubiquitous in human reasoning. In contrast, a forward progression seems to be more prevalent, for example in causal reasoning, where causes or decisions *lead* to possible effects. A well-known example of the persistency of causal reasoning is the *fundamental attribution error*, where causal explanations of observed behaviors are often dispositional despite more appropriate situational explanations (Kelley & Michela, 1980; Weber, Camerer, Rottenstreich, & Knez, 2001).

Given the prevalence of a forward direction in human reasoning, we expect that forward reasoning might also be a viable candidate strategy in Marble Drop games, even though backward reasoning is the game-theoretical favorite. However, forward reasoning would not always suffice to achieve the optimal outcome in Marble Drop. As explained above, a player might discover, while reasoning forwardly, that he or she unknowingly skipped the highest attainable outcome at a previous decision point. Thus, the player would need to jump back to inspect whether that outcome is indeed attainable. The procedure of jumping back to previous decision points is called *backtracking* (Brassard & Bratley, 1996). Backtracking superficially resembles backward reasoning, but it differs because jumping back to a previous decision point can be followed up with forward reasoning again. Note that our explanation of the Marble Drop game in Figure 5.1c followed the procedure of forward reasoning plus backtracking. Forward reasoning plus backtracking is less efficient than backward reasoning, because (at most stages in a game) multiple possible outcomes need to be retained to compare against next possible outcomes. Consequently, this strategy would cause high working memory load.

Besides the question which strategy is preferred (i.e., backward reasoning or forward reasoning plus backtracking), we investigate whether strategy preference can be influenced by task factors. The latter question is inspired by the work of Hedden and Zhang (2002). An important but also criticized aspect of that study was that each participant (assigned to the role of Player 1) was asked to predict the decision of Player 2 first, before making a decision (Colman, 2003). As this procedure *prompts* perspective taking, ToM reasoning might not have been completely spontaneous (Colman, 2003; Hedden & Zhang, 2002; Zhang & Hedden, 2003). In fact, we have shown that *prompting* participants for predictions indeed has a positive effect on performance (Meijering et al., 2010; 2011; see also Chapter 2 in this dissertation). In the current study, we investigate whether prompting may also have an effect on participants' preferences for any of the strategies.

Because Marble Drop has a predominantly visual interface and both strategies clearly predict a distinct succession in which the payoffs are to be compared, we employed eye tracking to measure the online (i.e., ongoing) process of ToM reasoning. Eye tracking has been used extensively in visual search tasks and reading tasks (Liversedge & Findlay, 2000; Rayner, 1998), and in complex visual problem solving tasks (Kong, Schunn, & Wallstrom, 2010; Nyamsuren & Taatgen, 2013). These studies have shown correlations between eye movements, on the one hand, and cognitive processes and higher-level strategies, on the other hand. For example, Kong et al. (2010) found a strong correlation between participants' visual working memory capacity and their eye movements while solving a nontrivial problem-solving task, the traveling salesman problem. Eye tracking has also been proven successful in exposing strategies in another complex (but non-social) reasoning task (Nyamsuren & Taatgen, 2013). Based on the eye movements of participants that played the game of SET, Nyamsuren and Taatgen (2013) were able to distinguish between bottom-up visual processes

and top-down planning processes. They were also able to detect in-game strategy shifts in participants.

An advantage of eye tracking is that it is an unobtrusive measure; participants were not constrained in any other way than in the original task setting. In contrast, other studies on online ToM reasoning required task modifications that may have influenced participants' strategies. For example, in Johnson, Camerer, Sen, and Rymon's computer task (2002), participants had to uncover task-relevant information that was hidden behind boxes displayed on the computer screen. The participants had to move the mouse cursor over a box to reveal the information behind it. Consequently, they might have felt disinclined to repeatedly move around the cursor to inspect each box's content. Tracking the eye movements (with a desk-mounted eye tracker) does not constrain participants so much.

In sum, the literature has identified one optimal strategy (backward reasoning), and we propose another (forward reasoning plus backtracking). Both strategies are clearly distinct from each other. This study aims to identify which strategy explains participants' performance in a ToM task best. It also investigates whether prompting participants for predictions has an effect on their strategies. We use eye tracking because it is an appropriate tool for showing whether the general direction of the eye movements, and thus reasoning, is either forward or backward.

## **Method**

### **Ethics Statement**

The Ethical Committee Psychology (ECP) of the University of Groningen approved this study. Written informed consent as approved by the ECP was obtained from each participant before conducting the experiment.

### **Participants**

Twenty-three first-year psychology students (14 female) with a mean age of 20.8 years (ranging from 18 to 24 years) participated in exchange for course credit. All participants had normal or corrected-to-normal visual acuity. None of the participants had difficulties distinguishing between the colors (blue and orange) presented in the experiment.

### **Stimuli**

Instead of using numerical payoffs, which are commonly used in strategic games, we chose for colored marbles to counter numerical but non-optimal strategies such as, for example, minimizing the opponent's outcomes, or maximizing the difference in Player 1 and Player 2 outcomes.

#### *Payoffs*

The payoffs were marbles of 4 different shades that could be ordered from light to dark. The

colors of the marbles were shades of orange and blue, taken from the HSV (i.e., hue, saturation and value) space. A sequential color palette was computed by varying saturation, for a given hue and value. This resulted in 4 shades (with saturation from .2 to 1) for both of the colors orange (hue = .1, value = 1) and blue (hue = .6, value = 1). The participants did not have any difficulties distinguishing between the shades of either color<sup>1</sup>.

### *Payoff structures*

The payoff structure (i.e., configuration of payoffs) and strategy preference determine the complexity of the reasoning required of Player 1, the participant. For example, a forward reasoning Player 1 immediately knows what to do if payoff-pair A contains his darkest marble: stop the game (i.e., remove the left-side trapdoor). In this case, Player 1 does not have to reason about Player 2's reasoning about Player 1. Therefore, we excluded this payoff structure, as it cannot inform us about second-order ToM. We only selected payoff structures that required Player 1 to reason about the decision at each of the three decision points (i.e., sets of trapdoors).

In line with Hedden and Zhang's criteria (2002), we considered payoff structures to be diagnostic of second-order ToM reasoning if, at the first set of trapdoors, second-order reasoning yielded a decision opposite to a decision based on first-order ToM reasoning. The payoff structures were balanced for the number of correct decisions to remove the left / right trapdoor, for both Player 1 and Player 2. The payoff structures are provided in Appendix A of Chapter 2.

## **Design**

The experiment consisted of three blocks: a training block and two test blocks. The training block was meant to familiarize participants with the rules of Marble Drop. In the first test block we manipulated whether participants were prompted to predict Player 2's decision. The first test block was followed by a second one, in which none of the participants had to make predictions anymore. This block was meant to measure the longevity of the effect of prompting participants for predictions.

## **Procedure**

Participants were seated in front of a 20-inch computer monitor, at 70 cm distance. An Eyelink 1000 eye-tracker was used to record the eye movements of the dominant eye, at a sample-rate of 500 Hz. The eye tracker was calibrated to each participant's dominant eye. Participants were always assigned to the role of Player 1. The target color, either blue or orange (marbles and trapdoors), was counterbalanced between participants. Participants were instructed that their goal was to maximize their payoffs, that is, to attain the darkest possible marble of their target color. Participants were told truthfully that they were playing against a computer-simulated

---

<sup>1</sup> The experiment was preceded by a block of 20 trials in which participants had to distinguish between the colors blue and orange, and between different shades of the colors blue and orange. They performed up to ceiling ( $M = 0.99$ ,  $SE < .01$ ), which implies that the participants did not have any difficulties distinguishing between colors and shades of colors.

Player 2<sup>2</sup>, whose goal was to maximize its payoffs. Participants were also instructed that the computer was programmed to look ahead and take into account the participant's last possible decision (i.e., Player 1's decision at the last set of trapdoors).

In the training block, participants were presented with 20 games of increasing difficulty. To familiarize the participants with the setup of the Marble Drop games, participants were first presented four trivial two-bin games that did not require ToM reasoning (Figure 5.1a). These two-bin games were followed by a set of eight three-bin games (Figure 5.1b), and a set of eight four-bin games (Figure 5.1c). The three-bin games require first-order ToM, because the participants have to reason about the decision of Player 2 at the second decision point (i.e., set of trapdoors). As discussed earlier, the four-bin games require second-order ToM. Each training game was played until either the participant or the computer decided to stop the game, by removing the left-side trapdoor, or until the last possible decision was made. After each game, participants were presented feedback displaying either "correct" if they obtained the darkest possible marble, or "incorrect" if they failed to do so. The feedback never indicated why a response was incorrect. Thus, participants had to find out themselves why an incorrect decision was incongruent with the other player's mental state. As the participants' performance on the eight four-bin games is indicative of their pre-experimental level of second-order ToM reasoning, we have included these items in the analyses.

Prompting participants for predictions was manipulated in the first test block, which consisted solely of second-order games. Participants were randomly assigned to either the so-called Prompt group (10 participants), or the so-called No-Prompt group (13 participants). Participants in the Prompt group were asked to enter their prediction of Player 2's decision at the second decision point before they were asked to enter their own decision at the first decision point. Participants in the No-prompt group were not explicitly asked to make any predictions. In this block, games stopped immediately after entering a decision. Feedback was presented after entering a prediction, if a prediction was queried, and after entering a decision. Feedback mentioned only whether a response was (in)correct. The first test block consisted of 32 trials; each of the 16 payoff structures was presented twice. The order was randomized.

The second test block was similar to the first one except that none of the participants were explicitly queried for a prediction anymore. This block also consisted of 32 trials.

## Results and discussion

### Behavioral results

Figure 5.2 depicts the mean accuracy of participants playing second-order Marble Drop games. The mean accuracy scores were analyzed by means of repeated-measures ANOVA. However, the scores were first arcsine-transformed to preserve homogeneity of variance. The analysis included the between-subjects factor prompting (No-prompt / Prompt) and the within-subjects factor block (Test Block 1 / Test Block 2).

In contrast to our earlier work (Meijering et al., 2011), the factor prompting was not significant,  $F(1, 21) = .1$ , ns. On average, asking participants to predict Player 2's decision did

<sup>2</sup> *Knowing whether the opponent was a computer player did not have an effect in Hedden and Zhang's (2002) study.*

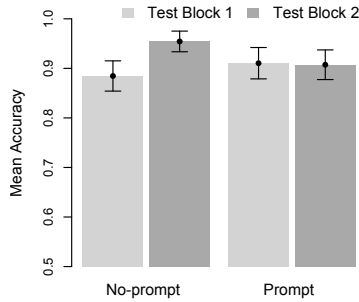


Figure 5.2: Mean accuracy in No-prompt and Prompt conditions, depicted separately for test blocks 1 and 2. Error bars represent standard errors.

not (positively) influence their performance. The lack of an overall effect of prompting might have been due to ceiling effects, as the mean accuracy was very high, around 90% in both test blocks.

The interaction between prompting and block was significant:  $F(1, 21) = 4.61, p = .044$ . On average, accuracy increased from Test Block 1 to Test Block 2,  $F(1, 21) = 5.09, p = .035$ , but that effect was mainly due to increasing accuracy in the No-prompt group. A possible explanation for the interaction might be that participants in the Prompt group, in contrast to participants in the No-prompt group, had to adjust to an experimental procedure that changed with each subsequent test block. This could have hindered their performance, which did not significantly differ between the two test blocks,  $t(9) = .12, ns$ .

### Eye tracking results

Eye movements were measured to distinguish between the strategies that participants may have used in second-order Marble Drop games, as backward and forward reasoning would clearly yield distinctive successions of fixations on each player’s payoffs. The default parameters of the Eyelink 1000 eye tracker were used to extract fixations from the eye movement data. Figure 5.3 gives an example of a participant’s succession of fixations in a particular game.

Each pair of payoffs was considered to be an area of interest (AOI). However, we did not define fixed AOIs with specific x and y coordinates. As the AOIs corresponding with the payoff-pairs are relatively small, a slightly inaccurate calibration of the eye tracker to a participant’s dominant eye would shift his or her fixations outside of the AOIs. Therefore,

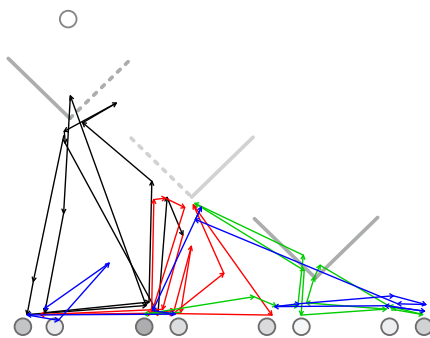


Figure 5.3: Example of a participant’s fixations in a particular game. The succession of fixations is indicated by arrows, which are superimposed on the payoffs and trapdoors (i.e., decision points). The first 15 fixations are depicted in black, fixations 16 – 30 in red, fixations 31 – 45 in green, and fixations 46 – 61 are depicted in blue. The succession of fixations on payoffs and trapdoors seems to indicate forward reasoning, followed by backtracking, which is indicated by the blue arrows that eventually go back to the first payoff pair.



cluster analysis was used to find four clusters of fixations in each participant's dataset, each cluster corresponding to a payoff-pair. The clustering algorithm used was a more robust version of k-means clustering (Kaufman & Rousseeuw, 1990). Fixations in the first (i.e., leftmost) cluster were labeled with the letter A; fixations in the second cluster were labeled with the letter B, and so forth. The labels are depicted above the payoff-pairs in Figure 5.1c. All following analyses solely include fixations that fall within these AOIs.

### *Onset times of fixations on payoff-pairs*

We analyzed the in-game times at which each cluster (i.e., payoff-pair) was first fixated, as these so-called onset times may indicate a general direction of reasoning in second-order Marble Drop games. The onset times were averaged across trials, separately for each participant (i.e., the 8 second-order trials from the practice block, and 32 trials from test block 2). The onset times were log-transformed, because their distribution was skewed to the right. The mean onset times (across participants) are depicted in Figure 5.4. We collapsed the data across the Prompt group and the No-prompt group, as there were no significant differences between these groups.

Figure 5.4a shows monotonically increasing onset times in the practice block, which indicates a forward (i.e., left-to-right) general direction of reasoning. All pairwise comparisons are significant, AB:  $p < .001$ ; AC:  $p < .001$ ; AD:  $p < .001$ ; BC:  $p < .001$ ; BD:  $p < .001$ ; CD:  $p = .028$ . The p-values are corrected by means of the Bonferroni-Holm method (Holm, 1979) to account for family-wise error rate.

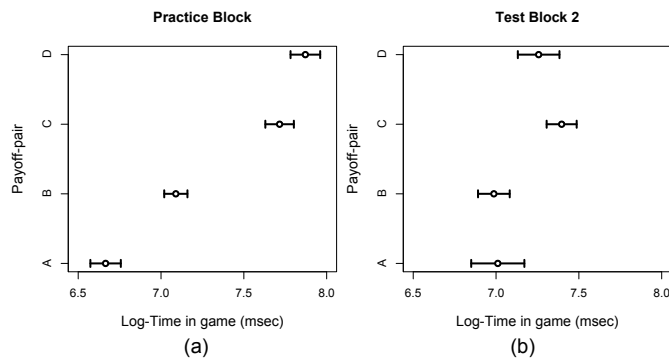


Figure 5.4: The logarithm of the onset times (in msec) of fixations on each payoff-pair. The onset times are depicted separately for the practice block (a) and Test Block 2 (b). The error bars represent standard errors.

Presumably, participants' strategies were most stable near the end of the experiment. However, the timing of the first fixations on each payoff-pair does not inform us on what these strategies might have been (see Figure 5.4b). The onset times do not increase monotonically anymore, in contrast to the onset times in the practice block. However, payoff-pairs A and B are still fixated earlier than payoff-pairs C and D. The average difference in onset times is significant,  $t(45) = -2.76$ ,  $p = .008$ .

As the onset times do not strongly correspond with either one of the candidate strategies, we analyzed the entire fixation sequences, which might reveal patterns corresponding to backward and/or forward reasoning.



*Fixation sequences*

Before presenting the statistics on the entire fixation sequences, we will first explain the statistical procedure, which involves several steps.

For each game, we predicted which payoffs would be fixated, and in which succession, given a particular strategy. The left panel of Figure 5.5 depicts an example game, the middle panel depicts fixation sequences that were predicted on the basis of backward reasoning, and the right panel depicts fixation sequences that were predicted on the basis of forward reasoning plus backtracking. For illustrative purposes, fixations on Player 2’s marbles were labeled with lowercase letters a, b, c, and d, and fixations on Player 1’s marbles with uppercase letters A, B, C, and D. Each line in the last two panels of Figure 5.5 represents a possible sequence of fixations given the corresponding strategy.

Backward reasoning yields eight possible fixation sequences for each individual game. Namely, a comparison between two payoffs can yield two possible successions of fixations, for example <D, C> versus <C, D>, and there are three comparisons to be made when applying backward reasoning. Thus, there is a total of two to the power of three, which is eight, possible fixation sequences. We granted forward reasoning plus backtracking the same degrees of freedom by applying the same procedure to the backtracking part, which is essentially the same as backward reasoning.

Given that we predicted fixations on individual marbles, we had to label each observed fixation for the specific marble that was fixated. We used cluster analyses to find two sub-clusters within each of the previously found payoff-pair clusters. Each left-side sub-cluster was considered to contain fixations on Player 1’s marbles, and each right-side sub-cluster was considered to contain fixations on Player 2’s marbles.

It is important to note that our implementations of the two strategies are idealizations, as we did not implement cognitive constraints such as, for example, working memory capacity. Consequently, the *predicted* fixation sequences did not contain repetitions. In contrast, the *observed* fixation sequences *did* contain repetitions, as participants would re-fixate payoffs if they had forgotten previously attended payoffs and comparisons. Figure 5.3 clearly shows an example of a participant repeatedly fixating payoffs. We accounted for these memory

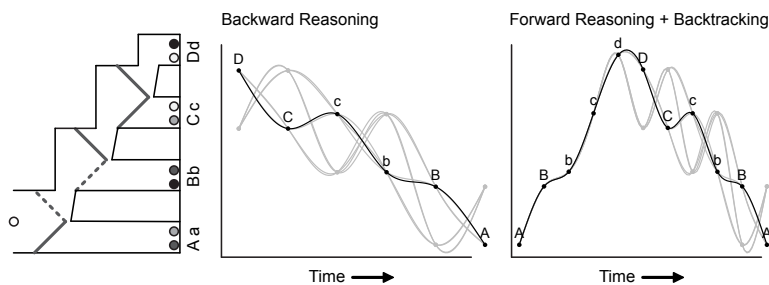


Figure 5.5: An example second-order Marble Drop game (left panel), and associated fixation sequences predicted on the basis of backward reasoning (middle panel) and forward reasoning plus backtracking (right panel). The fixation sequences represented by the black lines are annotated for AOI (A – D; a – d), and those represented by the grey lines are not. Player 1’s payoffs are labeled with uppercase A, B, C, and D. Player 2’s payoffs are labeled with lowercase a, b, c, and d. The sequences are depicted on “eye movement paths” for illustrative purposes.

effects by collapsing repeating patterns in the observed fixation sequences. For example, both *AAaBbCd* and *AaAaBbCd* would collapse to *AaBbCd*.

To evaluate how closely our predicted fixation sequences match the observed fixation sequences, we calculated the Levenshtein distance, which is the minimal number of insertions, deletions, and substitutions to get from one sequence to another. For example, if an observed fixation sequence for the game in Figure 5.5 would consist of AOIs  $\langle D, \mathbf{d}, C, c, b, B, A \rangle$ , we would find strong evidence in favor of backward reasoning, as it differs only one fixation (i.e., *d*) from one of the predicted sequences of AOIs  $\langle D, C, c, b, B, A \rangle$ . Importantly, the observed fixation sequence is compared with a set of eight predicted fixation sequences, thus eight Levenshtein distances are calculated, and the minimum Levenshtein distance is taken. To account for varying lengths of observed and predicted fixation sequences, the Levenshtein distance is normalized by dividing it by the length of whichever of the two sequences is longer, either the observed or the predicted one.

According to the procedure described above, the normalized Levenshtein distance was calculated for each individual trial (i.e., 32 trials per participant per test block). The normalized Levenshtein distance was averaged across trials, separately for each participant. Figure 5.6 depicts the mean normalized Levenshtein distance in Test Block 2, in which strategy preference is most stable.

We collapsed the data across the No-prompt group and the Prompt group, as the eye movement patterns did not significantly differ between these groups. Both the main effect of prompting and the interaction between strategy and prompting were not significant,  $F(1, 21) = .46$ , ns, and  $F(1, 21) = .71$ , ns, respectively. There are two possible explanations for this: Either prompting participants for predictions did not affect their strategy preference, or participants in the No-prompt group developed similar strategies on their own.

Figure 5.6 shows that, on average, the observed fixation sequences are most similar to the fixation sequences predicted on the basis of forward reasoning plus backtracking. The normalized Levenshtein distance is significantly larger for predictions based on backward reasoning,  $t(22) = 5.64$ ,  $p < 0.001$ . Figure 5.6 also depicts a baseline measure (dotted line), which is the average normalized Levenshtein distance between observed fixation sequences,

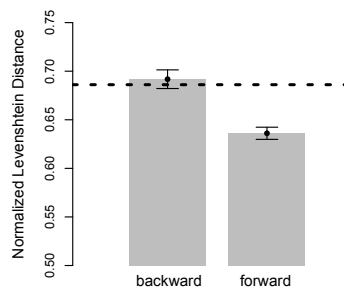


Figure 5.6: The average normalized Levenshtein distance between the observed sequence, on the one hand, and the closest of the set of predicted sequences, on the other hand. The dotted line is considered a baseline measure, which is the average normalized Levenshtein distance between an observed sequence and its randomized version. The error bars represent standard errors.

on the one hand, and each sequence randomized, on the other hand. Randomized sequences contain the same frequency of fixations as their observed counterparts, but nevertheless, forward reasoning plus backtracking fits the observed behavior significantly better than the baseline measure,  $t(22) = 4.91$ ,  $p < 0.001$ .

### *Sub-patterns in the fixation sequences*

To get a better idea of which specific components of the hypothesized strategies describe participants' reasoning best, we performed exploratory statistics on sub-patterns in the fixation data. The analysis concerns the fixation data from Test Block 2, as participants' strategies are assumed to be most stable in that test block. We will first describe the procedure of extracting sub-patterns from the fixation sequences, and then provide the results.

We analyzed sub-patterns of three subsequent fixations, as three is the minimal number of fixations that makes a pattern informative of either a forward or backward succession of comparisons between marbles. For example, subsequent fixations on payoff-pairs C, D, and B unambiguously indicate a backward succession of comparisons, even though the first two fixations seem to indicate a forward succession.

All subsequent triplets of fixations were extracted from each individual fixation sequence. If, for example, a trial consisted of fixations on payoff-pairs CDBCAB, sub-patterns CDB, DBC, BCA, and CAB were extracted. We considered fixations on payoff-pairs instead of fixations on individual payoffs (e.g., C versus c), as the latter would yield too many combinations with very low frequencies.

The results of the analyses are presented in Table 5.1, which shows the 50% most frequent *forward* and *backward* triplets. As can be seen in Table 5.1, the 50% most frequent triplets contain as many forward as backward triplets, and the frequencies of these triplets are quite similar. This seems to imply that, on average, participants made as many forward as backward comparisons between marbles.

We also analyzed the (in-game) onset times of forward and backward triplets, as these help us to determine whether forward and backward comparisons were made alternately, or forward comparisons first, followed by backward comparisons. Figure 5.7 depicts the relative

Table 5.1: The 50% most frequent forward and backward fixation triplets. The frequency of each triplet was divided by the total number of triplets,  $n = 6126$ , yielding the proportions given in the last column.

	Triplets	Proportion
<i>Forward triplets</i>	BCD	0.093
	ABC	0.058
	BDC	0.031
	ABD	0.024
	ACD	0.019
<i>Backward triplets</i>	DCB	0.079
	CDB	0.055
	CBA	0.047
	DCA	0.042
	DBC	0.022

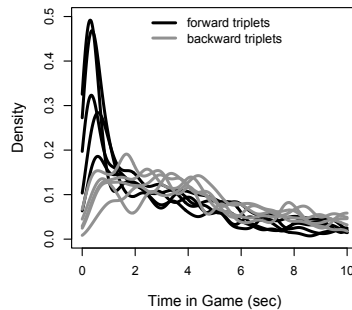


Figure 5.7: Densities of onset times of forward (black) and backward (grey) triplets.

likelihood (or probability density function) of onset times of all the triplets, and thus the likelihood of observing a particular triplet at a particular time in a game. Figure 5.7 clearly shows that all forward triplets have a relatively high likelihood of being observed early in a game, between zero and two seconds, whereas the highest likelihood of observing backward triplets is distributed over the entire range of 0 to 5 seconds.

Backward triplets correspond with either backward reasoning or the backtracking part of forward reasoning, depending on onset time. Early onset times would indicate backward reasoning, whereas late onset times would indicate backtracking. Figure 5.7 clearly shows that the densities of the backward triplets have less prominent peaks than the densities of the forward triplets. The flat likelihood distribution ranging from 0 to 5 seconds seems to imply that, at least in some games, backward reasoning was applied (indicated by early onsets). The finding that after 2 seconds the density functions of forward triplets are similar to those of the backward triplets implies that forward and backward comparisons were made equally often, presumably in alternating sequence. Figure 5.3 shows an example of such a pattern.

In sum, the densities in Figure 5.7 correspond best with the forward reasoning plus backtracking strategy. Given that the proportions of forward and backward triplets are quite similar, we can conclude that at early onset times forward triplets are more probable to be observed than backward triplets. In other words, payoffs are most likely to be compared in a forward succession until the last decision point is reached. Thereafter, backtracking takes place if the optimal outcome appears to be available at an earlier decision point. This succession, of forward comparisons followed by backward comparisons can be iterated multiple times until the highest attainable outcome is ascertained.

## General conclusions

We investigated strategy preference in a ToM task. Therefore, it was crucial that our task was successful at capturing ToM reasoning. Fortunately, mean accuracy was around 90%, close to ceiling, which means that the participants successfully applied (second-order) ToM in a large proportion of the trials (i.e., Marble Drop games).

Eye movements were measured to discern two candidate strategies with opposite general directions of reasoning: backward reasoning and forward reasoning plus backtracking. The onset times of the first fixations on each payoff-pair seem to imply that, in the practice block,

participants compared the payoffs in a forward succession. We analyzed the entire fixation sequences in the second test block and found that the forward reasoning plus backtracking strategy described the fixation sequences best. The observed fixation sequences were more similar to the fixation sequences predicted on the basis of forward reasoning plus backtracking than to the fixation sequences predicted on the basis of backward reasoning. Furthermore, by looking at sub-patterns in the fixation data, we found that, early in games, the likelihood of observing forward successions of comparisons between payoffs is higher than the likelihood of observing backward comparisons. These findings suggest that participants were applying forward reasoning, even though backward reasoning is the game-theoretical favorite strategy.

A possible explanation for a stronger preference for forward reasoning plus backtracking might be that backward reasoning requires deep structural knowledge of the task. Fu and Gray (2004) argued that in many interactive tasks, experts' behavior is rather dependent on, or even driven by, surface characteristics. Thus, the strong spatial and temporal structure of our task might have had a role in the adoption of forward reasoning (plus backtracking). Both the task display and the physics in Marble Drop games strengthen the intuitive and chronological direction of progressing decision points and comparing payoffs in a forward succession. Further research is needed to determine to what extent similar, or other, surface features might encourage the adoption of other strategies.

One could argue against forward reasoning by saying that the left-to-right (i.e., forward) fixations on the payoff-pairs merely represent a 'scanning phase' in which the payoffs are explored. However, this explanation does not hold since the participants kept fixating on the decision points (i.e., trapdoors) throughout the entire experiment. In fact, the fixations on the payoffs seemed to be interleaved with fixations on the trapdoors. Figure 5.3 provides an example of this pattern. For scanning purposes only, fixations on trapdoors are unlikely given that the trapdoors did not vary during the entire experiment (whereas the marbles did vary with each game). A more realistic and functional explanation for fixating trapdoors is reasoning, for example, about "*what would happen if the other player opened the left trapdoor*".

To conclude, ToM reasoning in games such as Marble Drop seems to progress in a forward succession, from causes (or possible decisions) to possible effects. Lacking a deep structural understanding of the logical problems posed in Marble Drop games, participants preferred to use a well-learned strategy, very similar to causal reasoning, even though it was not the most efficient strategy in this context.

## **Chapter 6**

# **Modeling inference of mental states: As simple as possible, as complex as necessary**

### **Abstract**

Behavior oftentimes allows for many possible interpretations in terms of mental states, such as goals, beliefs, desires, and intentions. Reasoning about the relation between behavior and mental states is therefore considered to be an effortful process. We argue that people use simple strategies and thus expend less effort as a way of dealing with limited cognitive resources. To test this hypothesis, we developed a computational cognitive model, which was able to simulate previous empirical findings: People start with simple strategies first, and only start revising their strategies when necessary. The model could simulate these findings by means of an interaction between factual knowledge and problem solving skills. At first, the model only considers its own goal, the most basic problem solving skill. Later, the model learns to attribute its problem solving skills to the other player, which only happens if its successes – stored as factual knowledge in declarative memory – do not increase anymore. The model was validated by means of a comparison with findings of a developmental study. This comparison showed that children use the same simple strategies that the model used. To conclude, the model was able to simulate two empirical findings: (1) People try to use simple strategies to infer mental states of others, and (2) they are able to improve such inference by attributing their own strategies to the other player.

This chapter was submitted to a journal and is currently under revision.

## Introduction

In social interactions, we try to understand others' behavior by reasoning about their goals, intentions, beliefs, and other mental states. Reasoning about mental states requires a so-called *theory of mind*, abbreviated ToM (Baron-Cohen, Leslie, & Frith, 1985; Wellman, Cross, & Watson, 2001; Wimmer & Perner, 1983). ToM has been implemented in computational cognitive models before (Hiatt & Trafton, 2010; Van Maanen & Verbrugge, 2010). However, these models either simulated one specific instance of ToM (Hiatt & Trafton, 2010) or attributed too much rationality to human reasoning (Van Maanen & Verbrugge, 2010). Here, we present a model that simulates application of various ToM strategies, ranging from simple strategies to full-blown recursive ToM. It is based on previous empirical results (Meijering, Van Maanen, Van Rijn, & Verbrugge, 2010; Meijering, Van Rijn, Taatgen, & Verbrugge, 2011) and is validated by means of a re-analysis of a previous developmental study by Flobbe et al. (2008). The model can explain why people use strategies that are relatively simple, while still being successful at inferring mental states of others.

Many studies have shown that people cannot always account for another's mental states in order to predict their behavior, particularly in the context of two-player sequential games (e.g., Flobbe et al., 2008; Hedden & Zhang, 2002; Raijmakers, Mandell, Van Es, & Counihan, 2013; Zhang, Hedden, & Chia, 2012). Sequential games require reasoning about complex mental states, because Player 1 has to reason about Player 2's subsequent decision, which in turn is based on Player 1's subsequent decision (Figure 6.1). Typically, performance is suboptimal and that is probably because players do not have a correct model of the other player's mental states (Johnson-Laird, 1983). By means of hypothesis testing, they may try to figure out which model works best in predicting the other player's behavior (Gopnik & Wellman, 1992; Wellman et al., 2001). However, a particular action or behavior can have many possible mental state interpretations (Baker, Saxe, & Tenenbaum, 2009), and testing all these interpretations strains our cognitive resources.

To alleviate cognitive demands, people generally start testing simple models or strategies that have been proven successful before (Todd & Gigerenzer, 2000). Because application of ToM and especially recursive ToM is an effortful process (Keysar, Lin, & Barr, 2003; Lin, Keysar, & Epley, 2010; Qureshi, Apperly, & Samson, 2010), reasoning about mental states probably also comprises the use of simple strategies. So where do these strategies come from? We hypothesize that they are a legacy of our childhood years. Raijmakers et al.'s (2013) findings corroborate this claim, as the children in their study consistently used strategies that were not fit to deal with the logical structure of the games presented to them. The strategies sometimes did yield the best possible outcome, however, which may be an explanation for why they still exist in adult reasoning: Simple strategies do not exhaust cognitive resources and are appropriate in a wide range of circumstances. Indeed, our computational cognitive model will show that the presence of simple strategies depends on the proportion of games in which they yield an optimal outcome.

In this study, we present a computational cognitive model that simulates inference of mental states in sequential games. The model initially uses a simple strategy that ignores many task aspects. However, if the model's strategy does not work, it learns to acknowledge that the other player has a role in its outcome. The model will therefore start attributing its own strategy to the other player. We will show that this process can account for the differential learning effects

in Meijering et al.'s study (2011; also see Chapter 2 in this dissertation), in which participants adopted distinct strategies based on the training regimen that was administered to them. To validate the model, the developmental study of Flobbe et al. (2008) was re-analyzed, searching for patterns that are indicative of the use of simple strategies in children.

Before we explain the model, we will first explain the empirical findings on which it is based.

## Empirical findings

Meijering et al. (2011) studied second-order ToM reasoning in two-player sequential games. Take the game in Figure 6.1 as an example game: Each end node contains a pair of payoffs, left-side payoffs belonging to Player 1 and right-side payoffs belonging to Player 2. The end node in which a game is stopped determines the payoff each player obtains in that particular game. Each player's goal is to obtain his or her greatest attainable payoff. As a player's outcome depends on the other player's decision, both players have to reason about one another's mental states. Participants are always assigned to the role of Player 1, and decide at the first decision point whether to stop the game at A or to continue to the next decision point, which is Player 2's decision between his payoff in B and his payoff in either C or D, which in turn depends on Player 1's decision between Player 1's payoffs in C and D. Thus, before making a decision at the first decision point, participants have to reason about Player 2, who in turn has to reason about Player 1's subsequent decision. In other words, participants have to apply second-order ToM when making a decision.

Meijering et al.'s study was based on the findings of Hedden and Zhang (2002; 2012) and Flobbe et al. (2008). Flobbe et al. had raised some concerns about Hedden and Zhang's

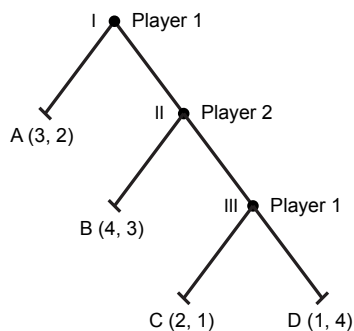


Figure 6.1: An extensive form representation of a two-player sequential game. Player 1 decides first, Player 2 second, and Player 1, again, third. The decision points are indicated in Roman numerals (I – III). Each end-node has a pair of payoffs, of which the left-side is Player 1's payoff and the right-side Player 2's payoff. Each player's goal is to obtain their highest possible payoff. In this particular game, the highest possible payoff for Player 1 is a 4, which is obtainable because Player 2's highest possible payoff is located at the same end node (i.e., B). Player 2's payoff of 4 is not obtainable because Player 1 would decide "left" instead of "right" at the third decision point (III).



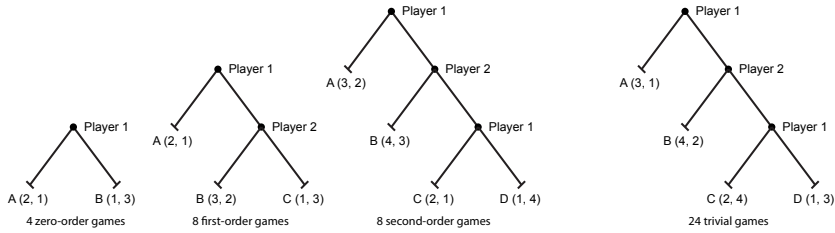


Figure 6.2: Extensive forms of example games (see Figure 6.1 for a detailed explanation). Stepwise Training consisted of 4 zero-order, 8 first-order, and 8 second-order games. Undifferentiated Training consisted of 24 trivial games. Each game had a unique distribution of payoffs.

training procedure, because it consisted of so-called trivial games (Figure 6.2; right panel), which are easier to play than truly second-order games such as in Figure 6.1. Trivial games are easier because Player 2 does not have to reason about Player 1's decision at III: Player 2's payoff in B is either lower or higher than both his payoffs in C and D. Consequently, Player 2 does not have to apply ToM, and Player 1 can suffice with first-order ToM. Flobbe et al. therefore argued that the training of Hedden and Zhang does not prepare people to play truly second-order ToM games. To test this claim, Meijering et al. administered two types of training procedures.

One group of participants was administered Hedden and Zhang's training procedure, which will henceforth be referred to as Undifferentiated Training, as all games had three decision points. The other group was administered Flobbe et al.'s training phase, but slightly modified (cf. Meijering et al., 2011). The latter training procedure will henceforth be referred to as Stepwise Training, as each additional decision point was introduced in subsequent blocks of games (Figure 6.2; left panel). Meijering et al. hypothesized that these training procedures would have distinct effects on strategy formation and thus performance. They predicted that Stepwise Training would facilitate participants to incorporate mental states of increasing complexity into their decision making process, yielding high accuracy. Undifferentiated Training, in contrast, would not motivate participants to develop recursive ToM, as they could suffice with application of first-order ToM. As expected, the participants that were assigned to Stepwise Training performed better than the participants assigned to Undifferentiated Training (see Figure 6.5).

One specific behavioral pattern is of particular interest to validate the model: The performance of participants assigned to Undifferentiated Training rose to ceiling during the training phase and dropped again when the experimental phase started (Figure 6.5). We hypothesize that the participants applied simple, child-like strategies during the training phase, because these strategies worked and did not consume much cognitive resources. At the start of the experimental phase, however, these strategies did not work anymore and accuracy dropped, because the games, while superficially similar, required more complex reasoning. Nevertheless, accuracy increased again over the course of the experimental phase, as the participants were able to revise their strategies. We will show that our computational cognitive model can simulate this process: The model's most important characteristic is that the complexity of its reasoning gradually increases by repeatedly attributing its own (evolving)

strategy to the other player.

## Computational cognitive model

The model<sup>1</sup> is implemented in the ACT-R cognitive architecture (Anderson, 2007; Anderson et al., 2004). ACT-R comprises a production system, which executes if-else rules, and contains declarative knowledge, which is presented as memory representations, or so-called chunks. In addition, ACT-R also includes modules that simulate specific cognitive functions, such as vision and attention, declarative memory, motor processing, et cetera. The results of these simulations appear as chunks in the modules' associated buffers, which the model continually checks (and manipulates) by means of its production system. ACT-R imposes natural cognitive constraints, as buffers can hold just one chunk at a time, and production rules can only fire successively, whenever their pre-specified conditions are matched. ACT-R does allow for parallel processing whenever a task induces cognitive processing in distinct modules. The model that we present here runs atop of ACT-R.

The model's behavior partially depends on memory dynamics. It needs to retrieve factual knowledge from declarative memory, and both the speed and success of retrieval depend on the so-called base-level activation of a fact (or chunk). The higher the base-level activation is, the greater the probability and speed of retrieval. The base-level activation in turn is positively correlated with the number of times a fact is retrieved from memory and the recency of the last retrieval.

The model simulates inference of mental states in sequential games. It uses a simple strategy at first and gradually revises that strategy until it can process recursive mental states. We consider the application of a particular strategy, and revising that strategy, to be deliberate processes. Therefore, application and revision are implemented by means of an interaction between factual knowledge and problem solving skills. Arslan, Taatgen, and Verbrugge (2013) successfully used a similar approach in modeling the development of second-order ToM in another ToM paradigm (i.e., the false-belief task). Van Rijn, Van Someren, and Van der Maas (2003) have successfully modeled children's developmental transitions on the balance scale task in a similar vein. Factual knowledge is represented by chunks in declarative memory, which store what strategy the model should be using. The problem solving skills, or strategy levels, are executed by (recursively) applying a small set of production rules. The model's goal is to make decisions that yield the greatest possible payoff. Decisions are either 'stop the game' or 'continue it to the next decision'. The model was presented with the same distributions of payoffs (i.e., items) as were presented to the participants.

The model's initial simple strategy is to consider only its own decision at the first decision point and to disregard any future decisions. The model's decision is based on a comparison between its (i.e., Player 1's) payoff in A and the maximum of its payoffs in B, C, and D. If the model's payoff in A is greater, the model will decide to stop. Otherwise, the model will decide to continue. By using this simple strategy the model seeks to maximize its own payoff, which can be considered a direct translation of the instructions given to the participants.

This strategy will work in some games but not in all. Whenever the strategy works, the model receives positive feedback and stores in declarative memory what strategy it is currently using. In fact, the model stores a strategy level, which is level-0 in the case of the

<sup>1</sup> The model can be downloaded from <http://www.ai.rug.nl/~meijering/iccm2013>

simple strategy described above. Whenever the strategy does not work, the model receives negative feedback and stores in declarative memory that it should be using a higher strategy level (e.g., level-1).

The higher strategy level means that the model should attribute whatever strategy it was using previously to the other player at the next decision point. In the case of strategy level-1, the model attributes the model's initial simple strategy (i.e., level-0) to Player 2. Accordingly, the model is applying first-order ToM, as it reasons about the mental state of Player 2, who considers only his own payoffs and disregards any future decisions.

Again, this strategy will work in some games but not in all. Whenever it does not work, the model receives negative feedback and stores in declarative memory that it should be using a higher strategy level. At a higher strategy level, the model will attribute whatever strategy level it was using previously to Player 2. At strategy level-2, the model attributes strategy level-1 to Player 2, who in turn will attribute strategy level-0 to the player deciding at third decision point: Player 1. Now the model is applying second-order ToM.

## Assumptions

The model is based on two assumptions. The first assumption is that participants, unfamiliar with sequential games, start playing according to a simple strategy that consists of one comparison only: Participants compare their current payoff, when stopping the game, against the maximum of all their future payoffs, when continuing the game. This strategy can be considered the simplest possible strategy, as participants who are using it ignore the consequences of any possible future decision, whether their own or the other player's.

If participants obtain expected outcomes, they do not have to revise their strategy. However, if participants obtain unexpected outcomes, they have to acknowledge that the unexpected turn of events was caused by the other player deciding at the next decision point. Reasoning about the other player, participants can only attribute a strategy they are familiar with themselves. This is our second assumption, which is based on variable frame theory (Bacharach & Stahl, 2000). Imagine a scenario in which two persons are asked to select the same object from a set of objects with differing shapes and colors but one person is completely colorblind. The colorblind person cannot distinguish the objects based on color, nor can he predict how the other would do that. Therefore, the colorblind person can only predict or guess what object the other would select based on which shape is the least abundant. The seeing person should account for the colorblind person's reasoning and also choose the object with the least abundant shape. This variable frame principle also applies to reasoning about others: We can only attribute to others goals, intentions, beliefs, and strategies that we are familiar with ourselves.

## Mechanisms

The simple strategy is implemented in two production rules. The first production rule determines what the payoff will be when stopping the game; the other production rule determines what the highest future payoff could possibly be when continuing the game. Both productions are executed from the perspective of whichever player is currently deciding (Figure 6.3). The model will attribute this simple strategy from the current decision point to

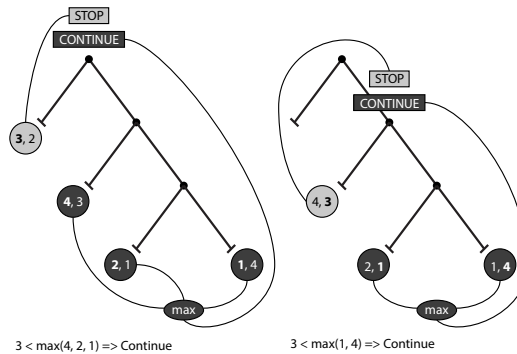


Figure 6.3: Depiction of the simple strategy. In the left panel, the model compares its payoff if it would stop (light grey) against its maximum possible payoff if it would continue (dark grey). In the right panel, the model compares Player 2's payoff if Player 2 would stop (light grey), against Player 2's maximum possible future payoff (dark grey). The left panel schematically represents the application of zero-order ToM, and the right panel the attribution of zero-order ToM to the other player.

the next, each time the model updates its strategy level (i.e., incrementing strategy level by one). The model will thus heighten its level, or order, of ToM reasoning.

### Zero-order ToM

Before the model starts applying its strategy, it needs to construct a game state representation to store the payoffs that are associated with a *stop* and *continue* decision, respectively. To construct a game state, the model first retrieves from declarative memory what strategy level it is currently using. At the beginning of the experiment, strategy level has a value of 0, which represents the simple strategy. After retrieving strategy level, the model constructs its current game state.

Starting with the simple strategy, the model will determine its own *stop* and *continue* payoffs (see Figure 6.3, left panel), which will be stored in the game state representation. The model will then compare these payoffs and make a decision. After the model has made a decision, it will update declarative memory by storing what strategy level the model should be playing in the next game: If the model's decision was correct, the model should continue playing its current strategy level; otherwise the model should be playing a higher strategy level.

After playing a couple of games in which the simple strategy (i.e., level-0) does not work, the higher strategy level (i.e., level-1) will have a greater probability of being retrieved, as its base-level activation increases more than the simple strategy's base-level activation. At the start of the next few games, before the model constructs its game state, it will begin retrieving strategy level-1 from declarative memory.

### First-order ToM

Playing strategy level-1, the model will first determine what payoff is associated with a *stop* decision at the first decision point (I). However, before determining what payoff is associated with a *continue* decision, the model needs to reason about the future and therefore consider the next decision point (II). It attributes strategy level-0 to Player 2, who is deciding at II.

Later, the model will return to the first decision point and determine what payoff is associated with a *continue* decision.

At II, the model will apply strategy level-0, but from the perspective of Player 2 (Figure 6.3, right panel). When reasoning about Player 2's decision, the model constructs a new game state, which references the previous one. The previous game state is referenced, because the model needs to jump back to that game state and determine what payoff is associated with a *continue* decision in that game state. At II, the model will execute the same production rules that it executed before when it was playing according to strategy level-0: It will determine what payoffs are associated with a *stop* and a *continue* decision, but from the perspective of Player 2.

The model will not produce a response whenever it determines the *stop* and *continue* payoffs at II, because the problem state at II references a previous one (i.e., I). The model will therefore backtrack to the previous game state representation, which did not yet have a payoff associated with a *continue* decision. That payoff can now be determined based on the current game state (i.e., Player 2's decision). The model will retrieve the previous game state from declarative memory.

After retrieving the previous game state representation, the model has two game states stored in two separate locations, or buffers: The current game state is stored in working memory, or the *problem state buffer* (Anderson, 2007; Borst, Taatgen, & Van Rijn, 2010), and the previous game state is stored in the *retrieval buffer*, which belongs to the declarative memory module. The model will determine what payoff is associated with a *continue* decision in the previous game state (stored in the retrieval buffer) given the decision based on the current game state (in the problem state buffer). It will update the previous game state and store it in working memory.

Playing strategy level-1 and being back in the previous game state, there is no reference to any previous game state and the model will make a decision based on a comparison between the payoffs associated with the *stop* and *continue* decisions. As explained previously, the model will stop if the payoff associated with stopping is greater; otherwise the model will continue.

Again, after the model has made a decision, it will update declarative memory by storing what strategy level the model should be playing in the next game(s). If the model's decision is correct, it will apply the current strategy level. Otherwise, the model will revise its strategy level by storing in declarative memory that it should be using strategy level-2 in the next game(s).

### *Second-order ToM*

The model will first determine what payoff is associated with stopping the game and then consider the next decision point. There, the model proceeds as if it were playing strategy level-1, but from the perspective of Player 2. In other words, the model is applying second-order ToM.

The strategy described above closely fits the strategy of forward reasoning plus backtracking (Meijering, Van Rijn, Taatgen, & Verbrugge, 2012; Chapter 5 in this dissertation). Meijering et al. (2012) conducted an eye-tracking study, and participants' eye movements reflected a forward progression of comparisons between payoffs, followed by backtracking to previous decision points and payoffs. Such forward and backward successions are present in strategy level-2 as well: Payoffs of *stop* decisions are determined one decision point after another, and this forward succession of payoff valuations is followed by backtracking, as payoffs of previous

*continue* decisions are determined in backward succession.

## Results

The model was presented with the same trials as in Meijering et al.'s (2011) study (see also Chapter 2 in this dissertation), with stepwise training versus undifferentiated training as a between-subjects factor. The model was run 100 times for each training condition. Each model run consisted of 20 (stepwise) or 24 (undifferentiated) training games, followed by 64 truly second-order games. The results are presented in Figures 6.4 and 6.5.

Figure 6.4 shows the proportions of models that apply strategy levels 0, 1, and 2, calculated per trial. The left panel of Figure 6.4 shows the output of the models that received 24 undifferentiated training games before playing 64 second-order games. As can be seen, initially all models apply strategy level-0, corresponding with zero-order ToM, but that proportion decreases quickly in the first couple of games. The proportion of models applying zero-order ToM decreases because that strategy yields too many errors, which can be seen in Figure 6.5. The models store in declarative memory that they should be using strategy level-1, but it takes a few games before the base-level activation of the level-0 chunk drops below the retrieval threshold. After it does, the models start retrieving level-1 chunks and will apply strategy level-1, which corresponds with first-order ToM. The proportion of models that use strategy level-1 increases up to 100% towards the end of the 24 undifferentiated training games. The models do not start applying strategy level-2 during the training phase, because strategy level-1 yields correct decisions in all undifferentiated training games, which can be seen in Figure 6.5. However, in the experimental games, which are truly second-order games, strategy level-1 yields too many errors, and accuracy drops. It takes approximately 40 games before the base-level activation of the level-1 chunk has dropped below the threshold in at least half of the models. The models gradually start using strategy level-2, and accuracy starts to increase again, as can be seen in Figure 6.5.

The right panel of Figure 6.4 shows the output of the models that were presented with 20 stepwise training games (4 zero-order, 8 first-order, and 8 second-order games) before playing

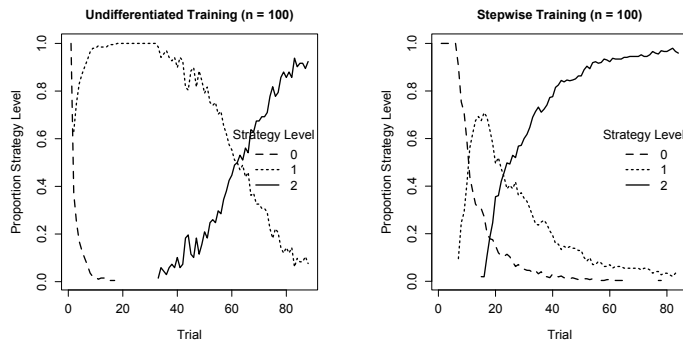


Figure 6.4: Proportion of models that apply strategy levels 0, 1, and 2; plotted as a function of trial. The left panel depicts these proportions for the model that received undifferentiated training; the right panel depicts the proportions for the model that received stepwise training.

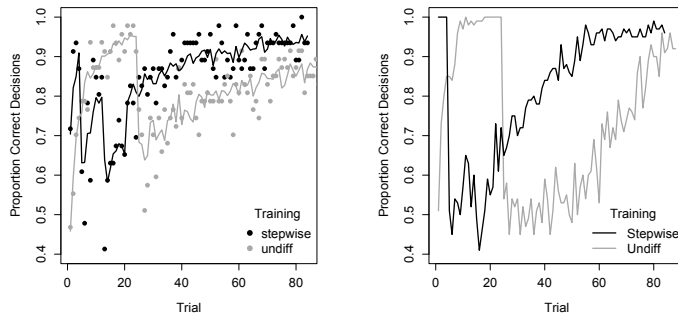


Figure 6.5: Proportion of correct decisions, or accuracy, across participants (left panel) and models (right panel). The solid lines in the left panel represent the fit of the statistical model, which is added to visualize the proportion trends.

64 second-order games during the experimental phase. As can be seen, all models start applying strategy level-0, and they use it longer than the models that received undifferentiated training. The reason is that strategy level-0 yields a correct answer in the first four games during stepwise training, because those are zero-order games. As can be seen in Figure 6.5 (right panel), accuracy is 100% in the first few games. In the next eight first-order training games (Trials 5 – 12), the proportion of models that apply strategy level-0 decreases, as strategy level-0 yields too many errors. Simultaneously, the proportion of models applying strategy level-1 increases, as the base-level activation of the level-0 chunk decreases and the models start retrieving the level-1 chunk. In the next eight second-order training games (Trials 13 – 20), the proportions of models that apply strategy level-0 and level-1 decrease, as both strategy levels yield too many errors. Simultaneously, the proportion of models that apply strategy level-2 increases. As strategy level-2 yields a correct decision in the remainder of the games, accuracy increases up to ceiling, which can be seen in Figure 6.5 (right panel).

The accuracy trends in the models' output qualitatively fit those of Meijering et al.'s study (2011). The quantitative differences are probably due to the fact that not all participants started out using the simple strategy, whereas all models did. One possible explanation is that some participants started with intermediate-level strategies and, due to large proportions of optimal outcomes, did not proceed to the highest level of reasoning. We could account for this by storing level-0, level-1, and level-2 chunks in declarative memory, and having the base-level activation of these chunks follow the distribution of zero-order, first-order, and second-order ToM in the adult population. A meta-review of (higher-order) ToM in adults and children may be a good starting point to find the appropriate distributions. Nevertheless, the qualitative trends in the model data, changing as a function of game complexity, correspond with the response patterns in the behavioral data. The trends suggest that people use simple strategies for as long as these yield expected outcomes.

In the introduction we hypothesized that simple strategies are a legacy of our childhood years, and that adults keep using those strategies that have proven themselves successful during development. To test this hypothesis, we have re-analyzed the data from Flobbe et al.'s (2008) developmental study. We expected that few children would have sufficient cognitive resources to apply second-order ToM, and that performance levels would therefore align well with lower and intermediate strategy levels. The most obvious prediction is that prevalence



of level-0, level-1, and level-2 strategies can be ranked, where level-0 is the most dominant strategy and level-2 is least frequent.

## Developmental study

Flobbe et al. (2008) studied the application of second-order ToM in children that were in between 8 and 10 years ( $M = 9;2$ ). They presented the children with sequential games, and performance was just above chance-level (57% correct). As children of age 9 are at the brink of mastering second-order ToM (Flobbe et al., 2008; Miller, 2009; Perner & Wimmer, 1985), we expect the lower and intermediate strategies to be most prevalent in Flobbe et al.'s study, which is thus perfect to validate our model.

We hypothesize that children apply the same simple strategies that are implemented in our computational cognitive model. We predict that the children start out with the simplest (i.e., zero-order) strategy, and that some will learn to attribute that strategy to the other player. Probably few children will learn that the other player, in turn, attributes the simple strategy to the player who decides next (i.e., to them). As each child was first asked to predict the other player's decision, before they were asked to make a decision themselves, we have a direct measure of the child's perspective of the other player's strategy. We will analyze both the predictions and the decisions.

## Predictions

We applied a binomial criterion to reliably categorize a participant's predictions as belonging to either level-1 or level-2: The predictions in at least 8 out of 10 consecutive games had to be congruent with one particular strategy level to label the predictions accordingly. This might seem strict, but 8 out of 10 is the minimum quantile that is still significant with a significance level of 0.05. As the experiment consisted of 40 second-order games, we categorized each child's responses in 4 sets of 10 games. Figure 6.6 depicts the proportion of children that applied either first-order or second-order ToM. These ToM-orders correspond with level-1 and level-2 in the computational model.

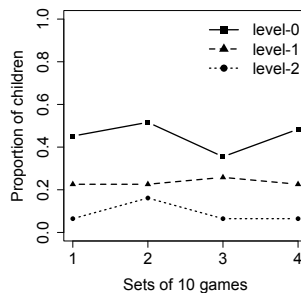


Figure 6.6: The proportion of children that applied zero-order ToM (level-0), first-order ToM (level-1) or second-order ToM (level-2) to the other player; depicted in 4 consecutive sets of 10 games.



Note that sets of predictions that could not be categorized level-1 or level-2 do not necessarily imply the use of level-0, because the predictions in those sets could have been completely random, or a mixture of the various strategy levels. The decisions are therefore analyzed to determine the prevalence of strategy level-0.

As can be seen in Figure 6.6, the proportion of children that applied first-order ToM by attributing strategy level-0 to the other player is greater than the proportion of children that applied strategy second-order ToM. Furthermore, many children's predictions could not be labeled according to one of the strategies at all (13 out of 40). These children probably switched frequently between multiple possible perspectives, and such switching is difficult to reliably capture by means of a statistical model. Nevertheless, most of the children whose responses could be categorized, were applying first-order ToM by attributing the simple (i.e., level-0) strategy to the other player. Almost none of the children was able to consistently attribute strategy level-1 to the other player, thereby applying second-order ToM.

## Decisions

As explained above, the predictions required application of first-order ToM at minimum and could therefore not be indicative of zero-order ToM. Therefore, the decisions were analyzed to determine how many children applied zero-order ToM, ignoring the other player entirely. Again, we categorized the decisions based on the binomial criterion that at least 8 out of 10 consecutive responses should be consistent with application of zero-order ToM (i.e., level-0 in the model). As can be seen in Figure 6.6, most of the children that consistently responded according to one of the strategies applied zero-order ToM when making a decision. This is remarkable, because each child that participated in the experimental phase successfully passed a training block in which they were required to apply first-order ToM. This finding suggests that the children could not see how first-order ToM would fit in the more complex games in the experimental blocks. They may have recognized that it did not work, but still could not revise their strategy to incorporate an additional ToM level.

To conclude, a re-analysis of Flobbe et al.'s (2008) study shows that few children were able to apply second-order ToM (level-2), and that most children used simple strategies. The most dominant strategy was the simplest one that did not account for any future decision points. Most children seemed to apply zero-order ToM (level-0) while making a decision. Some children, though, were able to attribute that simple strategy to the other player, thereby applying first-order ToM (level-1). These strategies are the same as those implemented in our computational cognitive model. The model is thus supported in two ways: (1) Its most simple strategies are found in children, and (2) it learns to revise its strategies as adults do.

## Conclusions

In this study we presented a computational cognitive model that simulates inference of mental states in sequential games. More specifically, the model was required to apply ToM recursively, a skill that appears to be unique to human intelligence. Many studies have shown that people oftentimes fail to apply ToM to interpret the behavior of others (e.g., Apperly et al., 2010; Keysar et al., 2003; Lin et al., 2010). In this study, in contrast, we show that people do not necessarily fail to apply ToM, but rather first apply simple strategies that are computationally

less costly. Only when necessary do people revise their strategies to account for complex mental states.

The model is based on previous empirical findings (Meijering et al., 2011) that seemed to imply that people exploit the possibility of using simple strategies for as long as these pay off. We implemented one such simple strategy that ignores any future decisions and simply compares the immediate payoff, when stopping a game, against the maximum of all future possible payoffs. By means of simple memory dynamics the model either retrieves a chunk that specifies that the model should continue using this strategy, or chunks that specify that the model should attribute the simple strategy to the player who decides next. Although this updating process may seem simplistic at first sight, the model does gradually master second-order ToM, but only because that is required in the games in this study. In other words, the model's most important dynamics are not task-specific, and because of that, the model is flexible and can accommodate many other two-player sequential games.

We found support for the model in the data from Flobbe et al.'s (2008) developmental study in which 9-year-old children were presented with similar sequential games. Most children used the simple, level-0, strategy when making a decision. The second-most prevalent strategy was the level-1 strategy. Using that strategy, the children attributed the simplest possible strategy (i.e., level-0) to the other player. Few children were able to apply second-order ToM mind. They did not recognize that the other player, in turn, attributed the simplest strategy (i.e., level-0) to them. These findings show that the children used the same simple strategies as the adults initially used in Meijering et al.'s study. However, the adults were able to revise their strategies to achieve the highest required level of ToM reasoning, whereas the children may not have had sufficient cognitive resources to achieve that same level of reasoning.

Our notion of zero-order ToM (i.e., strategy level-0) closely maps with the instruction given to the participants: to maximize their payoff. This strategy corresponds with a risk-seeking perspective, because it does not account for the fact whether higher future payoffs are actually attainable. There are other notions of a level-0 strategy, however. A risk-seeking strategy can be contrasted with a risk-averse strategy according to which one would stop if there were any lower future payoffs. There is still another notion of a level-0 strategy: Hedden and Zhang (2002; 2012) defined a so-called myopic level-0 strategy that only considers the current payoff and the closest future payoff. Player 1, for example, would only compare his payoffs in A and B, ignoring his payoffs in C and D. These strategies, however, are almost non-existent in Flobbe et al.'s dataset.

The findings from this study raise the question why younger children of 6 to 8 years are perfectly capable of accounting for second-order mental states in traditional false-belief studies (Coull, Leekam, & Bennett, 2006; Flobbe et al., 2008; Perner & Wimmer, 1985; Sullivan, Zaitchik, & Tager-Flusberg, 1994), as well as when they are asked to discriminate between ironic and deceptive speech acts (Winner & Leekam, 1991). One possible explanation is practice: Children have encountered false beliefs, irony, and deception more often than games such as in this study. Another explanation is that games can have a large space of possible outcomes, which requires extensive reasoning. False-belief stories and speech acts, on the other hand, are a given and thus require fewer computations. On a related note, children are better at reasoning about past events than about future possible outcomes (e.g., McColgan & McCormack, 2008; Suddendorf, Nielsen, & Gehlen, 2011). Reasoning about past events can be considered a linear traversal backwards in time, whereas reasoning about future events may follow an expanding tree-like structure.

This study has at least two methodological implications: One, experimenters should be careful in selecting ‘practice’ items, as participants exploit the possibility of using simple strategies when possible. Two, average proportions of correct answers, a popular statistic in most ToM studies, may not be as informative as a categorization of responses (also see Raijmakers et al., 2013). Flobbe et al., for example, reported that performance was just above chance-level (i.e., 57% correct), and the most common interpretation would be “on average children were able to apply second-order ToM in 57% of the games.” However, the current study shows that this score can be obtained if 1 or 2 children are applying second-order ToM and most of them below-optimal strategies such as zero-order and first-order ToM.

The theoretical implication of this study is that people do not necessarily perceive sequential games in terms of interactions between mental states. They know that there is another player making decisions, but they have to learn over time, by playing many games, that the other player’s depth of reasoning could be greater than initially thought. Learning takes place when people obtain unexpected outcomes and start recognizing that the other player has a role in their outcomes. They will have to attribute their own, simple, strategies to the other player, thereby developing increasingly more complex strategies themselves. Over time, reasoning will become as complex as necessary, as simple as possible.

## **Chapter 7**

### **Summary and discussion**

## Discussion

This dissertation details an investigation of higher-order theory of mind in adults. It is high time to study higher-order theory of mind, as it has not yet received as much attention as first-order theory of mind (ToM). Nevertheless, higher-order ToM is not some exotic cognitive function, and people need it to engage in complex social interactions. Communication, for example, may already require higher-order ToM if the wording is ambiguous, which is not that uncommon in language. To find the most probable meaning of an ambiguous utterance, the listener has to reason about the speaker's beliefs, and account for the fact that the speaker in turn may have reasoned about the listener's knowledge.

The study of ToM in adults also needs more attention, as ToM has mostly been studied in infants and children. That is not so surprising as first-order ToM develops around the age of 4, and second-order ToM between 6 and 8. Nevertheless, adults still frequently fail to account for the mental states of others, and it is not yet evident why that happens. For example, do adults not *have* a complete theory of mind, or do they not have sufficient cognitive resources to *apply* it? This dissertation provides new insights into why adults may not always apply ToM despite the fact that they have already mastered it.

In contrast to many other studies, theory of mind was investigated by means of two-player games, instead of false-belief tasks. One obvious advantage of these games is that they do not require language processing as much as the stories in false-belief tasks. Another advantage is that games can be presented many times, in various configurations, to the same participants. One concern has been that some games do not strictly require mental state reasoning. However, Chapters 3 and 4 have demonstrated that people did interpret the two-player games in terms of mental states. The study in Chapter 4, for example, shows that people do not consider a rational computer player to be the same as a completely deterministic device, even though the outcome of both was based on the same principle. In sum, the paradigm of two-player games has proven to be successful in examining various characteristics of theory of mind.

## Cognitive constraints

Some studies suggest that the development of ToM in children involves a conceptual change (Gopnik & Slaughter, 1991; Gopnik & Wellman, 1992; 2000; Wellman, Cross, & Watson, 2001; Wimmer & Perner, 1983). According to this account, children first experience desires and perceptions as simple causal links between them and the world, at age 2. Later, children learn about beliefs, and even think of these as representations, at age 3. However, they cannot yet incorporate belief representations into their 'core' theory of mind, which still reflects a rather direct causal link between them and the world. At the age of 4, they learn to perceive that representations (of mental states) are the basis of psychological function.

An important prediction of this account is that if development of ToM would involve a conceptual change, ToM would not be susceptible to improvement when it has already reached maturity. The findings in this dissertation (see Chapter 2), however, imply that application of ToM is a computational process that can benefit from supporting structure. Application of higher-order ToM improved when adults were trained to account for increasingly more

complex mental states, when they were prompted to take another's perspective, and when they were provided with visual cues as to how decisions are dependent on mental states. These findings show the ability to apply higher-order ToM is susceptible to improvement, and thus ToM may not be a fixed skill after all. However, adults and children do need sufficient cognitive resources to put that ability into practice.

The study in Chapter 3 shows that the application of ToM may involve a specialized cognitive function to infer the mental states of self and others. Previously, developmental studies suggested that unsuccessful inference of mental states reflects a broader problem with representations in general (Leekam, Perner, Healey, & Sewell, 2006; Perner & Leekam, 2008; Sabbagh, Xu, Carlson, Moses, & Kang, 2006; Todd & Gigerenzer, 2000). The findings in Chapter 3, however, show that reasoning about someone else's decision-making is more difficult than making the same decision oneself, even though the required reasoning steps are the same. Apparently, switching perspective makes the decision problem more difficult. One possible explanation is that the representation of a particular decision-making problem becomes more elaborate as the complexity of the involved mental states increases.

The findings in Chapter 4 corroborate the conclusions of the study in Chapter 3. Participants were presented with two-player games in which they had to reason about another player. The findings show that performance depended on whether the other player was reasoning about the participant's decision or, instead, about a mechanism. We hypothesized that the other player's mental states would be easier to infer if he would be reasoning about a mechanism, because the mechanism was completely deterministic. In contrast, if the other player would be reasoning about the participant, the participant would have to reason about a multitude of ideas the other player could be having about her. Importantly, both conditions were completely isomorphic with respect to the required reasoning steps. Still, the results showed that the response times were shorter in the mechanism games. We argue that the non-mechanism games were more difficult to solve because people had to test many possible mental state interpretations. In the mechanism games, in contrast, the other person's actions were dependent on a deterministic mechanism and people could therefore test fewer possible mental state interpretations and respond faster. Besides differential response times, the types of errors qualitatively differed between mechanism and non-mechanism games.

In sum, the studies in Chapters 3 and 4 showed that the complexity of mental states, all other task aspects controlled for, caused differential cognitive processes. The more intricate the involved mental states were, the worse the performance was, which suggests that application of ToM consumed cognitive resources. The studies in later chapters show how people try to preserve cognitive resources and still perform well.

## **Cognitive processes**

As application of ToM and especially higher-order ToM are considered to be effortful processes (see Chapters 2, 3, and 4), it is not surprising that people use simple strategies to reduce demands on cognitive resources. Todd and Gigerenzer (Meijering, Van Rijn, Taatgen, & Verbrugge, 2012; Szymanik, Meijering, & Verbrugge, 2013; see also Chapter 5 in this dissertation; Todd & Gigerenzer, 2000) already argued that people use simple strategies or heuristics to solve many (non-social) tasks. Chapters 5 and 6 show that this may be true

for inference of mental states as well: People start out reasoning about simple mental states, using basic strategies, and only account for more sophisticated mental states if their simple strategies do not yield desirable outcomes anymore.

In the study in Chapter 5, we tracked people's eye movements during a two-player game in which they had to infer the other player's mental states. The eye movements were analyzed for the presence of patterns, or eye movement sequences, that would indicate in what order people tend to construct a representation of recursive mental states. The most efficient strategy across all games would have been to construct these recursive mental states in a backward fashion, as each mental state depended on the next one. However, that strategy requires a deep understanding of the task domain, and most people tend to use more simple strategies that work across multiple domains (e.g., Gopnik et al., 2004; Todd & Gigerenzer, 2000). The eye movements indeed indicated that people inferred mental states in a more simple and forward progression, only tracking backward if previous (higher-order) mental states had to be revised.

The prevalence of a simple and forward approach can be explained by the principle of economy: Immediate decisions that ignore future possibilities *can* yield the optimal outcome in many cases (2008; Meijering et al., 2012; Szymanik et al., 2013). Furthermore, forward reasoning receives much practice across many domains, for example in causal inference (Apperly & Butterfill, 2009; Baron-Cohen, Leslie, & Frith, 1985; Gopnik et al., 2004; Gopnik & Wellman, 1992; Leslie, Friedman, & German, 2004; Onishi & Baillargeon, 2005; Premack & Woodruff, 1978; Saxe, Schulz, & Jiang, 2006; Wimmer & Perner, 1983). Thus, it is not surprising that forward reasoning, in its most simple form, is also used to construct representations of recursive mental states.

Chapter 6 provides a computational cognitive account of the use of simple strategies during inference of mental states. The model is based on previous empirical findings (reported in Chapters 2 and 5), and is validated by data from Flobbe et al.'s (but see Baker, Saxe, & Tenenbaum, 2009; Flobbe et al., 2008) developmental study. The model uses a simple strategy at first and only starts incorporating more complex mental states in the face of unexpected outcomes. The simple strategy is comparable to forward reasoning, which is later followed by backtracking, as the model starts considering future possible actions and underlying mental states. Investigating Flobbe et al.'s data, we saw response patterns that were indeed indicative of simple strategies. Few children were able to account for the fact that another person could be reasoning about them.

In sum, the findings in Chapters 5 and 6 corroborate the claim we made in earlier chapters: Application of (higher-order) ToM is a computational process that can either benefit from supporting structure, or be simplified by using simple strategies, thereby reducing cognitive demands. Given these findings, what new insights can future ToM research bring?

## Looking into the future

The studies in this dissertation show that application of (higher-order) ToM is a complicated task. The fact that higher-order ToM consists of multiple procedural and declarative building blocks almost poses something like an inverse problem in the sense that there are many possible sources that could yield the same behavioral patterns. Therefore, a multidisciplinary

approach to the investigation of ToM is desirable.

So far, ToM has mostly been studied in clinical settings, developmental studies, animal studies, and imaging studies, which all have produced many interesting insights and theories (Apperly & Butterfill, 2009; 2013; Baron-Cohen et al., 1985; Gopnik & Wellman, 1992; Leslie et al., 2004; Onishi & Baillargeon, 2005; Premack & Woodruff, 1978; Saxe et al., 2006; Wimmer & Perner, 1983). Most theories, however, exist only on paper, and do not directly translate to quantifiable predictions (but see Baker et al., 2009). This is why computational cognitive modeling needs to be employed more often, as a way of testing existing and new theories: Once implemented, a theory can yield quantifiable predictions that can be directly tested in one or more experiments.

This dissertation is a modest step towards a cognitive modeling approach, as we have yet to validate our model in many more domains. The model could be used to accommodate other ToM paradigms, and it could be used to generate hypotheses for clinical, developmental, imaging and other psychophysiological studies. The model could, for example, simulate application of ToM in higher-order false-belief tasks, in which one agent (e.g., Sally) has a false belief about another agent's (e.g., Anne's) beliefs. At first, the model would do the task from its own perspective, and later the model would attribute its knowledge to the other agent, Sally. Lastly, the model would reason about Sally as if she would attribute her own knowledge to yet again another agent, Anne. Scenarios such as these are currently being tested empirically by Arslan et al. (2013). The most obvious prediction is that children of 4 years old who have not yet mastered second-order ToM, but who do have first-order ToM, would not fall prey to the reality bias. Instead they would apply first-order ToM when, in actuality, they are asked to make a second-order inference.

It would also be interesting to look at transfer of ToM between various ToM domains by having one model play, for example, Marble Drop games and do higher-order false-belief tasks. Would experience in one task be beneficial in the other? How much overlap is there between the tasks with respect to demands on cognitive functions? Questions like these have recently been investigated by Taatgen (2013) in other (non-social) domains. His modeling approach has culminated in Actransfer, an extension of ACT-R (Anderson, 2007; Anderson et al., 2004), which is a theory about the nature and transfer of cognitive skills. Actransfer is particularly relevant to investigate the domain-specificity of ToM, which is still hotly debated: Does ToM require a specialized cognitive function, or does it involve general cognitive skills that are used across multiple domains? The empirical results in Chapter 3 suggest that ToM requires a specialized cognitive function, but it provides pointers to many possible explanations. Actransfer could help in finding the most primitive elements that comprise application of ToM.

Testing our model in various domains may not only yield insights that help us improve the model. Some insights will help improve the cognitive architecture (e.g., ACT-R) as a whole. If, for example, application of ToM indeed requires a specialized cognitive function to infer mental states of self and others, the architecture needs to accommodate such a module. Therefore, the study of ToM is particularly interesting for the entire cognitive sciences, as it will help constructing an integrated theory of cognition.





## Samenvatting

Het onderzoek in dit proefschrift gaat over hoe mensen redeneren over andermans denken; wat hun gedachten, intenties en doelen zijn. Dit onderzoek is belangrijk voor het vormen van een geïntegreerde theorie van cognitie, omdat de vraag is of redeneren over andermans denken een speciale cognitieve functie is. Zowel het onderzoek als het redeneren over andermans denken is ingewikkeld om één heel duidelijke reden: redeneren en denken zijn onzichtbaar processen. Toch zijn mensen verassend goed in het redeneren over andermans denken; we kunnen met enige zekerheid voorspellen op welke politieke partij vrienden en familieleden zullen stemmen, omdat we weten hoe zij denken, wat hun opvattingen zijn, et cetera. Maar er zijn limieten aan het redeneren over andermans denken. Met name bij jonge kinderen wordt dat snel duidelijk, bijvoorbeeld als ze niet begrijpen dat we ze nog steeds kunnen zien als ze hun handen voor hun ogen houden. Jonge kinderen vinden het moeilijk om zich in een ander te verplaatsen. Volwassenen hebben daar minder moeite mee, maar toch wordt tijdens een spel schaak (of poker) al snel duidelijk dat het ondoenlijk is om alle gedachten van de tegenspeler te anticiperen. Het is helemaal moeilijk om te anticiperen welke gedachten de tegenspeler heeft over onze gedachten. Dit recursieve denkproces, dat in theorie oneindig is, houdt in de praktijk al snel op. In dit proefschrift laat ik zien waarom het redeneren over andermans denken, wat ik vanaf nu meta-denken zal noemen, limieten heeft.

## Onderzoek naar meta-denken

Meta-denken is in het verleden voornamelijk bij kinderen onderzocht, met als achterliggend idee dat de ontwikkeling van een cognitieve functie iets kan vertellen over de uiteindelijke aard van die functie in volwassenen. Het meest populaire experiment om meta-denken bij kinderen te onderzoeken is de Sally-Anne taak. Het gaat als volgt: Sally en Anne spelen met knikkers. Als Sally eventjes weggaat, bergt ze eerst haar knikkers op in haar mandje. Terwijl Sally weg is, pakt Anne de knikkers en verstopt ze die in de speelgoeddoos. Na een tijdje komt Sally terug en de vraag aan kinderen is: Waar zal Sally naar haar knikkers zoeken? Kinderen die het meta-denken onder de knie hebben, weten dat Sally nog steeds denkt dat de knikkers in haar mandje zijn en dat ze daar zal zoeken, ook al zijn de knikkers in werkelijkheid in de speelgoeddoos. Kinderen die het meta-denken nog niet beheersen, zullen zeggen dat Sally naar de knikkers in de speelgoeddoos zal zoeken. Ze kunnen feiten en gedachten nog niet van elkaar onderscheiden.

Tot op heden is nog steeds niet duidelijk of dat komt omdat ze geen begrip hebben van mentale toestanden zoals kennis, gedachten en intenties of omdat ze nog niet de cognitieve vaardigheden hebben om over die mentale toestanden na te denken. Om dit probleem te ondervangen, heb ik onderzoek gedaan bij volwassenen, want die hebben een groter besef van onzichtbare denkprocessen. De algemene cognitieve vaardigheden zijn ook beter ontwikkeld bij volwassen en dus stelt onderzoek bij hen ons beter in staat om vast te stellen in hoeverre meta-denken een speciale cognitieve functie betreft.

Speciaal voor dit onderzoek heb ik een knikkerspel ontwikkeld waarin twee spelers om beurten een beslissing maken over het verloop van het spel. Het spel doet een beroep op het meta-denken, omdat de uitkomsten van elke beslissing afhankelijk zijn van de volgende beslissing. De ene speler moet dus nadenken over het beslisproces en de onderliggende

gedachten van de andere speler.

## De aard van het meta-denken

Om na te gaan of het meta-denken wordt gelimiteerd door cognitieve vaardigheden, heb ik onderzocht of training en andersoortige ondersteuning een gunstig effect hebben op het meta-denken. De uitkomsten van dit onderzoek zijn positief: Mensen presteerden beter als ze het meta-denken in een bepaalde taak stapsgewijs oefenden: Ze leerden eerst de taak vanuit het eigen perspectief te doen en later dat perspectief aan een ander toe te schrijven. Ook presteerden mensen beter als hen expliciet werd gevraagd om zich te verplaatsen in het perspectief van de ander. Deze resultaten laten zien dat suboptimale uitkomsten in een sociale interactie niet zozeer zijn toe te schrijven aan onbegrip als wel aan gebrek aan oefening. In die zin is het meta-denken dus een cognitieve vaardigheid die je kunt oefenen en niet perse een vaardigheid die men wel of niet beheerst.

De belangrijkste vraag in dit proefschrift is of meta-denken een speciale cognitieve vaardigheid is, in tegenstelling tot een vaardigheid die is opgebouwd uit meer algemene vaardigheden. De resultaten van een aantal ontwikkelingsstudies in 2006 doen vermoeden dat meta-denken niet speciaal is en bestaat uit algemene vaardigheden. Deze studies hebben laten zien dat kinderen die slecht scoorden op de Sally-Anne taak ook moeite hadden met redeneren in het algemeen. Echter, in dit proefschrift laat ik zien dat dergelijke conclusies niet zijn te trekken op basis van ontwikkelingsstudies. Het experiment in Hoofdstuk 3 laat duidelijk zien dat het meta-denken wel degelijk een speciale cognitieve vaardigheid is: volwassenen werden blootgesteld aan twee condities met als enige verschil de instructie om de taak vanuit het eigen perspectief of dat van een ander te doen. De resultaten laten zien dat zowel de kwaliteit als de snelheid van het meta-denken verschilde tussen de twee condities. Als de instructie was om de taak vanuit het eigen perspectief te doen, presteerden de volwassen sneller en beter dan wanneer de instructie was om dezelfde taak vanuit het perspectief van een ander te doen. Deze resultaten impliceren dat de aard van het menselijk redeneren afhankelijk is van het feit of een taak al dan niet over mentale toestanden gaat, los van de complexiteit van de taak.

In dit proefschrift laat ik ook zien dat in sociale interacties mensen de voorkeur geven aan simpele strategieën en pas beginnen met meta-denken als het echt niet anders kan. In een eerdere ontwikkelingsstudie werd kinderen gevraagd een computerspel te spelen, met de computer als tegenstander. De uitkomsten van die studie deden vermoeden dat de kinderen het meta-denken al enigszins beheersten. Maar door middel van computersimulaties laat ik zien dat het waarschijnlijker is dat de kinderen simpelere strategieën gebruikten als die ook tot juiste uitkomst leidden. Pas als duidelijk werd dat die strategieën niet altijd werkten, probeerden de kinderen het perspectief van de tegenspeler mee te nemen door hun eigen strategie aan de tegenspeler toe te schrijven. Zodoende ontwikkelden zij langzamerhand steeds complexere strategieën die steeds meer op meta-denken gingen lijken.

De bevindingen van dit proefschrift zijn belangrijk voor het formuleren van een geïntegreerde theorie van cognitie, omdat het meta-denken daarin een rol dient te hebben. Het is een speciale cognitieve functie die het denken in kwalitatieve zin beïnvloed. Daarnaast zijn de bevindingen in dit proefschrift van belang voor de praktijk, onder andere omdat het meta-denken is te trainen. Door het bieden van de juiste structuur kunnen mensen leren om optimalere uitkomsten te behalen in sociale interacties zoals onderhandelingen,

samenwerkingsverbanden en competities.

## Dankwoord

Toen ik pas begonnen was als PhD student en ik Rineke nog niet zo goed kende, vertelden andere PhD studenten me: 'Oh, she is so great!' en 'You are so lucky!' Inmiddels kan ik bevestigen dat dat waar is. Ik heb geluk gehad: behalve dat Rineke een inspirerende onderzoeker is, stond ze ook nog eens altijd voor me klaar.

Ik heb ook geluk gehad met mijn andere begeleiders, Hedderik en Niels. Van hen heb ik geleerd dat 'onderzoek doen' net zo gaaf kan zijn als detectivewerk.

In de tijd bij AI maakte ik deel uit van twee onderzoeksgroepen, de Vici Reading Group en de Cognitive Modeling Group. Ik wil huidige en voormalige deelnemers van beide groepen bedanken. De discussies in de Vici Reading Group waren bijzonder vermakelijk; wat konden we kritisch zijn op andermans papers :). Het was ook erg leerzaam om elkaars werk en presentaties te zien, in beide groepen.

Louter door organisatorische invloeden heb ik diverse kamergenoten gehad: Jelmer, Leendert, Enkhbold, Sujata, Jakob en Burcu. Ik wil Jelmer en Leendert bedanken voor hun inzichtelijke kijk op het wetenschappelijk spel; Enkhbold voor z'n gedeelde fascinatie voor (computer)spellen; Sujata en Jakob voor hun mooie verhalen over logica en logici; Burcu voor haar uitstekende richtingsgevoel in Berlijn.

Tussen de onderzoeksbedrijven door was er gelukkig vaak tijd voor ontspanning en afleiding, bijvoorbeeld tijdens koffiepauzes, waar we voor menig niet-bestaand probleem een oplossing hebben gevonden. Aswin, Bea, Burcu, Charlotte, Gyuhee, Harmen, Ioanna, Jean-Paul, Marijke, Michiel, Olarik en Trudy waren vaak van de partij, bedankt!

Hereby, I would like to express my gratitude to professors Petra Hendriks, Josef Perner, and Maartje Raijmakers for reading my dissertation. I'm proud that I can say that these great researchers have seen my work.

Ik wil ook graag mijn paranimfen bedanken, Jean-Paul en Durk. Ik ben blij dat twee goeie vrienden me bijstaan in het ceremoniële geweld van de verdediging. Beide beschikken namelijk ruimschoots over 'theory of mind'.

Tenslotte, pa, ma, Rob, Gonny, Charlotte: ik heb het gehaald!

## References

- Anderson, J. R. (2007). *How can the human mind occur in the physical universe?* New York, NY: Oxford University Press.
- Anderson, J. R., Bothell, D., Byrne, M. D., Douglass, S., Lebiere, C., & Qin, Y. (2004). An integrated theory of the mind. *Psychological Review*, *111*(4), 1036–1060.
- Apperly, I. A. (2011). *Mindreaders: The cognitive basis of “theory of mind.”* Hove, UK: Psychology Press.
- Apperly, I. A., & Butterfill, S. A. (2009). Do humans have two systems to track beliefs and belief-like states? *Psychological Review*, *116*(4), 953–970.
- Apperly, I. A., Back, E., Samson, D., & France, L. (2008). The cost of thinking about false beliefs: Evidence from adults’ performance on a non-inferential theory of mind task. *Cognition*, *106*(3), 1093–1108.
- Apperly, I. A., Carroll, D. J., Samson, D., Humphreys, G. W., Qureshi, A., & Moffit, G. (2010). Why are there limits on theory of mind use? Evidence from adults’ ability to follow instructions from an ignorant speaker. *The Quarterly Journal of Experimental Psychology*, *63*(6), 1201–1217.
- Arslan, B., Hohenberger, A., & Verbrugge, R. (2012). The development of second-order social cognition and its relation with complex language understanding and memory. In N. Miyake, D. Peebles, & R. P. Cooper (Eds.), *Proceedings of the 34th Annual Conference of the Cognitive Science Society* (pp. 1290–1295). Austin, TX: Cognitive Science Society.
- Arslan, B., Taatgen, N. A., & Verbrugge, R. (2013). Modeling developmental transitions in reasoning about false beliefs of others. In R. L. West, & T. C. Stewart (Eds.), *Proceedings of the 12th International Conference on Cognitive Modelling* (pp. 77–82). Ottawa, CA: Carleton University.
- Baayen, R. H. (2008). *Analyzing linguistic data: A practical introduction to statistics using R.* New York, USA: Cambridge University Press.
- Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed random effects for subjects and items. *Journal of Memory and Language*, *59*(4), 390–412.
- Bacharach, M., & Stahl, D. O. (2000). Variable-frame level-n theory. *Games and Economic Behavior*, *32*(2), 220–246.
- Baillargeon, R., Scott, R. M., & He, Z. (2010). False-belief understanding in infants. *Trends in Cognitive Sciences*, *14*(3), 110–118.
- Baker, C., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition*, *113*(3), 329–349.
- Baron-Cohen, S., Leslie, A. M., & Frith, U. (1985). Does the autistic child have a “theory of mind?” *Cognition*, *21*(1), 37–46.
- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, *68*(3), 255–278.
- Birch, S. A. J., & Bloom, P. (2004). Understanding children’s and adults’ limitations in mental state reasoning. *Trends in Cognitive Sciences*, *8*(6), 255–260.
- Birch, S. A. J., & Bloom, P. (2007). The curse of knowledge in reasoning about false beliefs. *Psychological Science*, *18*(5), 382–386.
- Bloom, P., & German, T. P. (2000). Two reasons to abandon the false belief task as a test of theory of mind. *Cognition*, *77*(1), 25–31.

- Borst, J. P., Taatgen, N. A., & Van Rijn, H. (2010a). The problem state: A cognitive bottleneck in multitasking. *Journal of Experimental Psychology: Learning, memory, and cognition*, 36(2), 363–382.
- Borst, J. P., Taatgen, N. A., Stocco, A., & Van Rijn, H. (2010b). The neural correlates of problem states: Testing fMRI predictions of a computational model of multitasking. *PLoS ONE*, 5(9).
- Brassard, G., & Bratley, P. (1996). *Fundamentals of algorithmics*. Englewood Cliffs, N.J.: Prentice Hall.
- Bull, R., Phillips, L. H., & Conway, C. A. (2008). The role of control functions in mentalizing: Dual-task studies of theory of mind and executive function. *Cognition*, 107(2), 663–672.
- Call, J., & Tomasello, M. (2008). Does the chimpanzee have a theory of mind? 30 years later. *Trends in Cognitive Sciences*, 12(5), 187–192.
- Carlson, S. M., Moses, L. J., & Breton, C. (2002). How specific is the relation between executive function and theory of mind? Contributions of inhibitory control and working memory. *Infant and Child Development*, 11(2), 73–92.
- Clements, W. A., & Perner, J. (1994). Implicit understanding of belief. *Cognitive Development*, 9(4), 377–395.
- Cohen, A. S., & German, T. P. (2010). A reaction time advantage for calculating beliefs over public representations signals domain specificity for “theory of mind.” *Cognition*, 115(3), 417–425.
- Colman, A. M. (2003). Depth of strategic reasoning in games. *Trends in Cognitive Sciences*, 7(1), 2–4.
- Coull, G. J., Leekam, S. R., & Bennett, M. (2006). Simplifying second-order belief attribution: What facilitates children’s performance on measures of conceptual understanding? *Social Development*, 15(2), 260–275.
- De Villiers, J. (2007). The interface of language and theory of mind. *Lingua*, 117(11), 1858–1878.
- De Villiers, J. G., & Pyers, J. E. (2002). Complements to cognition: A longitudinal study of the relationship between complex syntax and false-belief-understanding. *Cognitive Development*, 17(1), 1037–1060.
- De Weerd, H., Verbrugge, R., & Verheij, B. (2013). How much does it help to know what she knows you know? An agent-based simulation study. *Artificial Intelligence*, 199–200, 67–92.
- Dumontheil, I., Apperly, I. A., & Blakemore, S.-J. (2010). Online usage of theory of mind continues to develop in late adolescence. *Developmental Science*, 13(2), 331–338.
- 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.
- Fu, W., & Gray, W. D. (2004). Resolving the paradox of the active user: Stable suboptimal performance in interactive tasks. *Cognitive Science*, 28(1), 901–935.
- Gallagher, H. L., & Frith, C. D. (2003). Functional imaging of “theory of mind.” *Trends in Cognitive Sciences*, 7(2), 77–83.
- Gelman, A., & Hill, J. (2007). *Data analysis using regression and multilevel/hierarchical models*. New York: Cambridge University Press.

- 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.
- Ghosh, S., & Meijering, B. (2011). On combining cognitive and formal modeling: A case study involving strategic reasoning. In J. Van Eijck & R. Verbrugge (Eds.), *Proceedings of the Workshop on Reasoning About Other Minds: Logical and Cognitive Perspectives* (pp. 79–92). Groningen, NL: University of Groningen.
- Ghosh, S., Meijering, B., & Verbrugge, R. (2010). Empirical reasoning in games: Logic meets cognition. In T. Ågotnes, N. Alechina, & B. Logan (Eds.), *Proceedings Third Logics for Resource Bounded Agents Workshop* (pp. 15–34). Lyon, FR.
- 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.
- Gopnik, A., & Slaughter, V. (1991). Young children’s understanding of changes in their mental states. *Child Development*, 62(1), 98–110.
- Gopnik, A., & Wellman, H. M. (1992). Why the child’s theory of mind really is a theory. *Mind and Language*, 7(1-2), 145–171.
- Gopnik, A., Glymour, C., Sobel, D. M., Schulz, L. E., Kushnir, T., & Danks, D. (2004). A theory of causal learning in children: Causal maps and Bayes nets. *Psychological Review*, 111(1), 3–32.
- Gray, W. D., Sims, C. R., Fu, W., & Schoelles, M. J. (2006). The soft constraints hypothesis: A rational analysis approach to resource allocation for interactive behavior. *Psychological Review*, 113(3), 461–482.
- Hale, C. M., & Tager-Flusberg, H. (2003). The influence of language on theory of mind: A training study. *Developmental Science*, 6(3), 346–359.
- Hedden, T., & Zhang, J. (2002). What do you think I think you think?: Strategic reasoning in matrix games. *Cognition*, 85(1), 1–36.
- Hiatt, L. M., & Trafton, J. G. (2010). A cognitive model of theory of mind. In D. Salvucci & G. Gunzelmann (Eds.), *Proceedings of the 10th International Conference on Cognitive Modeling* (pp. 91–96). Philadelphia, PA: Drexel University.
- Hollebrandse, B., Hobbs, K., De Villiers, J., & Roeper, T. (2008). Second order embedding and second order false belief. In A. Gavarro & M. J. Freitas (Eds.), *Language Acquisition and Development: Proceedings of GALA 2007*. Cambridge Scholar Press.
- Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6, 65–70.
- Jansen, B. R. J., & Van der Maas, H. L. J. (2002). The development of children’s rule use on the balance scale task. *Journal of Experimental Child Psychology*, 81(4), 383–416.
- Johnson, E. J., Camerer, C., Sen, S., & Rymon, T. (2002). Detecting failures of backward induction: Monitoring information search in sequential bargaining. *Journal of Economic Theory*, 104(1), 16–47.
- Johnson-Laird, P. N. (1983). *Mental models: Towards a cognitive science of language, inference, and consciousness*. Cambridge, MA: Harvard University Press.
- Kaufman, L., & Rousseeuw, P. J. (1990). *Finding groups in data: An introduction to cluster analysis*. Hoboken, N.J.: John Wiley & Sons, Inc.
- Kelley, H. H., & Michela, J. L. (1980). Attribution theory and research. *Annual Review of Psychology*, 31(1), 457–501.



- Keysar, B., Lin, S., & Barr, D. J. (2003). Limits on theory of mind use in adults. *Cognition*, 89(1), 25–41.
- Kong, X., Schunn, C. D., & Wallstrom, G. L. (2010). High regularities in eye movement patterns reveal the dynamics of the visual working memory allocation mechanism. *Cognitive Science*, 34(2), 322–337.
- Kovács, Á. M., Téglás, E., & Endress, A. D. (2010). The social sense: Susceptibility to others' beliefs in human infants and adults. *Science*, 330(6012), 1830–1834.
- Leekam, S. R., Perner, J., Healey, L., & Sewell, C. (2006). False signs and the non-specificity of theory of mind: Evidence that preschoolers have general difficulties in understanding representations. *British Journal of Developmental Psychology*, 26(4), 485–497.
- Leslie, A. M., Friedman, O., & German, T. P. (2004). Core mechanisms in “theory of mind.” *Trends in Cognitive Sciences*, 8(12), 528–533.
- 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.
- Liversedge, S. P., & Findlay, J. M. (2000). Saccadic eye movements and cognition. *Trends in Cognitive Sciences*, 4(1), 6–14.
- McColgan, K. L., & McCormack, T. (2008). Searching and planning: Young children's reasoning about past and future event sequences. *Child Development*, 79(5), 1477–1497.
- McKelvey, R. D., & Palfrey, T. R. (1992). An experimental study of the centipede game. *Econometrica*, 60(4), 803–836.
- Meijering, B., Van Maanen, L., Van Rijn, H., & Verbrugge, R. (2010). The facilitative effect of context on second-order social reasoning. In R. Catrambone & S. Ohlsson (Eds.), *Proceedings of the 32nd Annual Meeting of the Cognitive Science Society* (pp. 1423–1428). Austin, TX: Cognitive Science Society.
- Meijering, B., Van Rijn, H., Taatgen, N. A., & Verbrugge, R. (2012). What eye movements can tell about theory of mind in a strategic game. *PLoS ONE*, 7(9).
- Meijering, B., Van Rijn, H., Taatgen, N. A., & Verbrugge, R. (2013). Reasoning about diamonds, gravity, and mental states: The cognitive costs of theory of mind. In M. Knauff, M. Pauen, N. Sebanz, & I. Wachsmuth (Eds.), *Proceedings of the 35th Annual Meeting of the Cognitive Science Society* (pp. 3026–3031). Austin, TX: Cognitive Science Society.
- Meijering, B., Van Rijn, H., Taatgen, N., & 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). Austin, TX: Cognitive Science Society.
- Miller, S. A. (2009). Children's understanding of second-order mental states. *Psychological Bulletin*, 135(5), 749–773.
- Nyamsuren, E., & Taatgen, N. A. (2013). Set as an instance of a real-world visual-cognitive task. *Cognitive Science*, 37(1), 146–175.
- O'Neill, D. K. (1996). Two-year-old children's sensitivity to a parent's knowledge state when making requests. *Child Development*, 67(2), 659–677.
- Onishi, K. H., & Baillargeon, R. (2005). Do 15-month-old infants understand false beliefs? *Science*, 308(5719), 255–258.
- Osborne, M. J., & Rubinstein, A. (1994). *A course in game theory*. Cambridge, MA: MIT press.

- Perner, J., & Leekam, S. R. (2008). The curious incident of the photo that was accused of being false: Issues of domain specificity in development, autism, and brain imaging. *The Quarterly Journal of Experimental Psychology*, *61*(1), 76–89.
- 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.
- Pinheiro, J. C., & Bates, D. M. (2000). Mixed-effects models in S and S-PLUS. Springer Verlag.
- Premack, D., & Woodruff, G. (1978). Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, *1*(4), 515–526.
- Qureshi, A. W., Apperly, I. A., & Samson, D. (2010). Executive function is necessary for perspective selection, not Level-1 visual perspective calculation: Evidence from a dual-task study of adults. *Cognition*, *117*(2), 230–236.
- Raijmakers, M. E. J., Mandell, D. J., Van Es, S. E., & Counihan, M. (2013). Children’s strategy use when playing strategic games. *Synthese*.
- Ramsey, R., Hansen, P. C., Apperly, I. A., & Samson, D. (2013). Seeing it my way or your way: Frontoparietal brain areas sustain viewpoint-independent perspective selection processes. *Journal of Cognitive Neuroscience*, *25*(5), 670–684.
- Rayner, K. (1998). Eye movements in reading and information processing: 20 years of research. *Psychological Bulletin*, *124*(3), 372–422.
- Ruffman, T., Garnham, W., Import, A., & Connolly, D. (2001). Does eye gaze indicate implicit knowledge of false belief? Charting transitions in knowledge. *Journal of Experimental Child Psychology*, *80*(3), 201–224.
- Sabbagh, M. A., Xu, F., Carlson, S. M., Moses, L. J., & Kang, L. (2006). The development of executive functioning and theory of mind: A comparison of Chinese and U.S. preschoolers. *Psychological Science*, *17*(1), 74–81.
- Samson, D., Apperly, I. A., Braithwaite, J. J., Andrews, B. J., & Bodley Scott, S. E. (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.
- Samson, D., Apperly, I. A., Kathirgamanathan, U., & Humphreys, G. W. (2005). Seeing it my way: A case of a selective deficit in inhibiting self-perspective. *Brain*, *128*(5), 1102–1111.
- Saxe, R. (2006). Uniquely human social cognition. *Current Opinion in Neurobiology*, *16*(2), 235–239.
- 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.
- Suddendorf, T., Nielsen, M., & Gehlen, von, R. (2011). Children’s capacity to remember a novel problem and to secure its future solution. *Developmental Science*, *14*(1), 26–33.
- Sullivan, K., Zaitchik, D., & Tager-Flusberg, H. (1994). Preschoolers can attribute second-order beliefs. *Developmental Psychology*, *30*(3), 395–402.
- Szymanik, J., Meijering, B., & Verbrugge, R. (2013). Using intrinsic complexity of turn-taking games to predict participants’ reaction times. In M. Knauff, M. Pauen, N. Sebanz, & I. Wachsmuth (Eds.), *Proceedings of the 35th Annual Conference of the Cognitive Science Society* (pp. 1426–1432). Austin, TX: Cognitive Science Society.
- Taatgen, N. A. (2013). The nature and transfer of cognitive skills. *Psychological Review*, *120*(3), 439–471.

- Thoermer, C., Sodian, B., Vuori, M., Perst, H., & Kristen, S. (2012). Continuity from an implicit to an explicit understanding of false belief from infancy to preschool age. *British Journal of Developmental Psychology, 30*(1), 172–187.
- Todd, P. M., & Gigerenzer, G. (2000). Précis of simple heuristics that make us smart. *Behavioral and Brain Sciences, 23*(1), 727–780.
- Van Maanen, L., & Verbrugge, R. (2010). A computational model of second-order social reasoning. In D. Salvucci & G. Gunzelmann (Eds.), *Proceedings of the 10th International Conference on Cognitive Modeling* (pp. 259–264). Philadelphia, PA: Drexel University.
- Van Rij, J., Van Rijn, H., & Hendriks, P. (2010). Cognitive architectures and language acquisition: A case study in pronoun comprehension. *Journal of Child Language, 37*(3), 731–766.
- Van Rij, J., Van Rijn, H., & Hendriks, P. (2013). How WM load influences linguistic processing in adults: A computational model of pronoun interpretation in discourse. *Topics in Cognitive Science, 5*(3), 564–580.
- Van Rijn, H., Van Someren, M., & Van der Maas, H. L. J. (2003). Modeling developmental transitions on the balance scale task. *Cognitive Science, 27*(2), 227–257.
- Verbrugge, R., & Mol, L. (2008). Learning to apply theory of mind. *Journal of Logic, Language and Information, 17*(4), 489–511.
- Weber, R., Camerer, C., Rottenstreich, Y., & Knez, M. (2001). The illusion of leadership: Misattribution of cause in coordination games. *Organization Science, 12*(5), 582–598.
- Wellman, H. M., Cross, D., & Watson, J. (2001). Meta-analysis of theory-of-mind development: The truth about false belief. *Child Development, 72*(3), 655–684.
- Wimmer, H., & Perner, J. (1983). Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception. *Cognition, 13*(1), 103–128.
- Winner, E., & Leekam, S. R. (1991). Distinguishing irony from deception: Understanding the speaker's second-order intention. *British Journal of Developmental Psychology, 9*(2), 257–270.
- Zhang, J., & Hedden, T. (2003). Two paradigms for depth of strategic reasoning in games: Response to Colman. *Trends in Cognitive Sciences, 7*(1), 4–5.
- 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.

## Publication list

### Journal publications

- Meijering, B.**, Van Rijn, H., Taatgen, N. A., & Verbrugge, R. (2012). What eye movements can tell about theory of mind in a strategic game. *PLoS ONE*, 7(9). Doi:10.1371/journal.pone.0045961
- Ghosh, S., **Meijering, B.**, & Verbrugge, R. (2014). Strategic reasoning: Building cognitive models from logical formulas. *Journal of Logic, Language, and Information*. Doi:10.1007/s10849-014-9196-x
- Szymanik, J., Robaldo, L., & **Meijering, B.** (2014). On the identification of quantifiers' witness sets: A study of multi-quantifier sentences. *Journal of Logic, Language, and Information*. Doi:10.1007/s10849-014-9197-9

### Peer-reviewed conference papers

- Bergwerff, G., **Meijering, B.**, Szymanik, J., Verbrugge, R., & Wierda, S. M. (2014). Computational and algorithmic models of strategies in turn-based games. In *Proceedings of the 36th Annual Conference of the Cognitive Science Society*. Austin, TX: Cognitive Science Society
- Meijering, B.**, Taatgen, N. A., Van Rijn, H., & Verbrugge, R. (2013). Reasoning about mental states in sequential games: As simple as possible, as complex as necessary. In R. L. West & T. C. Stewart (Eds.), *Proceedings of the 12th International Conference on Cognitive Modeling* (pp. 173–178). Ottawa: Carleton University.
- Meijering, B.**, Van Rijn, H., Taatgen, N. A., & Verbrugge, R. (2013). Reasoning about diamonds, gravity, and mental states: The cognitive costs of theory of mind. In M. Knauff, M. Pauen, N. Sebanz, & I. Wachsmuth (Eds.), *Proceedings of the 35th Annual Meeting of the Cognitive Science Society* (pp. 3026–3031). Austin, TX: Cognitive Science Society.
- Szymanik, J., **Meijering, B.**, & Verbrugge, R. (2013). Using intrinsic complexity of turn-taking games to predict participants' reaction times. In M. Knauff, M. Pauen, N. Sebanz, & I. Wachsmuth (Eds.), *Proceedings of the 35th Annual Conference of the Cognitive Science Society* (pp. 1426–1432). Austin, TX: Cognitive Science Society.
- Meijering, B.**, Van Rijn, H., Taatgen, N., & 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). Austin, TX: Cognitive Science Society.
- Ghosh, S., & **Meijering, B.** (2011). On combining cognitive and formal modeling: A case study involving strategic reasoning. In J. van Eijck and R. Verbrugge (eds.), *Proceedings of the Workshop on Reasoning About Other Minds: Logical and Cognitive Perspectives* (pp. 79–92). Groningen, 2011.
- Meijering, B.**, Van Maanen, L., Van Rijn, H., & Verbrugge, R. (2010). The facilitative effect of context on second-order social reasoning. In R. Catrambone & S. Ohlsson (Eds.), *Proceedings of the 32nd Annual Meeting of the Cognitive Science Society* (pp. 1423–1428). Austin, TX: Cognitive Science Society.

Ghosh, S., **Meijering, B.**, & Verbrugge, R. (2010). Empirical reasoning in games: Logic meets cognition. In T. Ågotnes, N. Alechina, & B. Logan (Eds.), *Proceedings Third Logics for Resource Bounded Agents Workshop* (pp. 15–34). Lyon, FR.

## Abstract

**Meijering, B.**, Van Rijn, H., Taatgen, N. A., & Verbrugge, R. (2011). Second-order theory of mind in strategic games. *Interdisciplinary Workshop on Cognitive Neuroscience, Educational Research, and Cognitive Modeling*. Delmenhorst, Germany.

## In preparation

**Meijering, B.**, Van Rijn, H., Taatgen, N. A., & Verbrugge, R. (Submitted). *Reasoning about self versus others: Changing perspective is hard*.

**Meijering, B.**, Van Rijn, H., Taatgen, N. A., & Verbrugge, R. (Under revision). *Integrating recursive theory of mind in decision making in sequential games*.

**Meijering, B.**, Taatgen, N. A., & Verbrugge, R. (Under revision). *Modeling inference of mental states: As simple as possible, as complex as necessary*.





