

# Negotiating with Incomplete Information: The Influence of Theory of Mind

Eveline Broers  
August 2014

Master's Thesis  
Human-Machine Communication  
Department of Artificial Intelligence  
University of Groningen, The Netherlands

First supervisor and reviewer:

Prof. Dr. L.C. Verbrugge (Artificial Intelligence, University of Groningen)

Second supervisor:

H.A. de Weerd (Artificial Intelligence, University of Groningen)

Second reviewer:

Prof. Dr. N.A. Taatgen (Artificial Intelligence, University of Groningen)



**university of  
 groningen**

**faculty of mathematics  
 and natural sciences**



# Abstract

The aim of this master's thesis was to investigate the reasoning behavior of people during negotiations with incomplete information. The question was whether people reason about the knowledge, intentions and beliefs of others in a negotiation setting with incomplete information; do they use so called 'theory of mind'? Participants played the negotiation game colored trails (for which the use of theory of mind has proven to be useful) against three types of computer agents, who all used a different order of theory of mind (zero, first or second). The negotiations were about the distribution of some resources of which a subset was needed to get to a certain goal location. The goal location of the computer agent was not public knowledge, which invited the participants to reason about the actions and possible goal location of the computer agent.

The results showed that people reasoned about the offers of the computer agent. They mainly used first-order theory of mind (reasoning about someone's mental states) and second-order theory of mind (reasoning about what ideas someone else has about someone's mental states). The scores of the participants were influenced by the order of theory of mind their opponent used. The participants also used more second-order theory of mind when the opponent used second-order theory of mind. Another outcome was that the participants mainly achieved higher scores when the opponent started the negotiation.

Furthermore, it was tested whether a training effect would occur when the participants first played marble drop: a game in which the actions of the opponent (opening a left or right trapdoor) need to be anticipated to get a marble at a preferred location. Unfortunately, no training effect was found, which might be due to the low accuracy of the participants on marble drop or because the game was not similar enough with colored trails.

Finally, it was investigated whether personality traits regarding empathy would influence a participant's results. People who reported that they tended to take into account the perspective of another person in daily life, did not perform better than others.



# Acknowledgements

First of all I would like to thank Rineke Verbrugge and Harmen de Weerd for the valuable meetings we had and their great ideas and support; they always encouraged me to continue. The assistance provided by Harmen de Weerd on the work with the computer agents was greatly appreciated as well. Furthermore I would like to thank Niels Taatgen for his useful comments.

Finally I would like to thank my fellow students with whom I could discuss my project. Special thanks go out to my family and to Ivo Brill for their support and for their help throughout the whole project concerning the content.



# Contents

<b>1</b>	<b>Introduction</b>	<b>9</b>
1.1	Research Questions . . . . .	10
1.2	Thesis Structure . . . . .	11
<b>2</b>	<b>Literature</b>	<b>13</b>
2.1	Theory of Mind . . . . .	13
2.2	Negotiations . . . . .	13
2.3	Colored Trails . . . . .	15
2.4	Training with Marble Drop . . . . .	17
2.5	Interpersonal Reactivity Index . . . . .	18
2.6	Basis for the Current Study . . . . .	19
2.6.1	Three Negotiating Agents with Complete Information . . . . .	19
2.6.2	Two Negotiating Agents with Incomplete Information . . . . .	20
2.6.3	Current Study . . . . .	21
<b>3</b>	<b>Method</b>	<b>23</b>
3.1	Introduction . . . . .	23
3.1.1	Procedure . . . . .	23
3.2	Color Test . . . . .	24
3.2.1	Materials . . . . .	24
3.2.2	Procedure . . . . .	24
3.3	Marble Drop . . . . .	24
3.3.1	Materials . . . . .	24
3.3.2	Procedure . . . . .	25
3.4	Colored Trails . . . . .	26
3.4.1	Rules . . . . .	26
3.4.2	Materials . . . . .	27
3.4.3	Procedure . . . . .	29
3.5	Questionnaires . . . . .	32
3.5.1	Materials . . . . .	32
3.5.2	Procedure . . . . .	33
3.6	Data Analysis . . . . .	33
<b>4</b>	<b>Results</b>	<b>35</b>
4.1	Participants . . . . .	35
4.2	Marble Drop . . . . .	35

4.2.1	Accuracy . . . . .	35
4.2.2	Questionnaire: the Participant's Experience . . . . .	38
4.3	Colored Trails . . . . .	39
4.3.1	Scores . . . . .	39
4.3.2	Used Orders of Theory of Mind by the Participant . . . . .	44
4.3.3	Questionnaire: the Participant's Experience . . . . .	49
4.3.4	Influence of Marble Drop . . . . .	52
4.4	Interpersonal Reactivity Index . . . . .	53
<b>5</b>	<b>Discussion</b>	<b>55</b>
5.1	Marble Drop . . . . .	55
5.2	Training . . . . .	56
5.3	Colored Trails . . . . .	57
5.3.1	Scores . . . . .	57
5.3.2	Used Orders of Theory of Mind by the Participant . . . . .	60
5.3.3	Comparison with Previous Research . . . . .	61
5.4	Interpersonal Reactivity Index . . . . .	61
5.5	Research Questions . . . . .	62
5.5.1	Subquestion 1: The Influence of the Opponent's Order of Theory of Mind . . . . .	62
5.5.2	Subquestion 2: The Influence of Training . . . . .	62
5.5.3	Subquestion 3: The Influence of Personality Traits regarding Empathy . . . . .	62
5.5.4	Main Question: Participant's Use of Theory of Mind . . . . .	63
5.6	Future Research . . . . .	63
5.6.1	Training or Transfer . . . . .	63
5.6.2	Adjusting the Colored Trails Set-up . . . . .	65
5.6.3	Colored Trails as Negotiation Practice . . . . .	65
<b>6</b>	<b>Bibliography</b>	<b>67</b>
<b>A</b>	<b>Marble Drop</b>	<b>73</b>
A.1	Colors . . . . .	73
A.2	Structures . . . . .	73
<b>B</b>	<b>Colored Trails</b>	<b>77</b>
B.1	Scenarios . . . . .	77
B.2	Control questions . . . . .	80
<b>C</b>	<b>Questionnaires: questions</b>	<b>83</b>
C.1	Questions about Marble Drop . . . . .	83
C.2	Questions about Colored Trails . . . . .	83
C.3	Interpersonal Reactivity Index . . . . .	84
<b>D</b>	<b>Questionnaires: participants' answers</b>	<b>87</b>
D.1	Colored Trails: Strategy Change by Opponent . . . . .	87

# Chapter 1

## Introduction

The world is a complex system. To function successfully in this world, it is often necessary to anticipate the actions of others. For example, when people play games, they often try to figure out the plans of their opponents in order to outsmart them. While playing a game is a setting in which people very consciously do this, they also use the so-called *theory of mind* in everyday life. Does he believe my story? Does the driver of that red car intend on stopping for me? Would my friend want a cat for her birthday? Trying to understand others, anticipating their actions and guessing what someone desires are all examples in which theory of mind can be applied.

In short, theory of mind is the ability to attribute mental states to others. There are different orders of theory of mind. Some examples to clarify each of them follow below (colors are from the perspective of the narrator):

<b>Zero-order</b>	<b>I know that</b> the money is in the drawer.
<i>First-order</i>	<b>I believe that</b> <i>James knows that</i> the money is in the drawer.
<i>Second-order</i>	<b>I believe that</b> <i>James hopes that Sandy does not know that</i> the money is in the drawer.
<i>n<sup>th</sup>-order</i>	etcetera

In the case of zero-order theory of mind, one does not ascribe knowledge and desires to other people (or animals). One only takes facts into account, like seeing that a coat is on a chair or knowing that someone just played a 3 of hearts in a card game. With first-order theory of mind, one can reason about the intentions and knowledge of another person. With second-order theory of mind, one can also reason about what ideas the other person has about one's own or someone else's thoughts.

In this master's thesis, we want to find out more about how people use theory of mind. Since this is a very broad statement, the current study will focus on a smaller question. It will investigate theory of mind in only one setting: negotiations. This is a setting in which using theory of mind is very important and it is something one does every day, often without realizing it. Negotiating is about reaching an agreement, so deciding in a social setting at what time to start a meeting or where to go for dinner are essentially negotiations.

In previous studies (see Section 2.6), this subject has been studied with computer agent simulations. In this study, this will be taken to the next level: humans will negotiate with computer agents. The question is if and how the participants make use of theory of mind. Another aim of this thesis is to investigate whether training to use theory of mind in a different setting influences the use of theory of mind in the negotiation setting.

## 1.1 Research Questions

The main research question this thesis will try to answer is:

*Do people use theory of mind when playing a negotiation game, for which the use of second order theory of mind has proven to be useful, against a computer agent and if so, what order of theory of mind do they use?*

To answer this question, participants have to play a negotiation game called colored trails, in which they need to negotiate about the distribution of some chips, against a computer agent. The agent will compute which order of theory of mind the participant is most likely using.

The use of theory of mind can be influenced by several things, of which the behavior of both the participant and the opponent is an important one. Therefore, the following sub-question is formulated:

*How is the use of theory of mind in the colored trails negotiation game influenced by the order of theory of mind the opponent uses?*

In order to answer this question, the participants have to play the negotiation game with different types of agents.

Training could also influence the use of theory of mind. This hypothesis is formulated as follows:

*What is the influence of training with the marble drop game on the use of theory of mind in the colored trails negotiation game?*

To be able to answer this question, half of the participants will be trained by playing marble drop, a game in which it is quite obvious to use theory of mind, and the other half will be the control group.

Personality traits regarding empathy, e.g. taking into account someone else's perspective, can influence which offers participants propose and which offers they accept. Based on this a third sub-question was formulated:

*What is the influence of personality traits regarding empathy on the performance on the colored trails negotiation game?*

This is tested via a questionnaire on four personality traits regarding empathy.

## **1.2 Thesis Structure**

This thesis starts with explaining the relevant concepts and describing the relevant literature in Chapter 2. Subsequently an overview of the used methods for all experiments is given in Chapter 3. The results are presented in Chapter 4, followed by a discussion in Chapter 5. Chapter 5 also contains the answers to the research questions stated above and options for further research.



## Chapter 2

# Literature

### 2.1 Theory of Mind

Theory of mind is the ability to attribute mental states to oneself, but also to others [1]. These mental states range from knowledge and beliefs, to intentions and desires. It is a system of inferences, which is called a theory because the states of the system are not directly observable and the system can be used to make predictions, in this case about the behavior of others.

Theory of mind is a concept which is studied in different fields, e.g. in philosophy by philosophers of mind and cognitive science, in biology by evolution theorists, and in psychology by animal psychologists and developmental psychologists [2].

There are different orders of theory of mind. In the case of zero-order theory of mind, one does reason about facts. With first-order theory of mind, one can reason about the mental states of another person. With second-order theory of mind, one can also reason about what ideas the other person has about one's own or someone else's mental states, and so forth.

A lot of research has been conducted to show the development of theory of mind in children, often via false belief tasks. The ability to use theory of mind develops gradually in early childhood: first-order theory of mind between the ages of three and five and second-order theory of mind around the age of five/six [2, 3, 4, 5, 6, 7]. Theory of mind experiments are also conducted with adults with, for example, strategic games [8] and negotiation settings [9]. Computer agents are also used to investigate theory of mind, in settings such as competitive settings [10] and negotiation settings [11, 12]. Other experimental studies on the topic of theory of mind are, amongst other things, about the influences of deficits like autism, mental handicaps and deafness [2, 13, 14] and about the controversial question whether animals use theory of mind or not [1, 15, 16].

### 2.2 Negotiations

People negotiate in order to reach an agreement. There are therefore many settings in which negotiations occur, a lot of them being quite trivial (shall

we buy the blue or the black car?), but some of them are very important, e.g. political (which budget to cut?) or economical negotiations (what wage to offer a future employee?). As a result, a lot of research on negotiations is conducted in the field of economics. For this research, subjects often have to play some sort of negotiation game. Some examples are: prisoners dilemma, ultimatum game, market game with proposer or responder competition, and trust game.

In a lot of these experiments, participants deviate significantly and consistently from the predictions of standard game theories [17]. For example, in ultimatum experiments, this might be caused by *other-regarding* behavior of the participants; this term includes concerns for fairness, the distribution of resources or the intentions of others [17]. The latter comes close to theory of mind, but Oxoby and McLeish [17] do not distinguish between the different orders at which one can reason about someone else.

In [18], another explanation for the differences is proposed. Fehr and Schmidt show that the seemingly contradictory evidence can be explained if one assumes that there is, in addition to purely selfish people, a fraction of the population that does care for equitable outcomes, i.e. that shows *inequity aversion*. The environment (i.e. the settings of the game and the distribution of the different types of players) influences how the players, intrinsically selfish or not, behave.

Another alternative can be given by using the *cognitive hierarchy* theory, where each player assumes that s/he uses the most sophisticated strategy of all players [19]. Each player assumes that the other players use less thinking steps. Many data sets and plausible restrictions suggest that the mean number of thinking steps is between 1 and 2. In [19], this parameter was set to 1.5. The model fitted data from different types of games as accurately as or even better than the Nash equilibrium.

When two or more parties negotiate, it can be very useful to think about what the other(s) want(s). If one does not consider the wishes of the other parties at all, it is very unlikely that they will accept one's offer. The other party could, of course, also think this. It might therefore be useful to reason about what the other party thinks you think. Negotiations are therefore interesting settings to test the use of theory of mind. Theory of mind research conducted in the negotiation setting up till now mainly used computer agent-based simulations. The current study lets computer agents and humans negotiate with each other, which is a set-up which has been used before in negotiation research on e.g. information revelation [20], cultural differences [21], and automated negotiation agents [22, 23]. This is done via the game colored trails, which is a multi-agent task environment for computer agents as well as for humans. An important difference with most other negotiation settings is that colored trails is a situated environment, while most other settings are very abstract [24]. This *situatedness* means that there is an interaction with the environment [25]. When a game is situated, it elicits stronger concerns with social factors, and when it is more abstract, people behave more in line with the Nash equilibrium play [26]. Situated games are therefore better if one wants to study real-life reasoning.

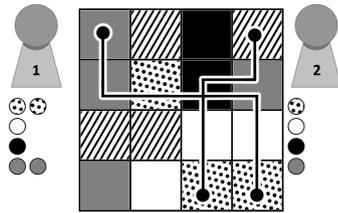


Figure 2.1: An example of the colored trails game. Player 1 starts at the left upper corner and its goal is the corner at the bottom right. Player 2 starts at the right upper corner and has to move to the corner at the bottom left. The lines show how close the players can get to their goal with the current distribution of the chips. (Adapted from [11].)

## 2.3 Colored Trails

*Colored trails* is a game which has been developed as a research test-bed [27, 28]<sup>1</sup> and can be played with various settings. Colored trails is a board game which is played on an  $n$  by  $n$  board which consists of colored tiles (see Figure 2.1). The game can be played by two or more players. The goal of the game is to move from a given start tile to a given goal location. Each player starts the game with a set of colored chips, which match the colors of the board. A player can only move to an adjacent tile (not diagonally) when s/he owns a chip with the same color as that tile. A chip can only be used once. To get as close to the goal as possible, the players need to negotiate about the distribution of the chips.

The game is abstract enough to represent many environments, which as a result does make the game situated. In general, it represents a complex negotiation situation [28]: the chips represent the skills and resources an agent owns. The board tiles are all different subtasks of which some need to be fulfilled in order to reach the goal. Matching colors between the chips and the board means that those skills and resources (chip) are necessary to complete the subtask (board tile). Not having all the chips needed to reach the goal at the start of the game represents that one is dependent on others. In the game the goal can be reached in several ways, which is usually also the case in real life.

There are many different settings for colored trails. Some examples of things which can be adjusted:

- complete or partial information (e.g. the goal or start location of the other player(s) can be unknown);
- the number of different tile colors (influences the complexity);
- the number of chips owned by the players (creates more or less opportunities per player);
- the size of the board (a larger board increases the complexity);
- the number of players (changes the dynamics of the game);

<sup>1</sup>See also <https://coloredtrails.atlassian.net/wiki/display/coloredtrailshome/>.

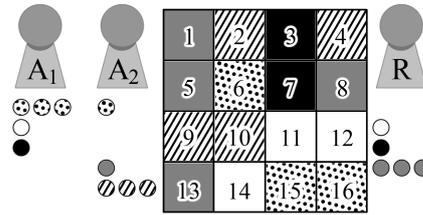


Figure 2.2: An example of the colored trails game with three players. Players  $A_1$  and  $A_2$  are Proposers (or Allocators) who, simultaneously, propose a distribution of their chips and those of  $R$ , the Responder. The Proposer then chooses the best offer. The Proposers try to move from square 1 to square 16, the Responder from 4 to 13. (From [11].)

- different scoring systems (e.g. whether there are bonus points for unused chips);
- number of rounds (changes the pressure on the negotiators and whether they need to give in a lot or not).

Some examples of how the game can be played follow below.

### Example 1: Three Players

Colored trails can be played with three players: two proposers and one responder. Both proposers simultaneously propose a distribution of the chips of the responder and their own chips (not those of the other proposer). The responder then chooses the best of the two offers (at random when they yield the same score). An example of this situation can be seen in Figure 2.2. Especially when it is a one-shot game, it is important that the proposers take the offer of the other proposer into consideration.

Dependent on the settings of the game, one can study different things. Some examples of the results one can gather: in [24] the game was played with three humans and it was found that humans are not reflexive, i.e. they do not base their decisions solely upon the options they have. They also reason about the other players in the game, both of them, also when there is uncertainty. The study by [11] used three agents and focused on the use of theory of mind: which order should one use? It was found that for a proposer, using second-order theory of mind is the superior tactic when the other proposer has a theory of mind as well. Otherwise first-order theory of mind suffices.

### Example 2: Two Players

The two-player setting is the setting used in this study. The two players need to negotiate about the distribution of their chips and they take turns in making a proposal. Therefore there is not a fixed proposer or responder and it is not a one-shot game. A more detailed description of a setting with two players can be found in Chapter 3.

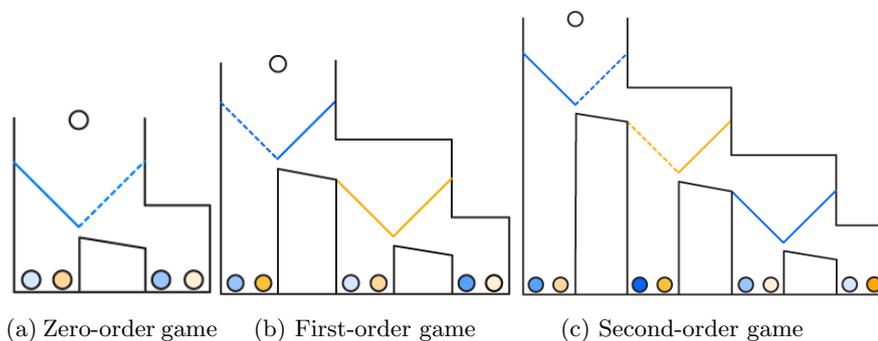


Figure 2.3: Examples of three types of marble drop. One player is blue, the other orange, and they decide which side of the trapdoors (diagonal lines) of their color to open to influence the trajectory of the white marble. The goal is to get the white marble in a bin with the darkest shaded marble of the right color. The dashed lines represent which side the player(s) should choose in order to get the best result. (Adapted from <http://www.ai.rug.nl/~meijering/MarbleDrop.html> ([29]).)

This set-up can also be used to investigate the reasoning of humans and the benefits of different orders of theory of mind. In [12], the benefits of the different orders of theory of mind were investigated by using computer agents, of which the results are shown in Section 2.6.2.

## 2.4 Training with Marble Drop

While negotiating, people might focus too much on their own goals and forget to use theory of mind, especially the higher orders. Training might compensate for this. In teaching people to negotiate, different tactics can be applied: principle-based (or didactic) learning, learning via information revelation, analogical learning, or observational learning. The study in [30] showed that the first two methods do not work very well and that participants in the observational learning group improved best, but they were not good in writing down the theory that was behind the tactics they used. They also showed that participants who received analogy training improved as well and, opposed to the observational group, these participants were able to write down what they did. Other studies also showed that analogical learning is a good method for learning to negotiate [31, 32]. Typically, the superficial structures of the base and target problem are different, but the underlying structures are similar. In the current study, the underlying similar structure is the need for the use of theory of mind. Our target problem was the colored trails negotiation game; as a base problem we used a game in which it is necessary and clear to use theory of mind: marble drop. The superficial structures of both games are different on two points: the appearances greatly differ and marble drop is a competitive game, as opposed to colored trails which is also cooperative.

In [29], people had to play *marble drop*. In this game, designed by Meijering,

two participants take turns in deciding in which direction a white marble will fall (see Figure 2.3). One player is orange, the other blue, and both try to get the white marble in the bin with the marble with the darkest shade of their color. They do this by deciding which trapdoor, left or right, to open, but they only control trapdoors of their own color (the colors of successive trapdoors alternate). When a trapdoor opens, the white marble falls into the underlying bin or it rolls to the next set of trapdoors. The number of trapdoors determines which order of theory of mind needs to be used. When there is only one trapdoor (Figure 2.3a), zero-order theory of mind suffices because the starting player can just choose the bin with the darkest color-graded marble. When there are two trapdoors (Figure 2.3b), one has to reason about what the other player will do at the second trapdoor. When there are three trapdoors (Figure 2.3c), one also has to reason about what the other player will think one will do at the third trapdoor.

To play marble drop, it is very clear that one has to use theory of mind. The study by [29] also revealed that, when knowingly playing against a rational computer agent opponent, the players used second-order theory of mind most of the time (94%) when it was necessary. Marble drop is therefore a good game to make people aware of theory of mind.

## 2.5 Interpersonal Reactivity Index

As a reminder, theory of mind is the ability to attribute mental states to oneself and to others. Empathy is the ability to infer *emotional* experiences about a person’s mental states and feelings, to attribute emotion to others [33]. In both theory of mind and empathy, perspective taking is involved and one needs to make a distinction between one’s own thoughts and those of others. Several studies showed that when people use theory of mind or empathy, part of the activated brain networks overlap [34, 35]. Schulte-Rüther and colleagues [35] concluded that theory of mind mechanisms are involved in empathy. The level of empathy a participant displays might therefore be related to his/her use of theory of mind in the colored trails negotiation game.

To test for this, the *Interpersonal Reactivity Index* (IRI) [33] was used. The IRI evaluates four aspects of empathy:

- **Perspective Taking scale:** tendency to spontaneously view something from someone else’s point of view;
- **Fantasy scale:** tendency to identify oneself with imaginative characters from e.g. books and movies;
- **Empathic Concern scale:** tendency to have feelings of compassion and concern for (unfortunate) others;
- **Personal Distress scale:** tendency to have feelings of anxiety and discomfort when viewing someone else’s negative experience.

The interpersonal reactivity index consists of 28 questions in total, seven questions per empathy aspect (see Appendix C.3). The questions are not asked per aspect of empathy but are mixed. There are five answer options ranging from ‘does not describe me well’ to ‘describes me very well’. Questions are

formulated in such a way that the answer option ‘does not describe me well’ results in a high score for some questions and in a low score for other questions. Via a formula, the final scores are calculated per empathy aspect, where a higher score means that someone has a higher tendency towards the behavior of that scale. Davis, the developer of the IRI, found that women generally score higher on all four scales when compared to men [33].

## 2.6 Basis for the Current Study

The current study is based on research conducted by De Weerd and colleagues [11, 12]. They conducted agent-based simulation studies in order to find an explanation for the evolutionary pressure of the development of theory of mind in humans. Both studies made use of colored trails and agents of different orders of theory of mind. An  $n$ -order theory of mind ( $ToM_n$ ) agent initially always believed that the other player was a  $ToM_{n-1}$  agent. Based on the observed behavior, this assumption could be adjusted downwards to the model which best predicted the other player’s behavior.

### 2.6.1 Three Negotiating Agents with Complete Information

In [11], a setting of colored trails with three computer agents was used (see Section 2.3, Example 1). De Weerd and colleagues simulated repeated single-shot games with complete information. The agents only had one goal: reaching the goal tile, because there was no bonus for unused chips. The responder always used zero-order theory of mind, because learning across games was not considered. The proposer agents did have theory of mind: zero-, first-, second-, third- or fourth-order. The responder always selected the best offer, unless it decreased her own score, without taking the scores of the proposers into consideration.

**Zero-Order Theory of Mind Proposer** This type of agent looks at all the chips with which it can make a proposal and then offers a trade that gets it to its goal or at least as close as possible. When there is more than one optimal option, one of those offers is selected at random. Since this agent does not take into account the desires of the responder, its strategy is not very successful.

**First-Order Theory of Mind Proposer** A  $ToM_1$  agent does take into account what the responder wants. It will therefore never propose a distribution which causes the responder to end up with fewer points, since such an offer will never be accepted. This agent also reasons about what the other proposer will offer to the responder, but assumes that the other proposer has zero-order theory of mind and will thus offer something which maximizes his own score. It makes the best offer possible which does not decrease the responder’s score and is better than the offer of the other proposer.

**Second-Order Theory of Mind Proposer** This agent does assume that the other proposer uses theory of mind. Therefore, the agent does not expect the other proposer to make an offer without taking into account the wishes of

the responder and the possible actions of the agent itself. The agent bases its own proposal on this information.

**Higher-Order Theory of Mind Proposer** The strategies of these agents are similar to those of the second-order theory of mind agent, but then with deeper nesting of beliefs.

The results of [11] showed that using first- and second-order theory of mind enhanced performance. A  $ToM_1$  proposer was always better than a  $ToM_0$  proposer, irrespective of the order of theory of mind used by the other proposer. When the competing proposer did not use theory of mind, using first-order theory of mind was the best tactic. When the other player did use theory of mind, second-order theory of mind yielded the best results.

## 2.6.2 Two Negotiating Agents with Incomplete Information

In [12], two computer agents played colored trails with incomplete information: they did not know the goal of the other player. This was not a one-shot game, instead the agents negotiated by alternately proposing chip distributions until an offer was accepted or until a player quit, in which case the initial distribution became final. By using theory of mind, the agents tried to figure out which tile the goal of the other player was. However, there was a penalty of one point per round of play. Another goal for the players was to own as many chips as possible, since there was a bonus for unused chips.

The key to good play in this variant of colored trails is to not only maximize one's own score, but to also try to enlarge the score of the other player. Then the other player will be much more inclined to accept an offer. In [12] this is called 'enlarging the shared pie'. With a larger pie there is a larger piece for both players. So some cooperation is necessary for an optimal result, but both players will of course try to get the larger piece. Agents with zero-, first-, and second-order theory of mind were used.

**Zero-Order Theory of Mind Negotiator** These agents base their beliefs and thus their offers solely on the behavior of the other player. For example, if an offer of four chips is declined, the agent believes that an offer with fewer chips will be declined as well. A learning parameter determines how much influence the behavior of the other player has on the beliefs of the zero-order theory of mind agent. (Learning speed is the degree to which an agent adjusts his beliefs, based on the observed behavior of the other player.)

**First-Order Theory of Mind Negotiator** A  $ToM_1$  agent considers what its proposal would look like from the perspective of the other player. He also forms ideas about the possible goal location of the other player and his beliefs; he can identify the interests of the other player. With this information it can adjust his own offers to make the other player believe things which will let him make an offer which is actually better for the  $ToM_1$  agent. This is, however, not a watertight strategy, because the  $ToM_1$  agent does not take into account the learning speed of the other agent but uses its own learning speed as an estimate instead. So its representation of the other player will only be correct if they have the same learning speed.

**Second-Order Theory of Mind Negotiator** A  $ToM_2$  agent believes that the other player might be a  $ToM_1$  agent. Therefore it thinks that the other player tries to interpret its offers to figure out what its goal tile is. So besides identifying the interests of the other player, it can also propose distributions of chips in order to ‘tell’ the other player what its own goal tile is. It could also use this to communicate other things, for example to manipulate the other player.

The study by De Weerd and colleagues [12] showed that when two  $ToM_0$  agents negotiate, there is an incentive not to give in to the other player, which often leads to an impasse. When a  $ToM_1$  agent negotiates with a  $ToM_0$  agent, the results are much better. The  $ToM_0$  agent benefits most from this: the  $ToM_1$  agent has to pay the costs for the cooperation. A  $ToM_2$  agent can negotiate successfully with a  $ToM_1$  agent, and since it can control the situation better than the  $ToM_1$  agent can, it benefits most from the cooperation. Negotiations between two  $ToM_2$  agents also work well.

### 2.6.3 Current Study

The study described in this thesis resembles the above-mentioned study [12] most. The difference between our study and [12] is that in our case only one of the players is a computer agent, the other is a human. Furthermore, the complexity of the settings is reduced, since humans have less processing power and speed than computers. This means that fewer colors will be used. To support the human player, a history panel (based on [28]) will be provided which shows all previous offers, categorized per game. The human players play against three different types of computer agents:  $ToM_0$ ,  $ToM_1$ , and  $ToM_2$  agents, as developed by De Weerd. A more elaborate description of the set-up is given in Chapter 3.

We hypothesize that the effectiveness of the different orders of theory of mind will remain the same as in [12]. The question is which order people will use. We hypothesize that at least part of the participants will use first-order theory of mind, since it is an order which one also uses in everyday life. They have experience with it so it will not be too hard to use. Using second-order theory of mind is already more difficult and people are less familiar with it. Therefore we expect that this order will mainly be used by participants in the training group, because they have actively used it just before they start the negotiations.

The agents in [12] adjust their own order of theory of mind based on the behavior of the opponent. They do this by matching their predictions with the actual outcomes. We hypothesize that this is harder to do for humans, since they are not such infallible calculators. We therefore expect that most of the participants will use the order of theory of mind they think is best and are able to use, irrespective of their opponent.



# Chapter 3

## Method

**Participants** 27 students of the University of Groningen (Groningen, the Netherlands) participated in the experiment (10 female, 17 male; age range: 18-27; mean age = 21.1). Of the participants, 18 (had) studied artificial intelligence, 3 computing science and 4 (had) studied other studies at the University of Groningen; 1 participant had studied at the Hanze University of Applied Sciences. The three participants who scored best on the Colored Trails part received €15,- / €10,- / €5,-. All participants gave informed consent before the experiment started.

**Apparatus** The experiment was conducted on a laptop that was running Windows 7 which was connected to a screen with a resolution of 1920 x 1080 pixels. The experiment was built in Java with Swing.

**Design** The experiment was a between-subjects design. One half of the participants (13) participated in the control group, which only had to do the zero-order variant of the marble drop game. The other half of the participants (14), the test group, also had to do the first- and second-order games of marble drop.

For the marble drop part, the independent variable was the variant of the game (zero-, first- or second-order). The dependent variable was the accuracy. For the colored trails part, the independent variable was the order of theory of mind used by the agent (zero, first or second). The dependent variables were the order of theory of mind most likely used by the participant and the score.

### 3.1 Introduction

#### 3.1.1 Procedure

The experiment took place in a quiet room. It started with a general instruction which told that the experiment consisted of several parts. The first part was a short questionnaire to gather demographic data. Since one should be able to distinguish between orange and blue for the marble drop part, the experiment then continued with a color test.

## 3.2 Color Test

### 3.2.1 Materials

The color test was based on [29] and consisted of two blocks with ten questions each. The first block tested whether the participant could distinguish between two different colors and the second block tested whether participants could distinguish between two different shades per color.

The colors were blue and orange, both with four different shades. Appendix A.1 shows the HTML-codes for all the colors. These colors were used throughout the whole experiment.

For block one, there were 16 possible color combinations (4 blue shades x 4 orange shades) which were all generated. Since this block consisted of only ten questions, per participant it was randomly determined which color combinations were used.

For the second block, there were 12 possible combinations per color, resulting in 24 questions in total (each shade of one color could be matched with three different shades of the same color). The block consisted of ten questions, five per color. Which combinations of shades were used was again randomly determined per participant.

### 3.2.2 Procedure

The participants read a short instruction, stating that there would be two blocks of ten questions each. They were told that in the first block they had to indicate which colored square was blue (or orange, this was randomly distributed among participants) and in the second block which colored square was darkest (or lightest, this was randomly distributed among participants). They had to indicate this by clicking on the correct colored square.

In the first block the participants received feedback after every question (“correct” or “incorrect”). At the end of the block it was stated how many questions they answered correctly. If this was eight or lower, the experiment stopped. In the second block, this procedure was the same. When they answered nine or ten questions correctly in this block, it was stated that they would continue to the next part of the experiment, otherwise the experiment would stop. The next part was the marble drop experiment.

## 3.3 Marble Drop

### 3.3.1 Materials

Set-ups consisting of bins, trapdoors and marbles were used, as described in Section 2.4. They were based on the stimuli used in [29]. Two colors were used (one for each player): orange and blue, both with four different shades. These were the same colors as were used for the color test.

There were three different types of marble drop games: with one, two or three trapdoors; each with two, three, and four bins, respectively. The color of the trapdoors, orange or blue, indicated who had to make a decision at those

trapdoors; the first set of trapdoors always matched the color of the participant. In each bin there was an orange and a blue marble, the marbles of the participant were always at the left side in a bin. Between bins, the shades of the colors differed. The marble with the darkest shade of the participant's color was the best marble, the one with the lightest shade of the participant's color the worst one. The computer always played optimally (maximizing its own score).

**Zero-order level** The games with one trapdoor did not require the use of theory of mind. All different permutations of the distribution of the marbles were used (four distributions).

**First-order level** The games with two trapdoors should be solved with first-order theory of mind. If there would be a marble of the best shade of the participant's color in the first bin, one would have to choose the first bin without taking into account the behavior of the other player. If there would be a marble of the worst shade, one would always continue to the other bins, without taking the other player into account either. To be able to check whether people used first-order theory of mind, those two kinds of pay-off structures were not used.

**Second-order level** Those two types of pay-off structures were not suitable for the games with three trapdoors either, for the same reasons. Furthermore, settings in which the shade of the computer's marble in the second bin is better or worse than *both* of his marbles in bins three and four were excluded as well. In those cases, the computer does not need to use first-order theory of mind to determine which side of the trapdoor to choose, so the participant does not need to use second-order theory of mind. Therefore, such settings are not indicative for the use of second-order theory of mind by the participants. Eight of the remaining possible pay-off structures were used.

Four zero-order games, eight first-order games and eight second-order games were created, all with different pay-off structures (Appendix A). They were balanced for the number of correct left/right trapdoor removals (for the predictions about the computer and for the decisions of the participant).

### 3.3.2 Procedure

The marble drop part was different for participants in the control and test group. Both groups were presented with a screen with instructions for the zero-order level marble drop games. It was stated whether they were blue or orange. It was explained that they had to try to let the black marble drop into the bin where the marble with the darkest attainable shade of their color was. It stated that they had to click on the trapdoor they wanted to open in order to do this. An example of a zero-order level game was presented (see Figure 3.1a) of which the answer was given.

The next two screens were only shown to participants in the test group. On these screens it was explained that there are also more complex games (first- and second-order) in which there would be interaction with the computer. It was stated that the computer played optimally (maximizing its own score) and that the game ended when the black marble had fallen in one of the bins. The explanations were accompanied by a picture of a first-order level game (Figure 3.1b) and a picture of a second-order level game (Figure 3.1c) and the correct

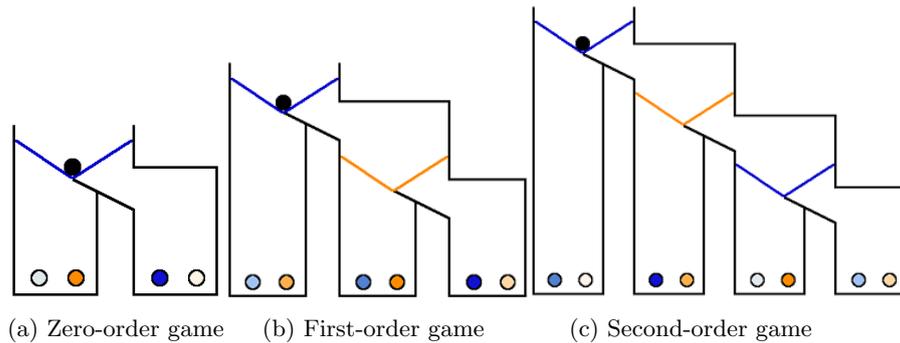


Figure 3.1: The marble drop structures that were used in the instructions of marble drop. (Participants who were the orange player received an orange version.)

answers. Both example settings were not used in the experiment itself.

All participants then started with four zero-order level games which were presented in a random order. After each game, “correct” or “incorrect” was displayed, according to the correctness of the given answer. If the answer was wrong, an arrow indicated what the correct answer would have been. After the four games were played, it was stated how many times they had given the correct answer.

For the control group, this was the end of the marble drop part. A screen was presented stating that they would continue to the next part of the experiment.

For the test group, eight first-order level games followed, in a random order. The participants had to make a decision at the first set of trapdoors, the computer at the second set. Again, after each game “correct” or “incorrect” was displayed to indicate whether they gave the right answer. When they did not, an arrow indicated what the right bin would have been. After the eight games were played, it was stated how many times they had given the correct answer.

For the test group, this part was followed by eight second-order level games which were presented in the same manner as the first-order games. After the participants made a decision at the first set of trapdoors, they were shown the action of the computer at the second set of trapdoors. Then they could make a decision at the third set of trapdoors (if the black marble was not in a bin yet). After the eight games were played, it was stated how many times they had given the correct answer.

Finally, a screen was shown to the test group that stated how many questions they had answered correctly in total and that they would continue to the next part of the experiment.

## 3.4 Colored Trails

### 3.4.1 Rules

As mentioned in Section 2.3, colored trails is a game which can be played with various settings. In the variant that was used for the current setting, the rules



Figure 3.2: The four different types of tiles used for the colored trails game.

were as follows.

Every player starts with four chips. The start position on the board is the center square. The goal tile is always at least three steps away from the start tile and can differ between participants. To move on the board, a chip of the correct texture is necessary, chips can only be used once. One can only move to adjacent tiles, but not diagonally.

A negotiation consists of maximally six rounds, which means every player can create at most three offers. Players take turns in making an offer. The starting player of a negotiation switches between games. Every round has a time limit of one minute, after which the turn switches to the other player. Per round, one can choose from three actions: accepting the last offer, creating a new offer, or withdrawing from the negotiation (after which the initial distribution becomes final). During the last round, it is only possible to accept or to withdraw. The negotiation is over when someone accepts an offer or when someone withdraws.

The score system was as follows: every participants starts a negotiation game with fifty points. If the goal is reached, one will earn another fifty points. When the goal is not reached, ten points will be deducted per missing step. Per chip which is not needed to get closer to the goal, five points are awarded.

Both players know that the rules and scoring system are common knowledge. Players only know their own goal tile.

### 3.4.2 Materials

The size of the board used for the colored trails experiment was 5 by 5 tiles. Each tile consisted of one of the four textures presented in Figure 3.2. The center tile, which was always the start tile, was black with an ‘S’ in it. The goal tile of the participant had a border in the color of the participant (blue/orange), the goal tile of the computer was not marked since this information was not given to the participant.

Both the participant and the computer received four chips at the start of a game. Each chip had one of the four textures from Figure 3.2. Via ‘spinners’ (graphic control element to adjust a value), the distribution of chips between players could be changed so the participant could form a new proposal. During the sixth round, the spinners were hidden so the participant was forced to choose for acceptance or withdrawal. This was done to make it more obvious that it was the last round.

Since humans do not have infallible memories and the computer agents do, the participants were provided with a ‘history panel’ (based on [28]). This panel showed all previous offers of the current game and of the previous games, categorized per game, with the most recent games at the top.

There was a time limit of one minute for each round and each game consisted

of maximally six rounds (a round was the action of one negotiator). The current round was indicated at the top of the screen and to indicate how much time was left, a timer was presented in the form of a countdown from 60 to 0 seconds.

### Agents

The agents used in this experiment were created by De Weerd and had been used before in the experiments presented in [12]. They have an internal model which evaluates the offer of an opponent. The agent matches the opponent's offer with what it thinks a  $ToM_0$ ,  $ToM_1$  and  $ToM_2$  agent would offer. Based on this, it creates beliefs and the beliefs help in deciding what to do. For example, a  $ToM_2$  agent which believes that his opponent is a  $ToM_0$  agent, might start behaving like a  $ToM_1$  agent. All agents had some basic knowledge of colored trails, gathered by playing 200 random games against other agents (training).

Due to the nature of the experiments in [12], the agents had no knowledge of the number of rounds. Therefore the agents did not take into account that the game ended after six rounds. To overcome this shortcoming, the agents were adapted in the following way.

**Round 6 (last round), turn of the agent** In this case the agent compares the initial distribution with the last offer. When the initial distribution is the best one of the two for the agent, the agent withdraws. Otherwise it will accept the last offer.

**Round 5, turn of the agent** The agent generates an offer. If this offer is equal to a previous offer from the participant, the agent proposes it, since it is very likely that the participant will accept it. If it does not equal one of the participant's previous offers, the agent will do the following to stay in control. It compares the initial distribution with the last offer. If the initial distribution is better for the agent, it will propose its generated offer: if the participant does not accept the offer, the agent will get the initial distribution and when the participant does accept it, the agent will get an even higher score. If the last offer is better than the initial distribution, the agent will not take the risk and will accept the last offer.

### Colored Trails scenarios

To make sure the scenarios used were relevant for the research, a selection was made from all the possible initial distributions of chips and board tiles. There were six categories for each of which four games needed to be found. The categories were:

1. An agent with zero-order theory of mind, who starts the first round of the negotiation
2. An agent with zero-order theory of mind, the participant starts the first round of the negotiation
3. An agent with first-order theory of mind, who starts the first round of the negotiation

4. An agent with first-order theory of mind, the participant starts the first round of the negotiation
5. An agent with second-order theory of mind, who starts the first round of the negotiation
6. An agent with second-order theory of mind, the participant starts the first round of the negotiation

To check whether a certain scenario is relevant for category 1, a  $ToM_0$  agent ('constant agent') played that scenario three times: against a  $ToM_0$ , a  $ToM_1$  and a  $ToM_2$  agent ('variable agent').

The game could end in three ways: 'acceptance', 'withdrawal' or 'out of time' (no agreement/withdrawal after six rounds). A scenario is indicative when playing a game with that setting against a  $ToM_0$ ,  $ToM_1$  and  $ToM_2$  agent all result in different end states. This could mean that the games ended in different ways. However, it could also mean that two or more games ended in an agreement, but with different final distributions.

Furthermore, there were also a few other criteria.

- The game should not be finished before the 'variable agent' performed two actions. In the real experiment, the participant is the 'variable agent'. For a better estimation of the strategy of the participant, it is better to have more actions available.
- It should be possible to end the game in six rounds. Therefore, only games in which an agreement was reached or withdrawal occurred within six rounds against all three types of opponents were selected. Note: when selecting the games, the agents did not have the extra rules about ending the game as described in Section 3.4.2.
- Ending a game with 'withdrawal' was only allowed against one of the three 'variable agents'. When 'withdrawal' is reached in different rounds, then the games are different, but ending up with different distributions is a 'stronger' difference and was therefore preferred.
- The goal was always reachable for the participant with the eight chips in play. However, only in four games (out of twenty-four) it was possible that both the agent and the participant could reach their goal at the same time, otherwise the negotiations would be too easy.

The scenarios that were used in the experiment are shown in Appendix B.1.

### 3.4.3 Procedure

This part of the experiment was the same for both groups. A first screen told the participants that they had to imagine they were an attorney for a big company. They would get involved in different negotiations for different clients. It was explained that the trading partner was played by the computer, 'Alex'. Alex would always react on a proposal as fast as possible and played optimally (maximizing its own score). This screen also stated that there were different means in play and that the participant and Alex had to negotiate about the distributions of those means via a game board.

The second screen explained the different aspects of the game board (example as shown in Figure 3.3), the chips (i.e. the means) and the basic principles of

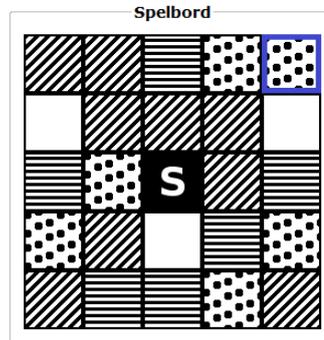


Figure 3.3: The example board used in the instructions of colored trails. (Participants who were the orange player received an orange version.)



- (a) The chips owned by the participant in the example.
- (b) The goal tile of the participant is marked by a border in the color of the participant (blue/orange). The lines are possible paths from the start tile towards that goal (only shown in the instructions).

Figure 3.4: Example in the instructions of colored trails to indicate possible paths from the start tile towards the goal with the chip distribution at the left. Only the most optimal paths are shown. (Participants who were the orange player received an orange version.)

the game. The participant was told that s/he could choose from three actions: counter-proposal, accept offer, withdrawal. It was also explained that the game could end in three ways: accept offer, withdrawal, after 6 rounds.

The next screen showed an example with possible paths. This can be seen in Figure 3.4. The fourth screen repeated the rules from the second screen and added the following: Alex and the participant take turns in executing one action. Per turn, there is a time limit of one minute and the game ends after six rounds (meaning that the initial distribution becomes final). It was clarified that six rounds means three actions from Alex and three from the participant.

The following screen told the participants that they and Alex would both start at the center of the game board and that only the goal of the participant was indicated (with a border in the color of the participant). A picture was shown to explain which tiles could be possible goal tiles (see Figure 3.5). It was also explained that all the rules were common knowledge, as were the chip

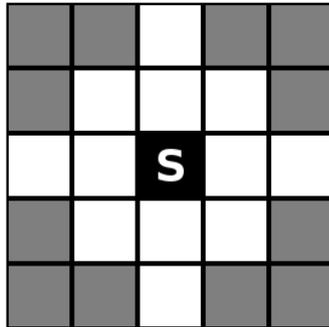


Figure 3.5: The black center square is the start tile for every player. Possible goal tiles are at least three tiles away from it, here indicated with gray tiles. (Adapted from [12].)

distributions. However, the goals of both players were only known by those players themselves.

The sixth screen contained information about how to make an offer, how to accept an offer and how to withdraw, as can be seen in Figure 3.6.

Then a screen with the scoring system (based on [28]) followed:

- you start with 50 points;
- if you reach your goal, you receive 50 bonus points;
- if you do not reach your goal, you get 10 points deduction per step that you are apart from the goal;
- per chip you do not use to get closer to your goal, you receive 5 bonus points.

This screen also included an example (see Figure 3.4) and the score one would get in this situation ( $50-10+5=45$ ). It was also mentioned that the best three negotiators would receive a monetary reward.

Then a screen was shown that explained how the history panel worked, which was accompanied by an example (see Figure 3.7). The participants were told that they could not make notes during the experiment. This screen also indicated that at the top of the negotiation screen, the number of the current game and round would be displayed and at the left, how much time was left and the participant's current score.

The next screen showed a screen shot from the complete negotiation game (Figure 3.8), so the participants could get familiar with the locations of all elements.

What followed was a short test to see whether the participant understood everything. In the introductory part it was not just possible to coordinate to the next screen, but also to previous screens so participants could go back to look up previous information. From this point on, that was not possible any more. The questions can be found in Appendix B.2. The next screen indicated whether the given answers were correct. If not, the correct answer was given with a short motivation.



Figure 3.6: Picture used in the instructions of colored trails. (Participants who were the orange player received an orange version.) In the top part, the participant could see the current offer. With the buttons ‘Accepteer verdeling’ and ‘Staak onderhandeling’, the participant could accept the offer or withdraw from the negotiation, respectively. Below that, the participant could adjust the distribution and make a counter-proposal (‘Doe voorstel’).

The last screen before the experiment started stated who would start the first game and that the start player would switch per game.

Then each participant had to finish three blocks of eight games each, but it was not indicated that there were different blocks. In one block the computer was a  $ToM_0$  agent, in another a  $ToM_1$  agent and in the third block a  $ToM_2$  agent. The order in which those blocks were presented was randomized between participants. The order of the eight games within each block was also randomized between participants. It was also randomized between participants who would start the very first negotiation, after that the starting player was alternated.

At the end of each game, the participants received their score and continued to a next negotiation after clicking on a button. After the last game, it was stated that the experiment was over and the total score was presented. The participant then continued to a questionnaire.

## 3.5 Questionnaires

### 3.5.1 Materials

The questionnaire consisted of three parts:

- marble drop experiment: questions about the perceived difficulty and rea-

Geschiedenis				
Ronde	Bod van	Jouw deel	Alex' deel	Wat
1	Alex			Tegenbod
Start				
5	Jou			Terugtrekking
4	Alex			Tegenbod
3	Jou			Tegenbod
2	Alex			Tegenbod
1	Jou			Tegenbod
Start				
3	Alex			Accepteert bod
2	Jou			Tegenbod
1	Alex			Tegenbod
Start				

Figure 3.7: An example of a history panel used in the instructions, where Alex is the computer. (Participants who were the orange player received an orange version.) ('Ronde' = round, 'Bod van' = offer from, 'Jouw deel' = your part, 'Alex' deel' = Alex' part, 'Wat' = what, 'Tegenbod' = counter-proposal, 'Terugtrekking' = withdrawal, 'Accepteert bod' = accepts offer.)

- soning strategies on all three levels of marble drop games;
- colored trails experiment: questions about the perceived difficulty and reasoning strategies;
- interpersonal reactivity index: questions to assess someone's score on four aspects of empathy;

All questions can be found in Appendix C.

### 3.5.2 Procedure

All participants started with questions about the  $ToM_0$  level marble drop games. The test group then continued with questions about the  $ToM_1$  and  $ToM_2$  level marble drop games. Then all participants answered questions about colored trails. This was followed by the questions from the interpersonal reactivity index. At the final screen the participants could leave general remarks on the experiment, after which the experiment ended and the participants were thanked for their participation.

## 3.6 Data Analysis

The significance level used for the tests presented in Chapter 4 was .05, unless stated otherwise.

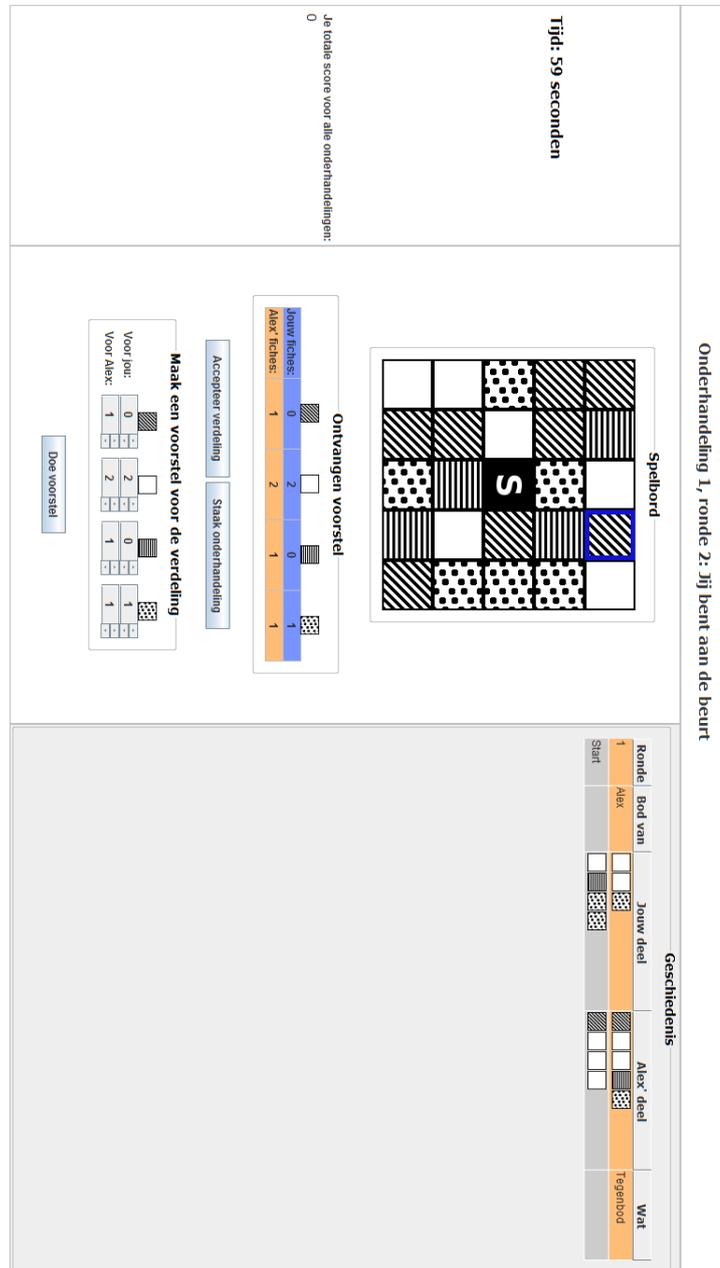


Figure 3.8: The interface for the colored trails game. (Participants who were the orange player received an orange version.)

# Chapter 4

## Results

### 4.1 Participants

There were 13 participants in the control group and 14 in the test group. All participants passed the color test. There was no correlation between the score (on marble drop, colored trails and the interpersonal reactivity index) of a participant and his/her age, gender, or field of study.

### 4.2 Marble Drop

#### 4.2.1 Accuracy

The games of the zero-order theory of mind level were played by all participants. Only one mistake was made in total, by a participant from the test group, which led to an overall accuracy of 99%. The games of the first-order theory of mind level were only conducted by the participants from the test group. The results are shown in Figure 4.1. The accuracies ranged from 63%-100% per person, with an average of 93% overall. The games of the second-order theory of mind level were also only played by the participants from the test group. The results are shown in Figure 4.2. The accuracies ranged from 38%-100% per person, with an average of 67% overall.

The accuracies found on the zero- and first-order theory of mind levels were as expected. The accuracy on the second-order theory of mind level, however, was lower than expected. The accuracy that was found in [29], the study on which the current marble drop setup was based, was much higher: participants used  $ToM_2$  90%-94% of the time when it was necessary.

Table 4.1 shows the average accuracies per second-order marble drop game. The table shows that the participants scored particularly poorly on game 3 (29% accuracy). This game is shown in Figure 4.3. The second worst game was game 5 (50% accuracy), which is also shown in Figure 4.3. These were the only games with pay-off structures where the black marble had to end up in the second bin from the left. In Section 5.1, some possible explanations are presented for the poor accuracies on this type of pay-off structure.

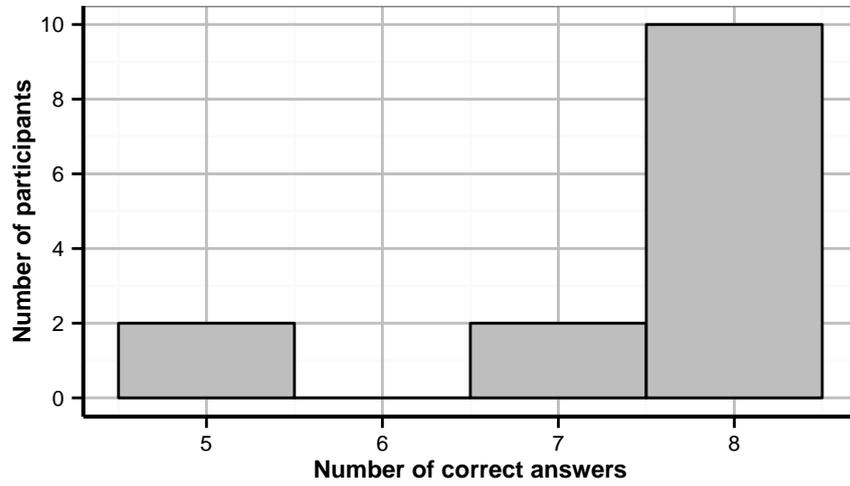


Figure 4.1: Score distribution of the  $ToM_1$  level marble drop games. The  $ToM_1$  level consisted of eight questions corresponding to different pay-off structures.

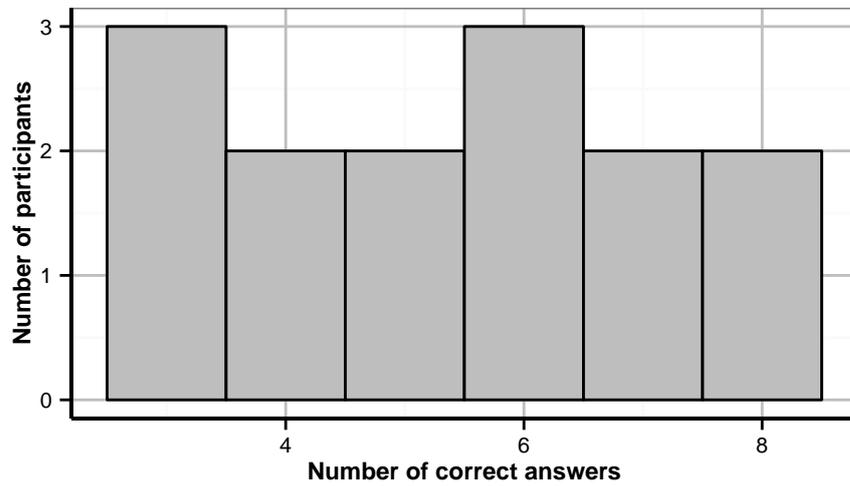
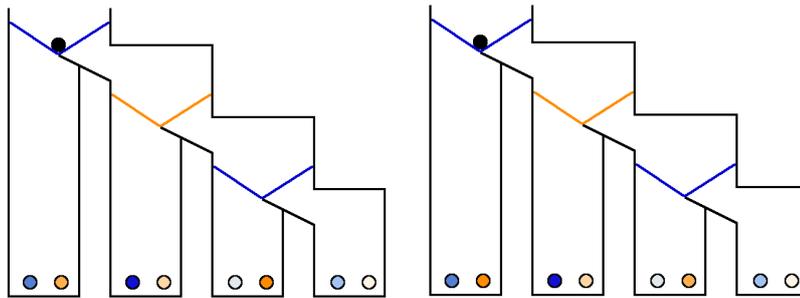


Figure 4.2: Score distribution of the  $ToM_2$  level marble drop games. The  $ToM_2$  level consisted of eight questions corresponding to different pay-off structures.

Table 4.1: The average accuracies over participants per marble drop game in the  $ToM_2$  level of marble drop. (A description of each game can be found in Appendix A.)

Game	Accuracy (%)
1	79
2	57
3	29
4	79
5	50
6	100
7	64
8	79



(a) The accuracy on this game was only 29%. (b) The accuracy on this game was only 50%.

Figure 4.3: Two  $ToM_2$  marble drop games. For both games, the most optimal reachable bin for the participant was the second bin from the left.

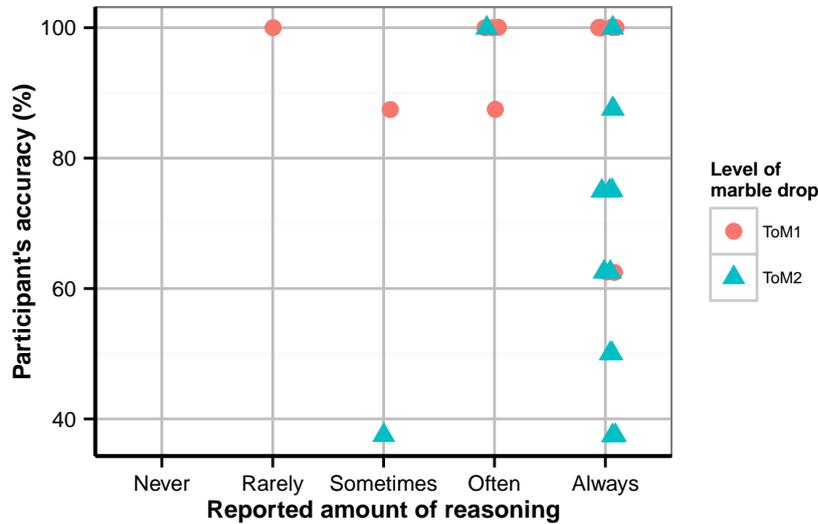


Figure 4.4: Scatter plot of the accuracies on the  $ToM_1$  and  $ToM_2$  level marble drop games and how much the participants reasoned about the opponent (the computer agent). The data comes from the test group only.

#### 4.2.2 Questionnaire: the Participant's Experience

In the questionnaire, the participants were asked how much they reasoned about the opponent for every level of marble drop. During the  $ToM_0$  level marble drop games, participants said not to have reasoned about the other player. Figure 4.4 shows how much the participants from the test group said they reasoned during the  $ToM_1$  and  $ToM_2$  level marble drop games and what their score was on those two levels. During the  $ToM_2$  level games, nearly all participants said they always reasoned about the opponent, while during the  $ToM_1$  level games this was less.

The comments the participants gave accompanying their answer on this question indicated that for the  $ToM_1$  level games of marble drop, everyone understood the strategy: looking at what the opponent would do at the next trapdoor. For the  $ToM_2$  level games, most of the participants understood how they should end up at the right conclusion, some did not, and for some it was not possible to tell. One participant indicated that at first s/he thought the opponent would always go in the direction of its darkest color, like in the  $ToM_1$  level games. Then s/he realized that the opponent reasoned about what s/he (the participant) would do.

The participants also answered questions about how hard they found the marble drop games. Figure 4.5 shows their answers and scores for the three marble drop levels: participants tended to report higher difficulties for the games of higher orders of theory of mind. There was a negative relation between score and reported difficulty for the  $ToM_2$  level games ( $r_{\tau}(12) = -.53$ ,  $p = .022$ ):



Figure 4.5: Scatter plot of the accuracies on the  $ToM_0$ ,  $ToM_1$  and  $ToM_2$  level marble drop games and how difficult the participants reported the games were. (Note: in the  $ToM_0$  setting there were 27 participants, in the  $ToM_1$  and  $ToM_2$  settings 14.)

participants with a lower accuracy tended to report higher difficulty.

## 4.3 Colored Trails

### 4.3.1 Scores

The distribution of the overall scores on colored trails is shown in Figure 4.6. The average score of the participants was 1521.67 points ( $SD = 195.89$  points) and the average score of the agents was 1516.67 points ( $SD = 143.69$  points). The correlation between the scores of the participants and agents was significant:  $r(25) = -.73$ ,  $p < .001$ . The higher the score of the participant, the lower the score of the agent tended to be and vice versa. Since the participant and agent often needed the same chip(s) and because ‘unused’ chips were also worth points, it is not strange to see that the success of one player resulted in losses for the other player.

The mean score of the test group on the colored trails experiment was 1515.00 points ( $SD = 207.81$  points), the mean of the control group 1528.85 points ( $SD = 190.40$  points). The difference between the scores was not significant ( $t(25) = 0.18$ ,  $p = .86$ ), opposed to what was expected. The hypothesis was that participants from the test group would outperform the participants from the control group.

The Pareto front, the solid line in Figure 4.6, is the optimal score on the colored trails experiment for a player, given the score of the opponent. It shows how much room everyone had to improve. The score a participant would have

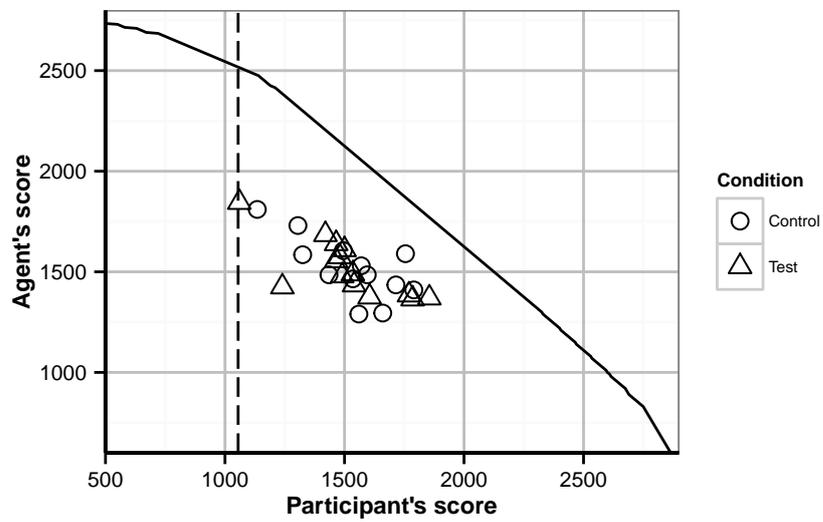


Figure 4.6: Scatter plot of the overall scores on the colored trails experiment for both the control and test group versus the agent's scores. The solid line depicts the Pareto front, the optimal score for the participant, given the score of the opponent. The dashed line shows the score participants would have gotten when they would have withdrawn every game: 1055 points. The lowest score of the participants was 1060 points.

Table 4.2: Average scores and standard deviations for both the participants and agents, per block and starting player.

Block	Starting player	Participants		Agents	
		<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
<i>ToM</i> <sub>0</sub>	Participant	58.06	27.87	51.30	16.64
<i>ToM</i> <sub>0</sub>	Agent	80.23	26.88	75.97	26.49
<i>ToM</i> <sub>1</sub>	Participant	49.72	20.09	48.98	17.30
<i>ToM</i> <sub>1</sub>	Agent	61.11	24.57	72.13	27.68
<i>ToM</i> <sub>2</sub>	Participant	70.93	27.45	60.79	24.48
<i>ToM</i> <sub>2</sub>	Agent	60.37	23.56	70.00	27.78

gotten if s/he would have chosen for withdrawal in every game is also shown, as a dashed line. Every participant scored higher than this ‘withdrawal score’ of 1055 points. The absolute minimum score a participant could get on the colored trails experiment was 380 points (when the opponent would always get all chips), the absolute maximum 2950 points (when the participant would always get all chips).

The more games one plays, the better one might get at colored trails. Figure 4.7 shows all participants’ scores on the colored trails experiment over time. The regression lines show that for most participants at least some improvement is visible. The standard error, however, is quite large. Therefore it is not possible to conclude that participants got better over time, although that is what the current data suggests.

It mattered for the participants and agents who started a colored trails game. The mean scores per block can be found in Table 4.2. Against the *ToM*<sub>0</sub> and *ToM*<sub>1</sub> agent participants generally scored better when the opponent started than when they started themselves. Against the *ToM*<sub>2</sub> agent this was the other way around. For the agent it was always better to start a game itself, regardless of the order of theory of mind it used. In every block the scores between games started by the agent and scores from the games started by the participant were significantly different, for both players:

- **Block *ToM*<sub>0</sub>, participant’s score:**  $Z = -6.05, p < .001$
- **Block *ToM*<sub>1</sub>, participant’s score:**  $Z = -4.45, p < .001$
- **Block *ToM*<sub>2</sub>, participant’s score:**  $Z = -2.73, p = .0063$
- **Block *ToM*<sub>0</sub>, agent’s score:**  $Z = -7.72, p < .001$
- **Block *ToM*<sub>1</sub>, agent’s score:**  $Z = -6.68, p < .001$
- **Block *ToM*<sub>2</sub>, agent’s score:**  $Z = -3.32, p < .001$

The highest obtained scores for the participants (more than 100 points, which means the participant reached the goal and had (an) unused chip(s)) all came from games where the agent started the game. The influence of the starting player might come from the first offers.

The first offers of the participants and agents show similarities and differences (see Table 4.3). With their first offer both players generally ascribed a high score to themselves. There was no significant difference found between the score a participant would get from his own first offer and the score the agent

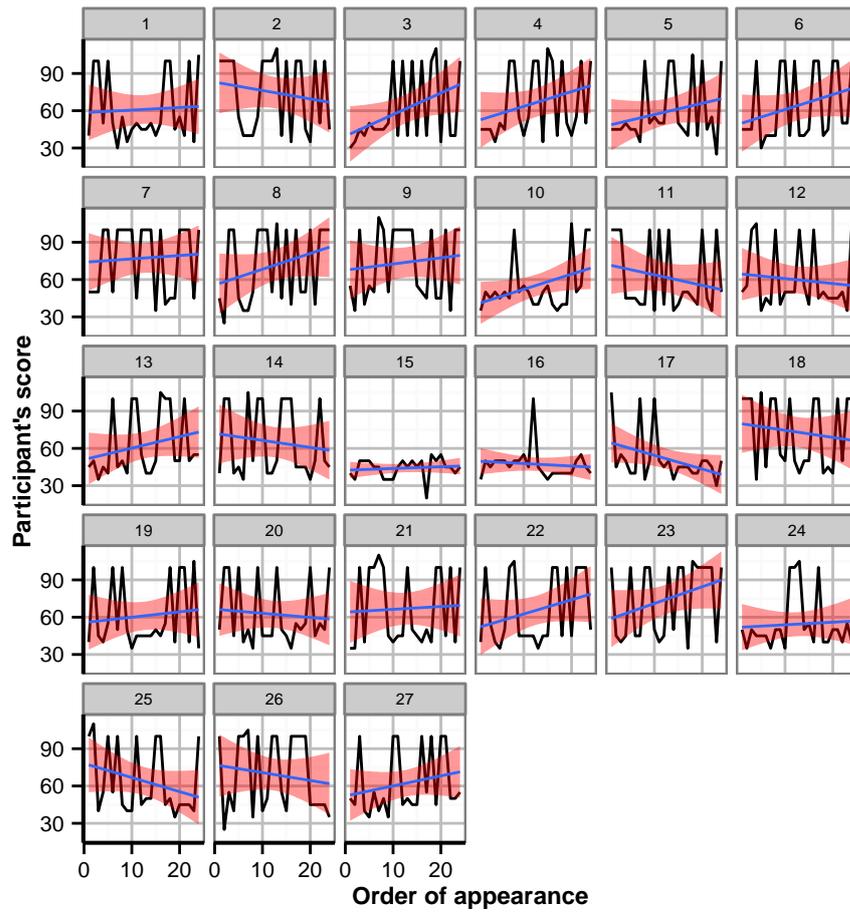


Figure 4.7: Every window shows the scores on the colored trails games for one participant. The scores are ordered in order of appearance of the games (which was different per participant). The blue lines depict the regression lines and the red areas their standard errors.

Table 4.3: Average scores and standard deviations from the first offers for both the participants ('part.') and agents, per block and starting player.

Block	Starting player agent				Starting player part.			
	Score agent		Score part.		Score part.		Score agent	
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>
<i>ToM<sub>0</sub></i>	88.75	22.57	48.75	6.53	93.84	21.76	44.81	12.11
<i>ToM<sub>1</sub></i>	88.75	19.58	45.00	8.70	90.79	23.37	36.94	8.93
<i>ToM<sub>2</sub></i>	103.75	2.18	46.25	6.53	96.53	19.71	47.73	16.43

would get from its own first offer for block *ToM<sub>0</sub>* and block *ToM<sub>1</sub>*. However, there was a significant difference between the scores the participants and agents ascribed to themselves in the *ToM<sub>2</sub>* block: the agents ascribed higher scores to themselves than the participants did.

- **Block *ToM<sub>0</sub>*:**  $Z = -1.85, p = .064$
- **Block *ToM<sub>1</sub>*:**  $Z = -0.80, p = .43$
- **Block *ToM<sub>2</sub>*:**  $Z = -3.36, p < .001$

The participant's score from the agent's first offer and the agent's score from the participant's first offer did significantly differ in block *ToM<sub>0</sub>*: the agent tended to assign a higher score to the participant with its first offer than the participant to the agent. The same was found in block *ToM<sub>1</sub>*. The difference in block *ToM<sub>2</sub>*, however, was not significant.

- **Block *ToM<sub>0</sub>*:**  $Z = -5.95, p < .001$
- **Block *ToM<sub>1</sub>*:**  $Z = -7.33, p < .001$
- **Block *ToM<sub>2</sub>*:**  $Z = -0.97, p = .33$

In general, most of the time the first offer was better when the agent started as opposed to when the participant started.

Figure 4.8 indicates that the final scores of the participants and agents on colored trails games differed per block (*ToM<sub>0</sub>*, *ToM<sub>1</sub>* or *ToM<sub>2</sub>* agent) and this was supported by significant Friedman tests:  $\chi^2(2) = 21.87, p < .001$  for the participants' scores and  $\chi^2(2) = 12.05, p = .0024$  for the agents' scores. Post hoc analysis with Wilcoxon signed-rank tests was conducted with a Bonferroni correction applied, resulting in a significance level of .017.

- **Blocks *ToM<sub>0</sub>* & *ToM<sub>1</sub>*, participant's score:**  $Z = -4.01, p < .001$
- **Blocks *ToM<sub>0</sub>* & *ToM<sub>2</sub>*, participant's score:**  $Z = 0.21, p = .58$
- **Blocks *ToM<sub>1</sub>* & *ToM<sub>2</sub>*, participant's score:**  $Z = -5.14, p < .001$
- **Blocks *ToM<sub>0</sub>* & *ToM<sub>1</sub>*, agent's score:**  $Z = -2.95, p = .0032$
- **Blocks *ToM<sub>0</sub>* & *ToM<sub>2</sub>*, agent's score:**  $Z = -0.14, p = .89$
- **Blocks *ToM<sub>1</sub>* & *ToM<sub>2</sub>*, agent's score:**  $Z = -2.54, p = .011$

The participants' scores on games with a *ToM<sub>1</sub>* opponent were different from the scores on *ToM<sub>0</sub>* opponent games and the *ToM<sub>2</sub>* opponent games. The participant's scores on *ToM<sub>0</sub>* opponent and *ToM<sub>2</sub>* opponent games were not significantly different. The same was true for the scores of the agent: for both players the scores in the block with the *ToM<sub>1</sub>* agent were significantly lower than the scores in the *ToM<sub>0</sub>* and *ToM<sub>2</sub>* agent blocks.

The participants' and agents' final scores did not significantly differ in the block with the  $ToM_0$  agent, nor in the block with the  $ToM_2$  agent. The difference was significant in the block with the  $ToM_1$  agent.

- **Block  $ToM_0$ :**  $Z = -0.55, p = .58$
- **Block  $ToM_1$ :**  $Z = -2.97, p = .0029$
- **Block  $ToM_2$ :**  $Z = -0.73, p = .47$

The agents scored higher than the participants (see Figure 4.8). This might be due to the nature of the  $ToM_1$  agent, which tried to manipulate the participant. This could also explain the low scores found for the participant in the  $ToM_1$  block.

Figure 4.8 also shows that high scores (100-110 points) on colored trails games were only obtained when one of the players accepted an offer. This was due to the nature of the selected games: it was not possible to reach the goal with the initial set of chips (which became final after a withdrawal). Most of the time when someone obtained a high score, the other player accepted the offer. Apparently, players only rarely created an offer which led the other player to the goal. Maybe the players did not figure out what the other player's goal tile was or just did not want to give the chips the other player needed because this would result in a low score for the player itself.

Scores higher than 100 points for the participant were only obtained in the  $ToM_0$  opponent games (except for one score in the  $ToM_1$  block). In every game with a score of 110 points for the participant, the participant accepted the final offer of the agent. This might be due to the 'naive' strategy of the  $ToM_0$  agent. Very high scores for the agent were obtained via acceptance by both the agent and the participant and occurred in all blocks. This could indicate that the participants were less demanding or were more altruistic than the agents.

There was a relation between the participants' overall scores on the colored trails experiment and the total time they spent on the experiment:  $r_{\text{tau}}(25) = .29, p = .032$ . The more time a participant invested in the colored trails experiment, the higher the final score tended to be, which is not a surprising result. This is shown in Figure 4.9. Grubb's test showed that there were two outliers in one tail ( $G = 0.52, p = .020$ ). When these outliers were removed, there was still a significant relation found between the participant's overall scores on the colored trails experiment and the total time they spent on the experiment:  $r_{\text{tau}}(25) = .34, p = .016$ .

### 4.3.2 Used Orders of Theory of Mind by the Participant

According to theory [12], if the agent uses  $ToM_n$ , the participant should use  $ToM_{n+1}$  to be able to fairly accurately predict the agent's actions. So when Alex was a  $ToM_0$  agent, it was sufficient for the participants to use  $ToM_1$ . Since Alex used  $ToM_2$  as his highest order, the highest order of theory of mind the participants would have needed was  $ToM_3$ . This order, however, was not evaluated. The first reason was that there was no indication at all that the participants had used this, based on the participants' comments. Secondly, as stated in Section 2.2, people usually use one to two thinking steps in playing (strategic) games [19]. Given these factors and the fact that the participants

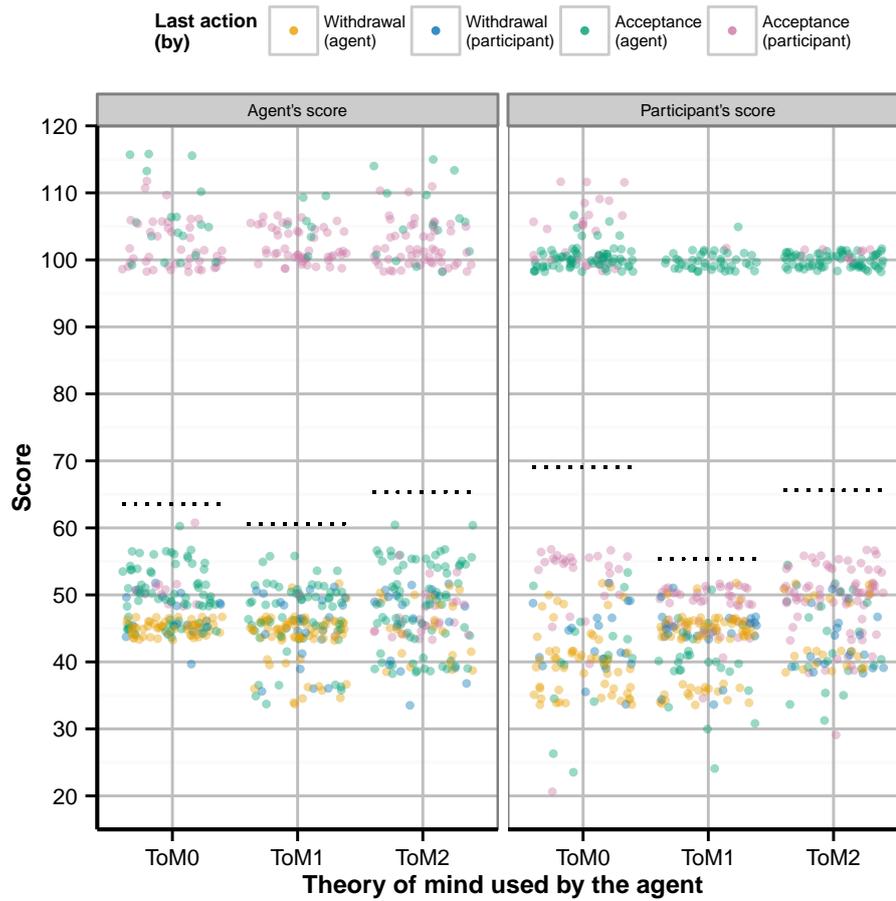


Figure 4.8: Scores of the participants (right) and agents (left), per colored trails game, split by the order of theory of mind used by the agent. The colored points represent all data points: the final score per player, per game. The colors indicate what the last action was and who performed that action. The dotted lines depict the mean scores.

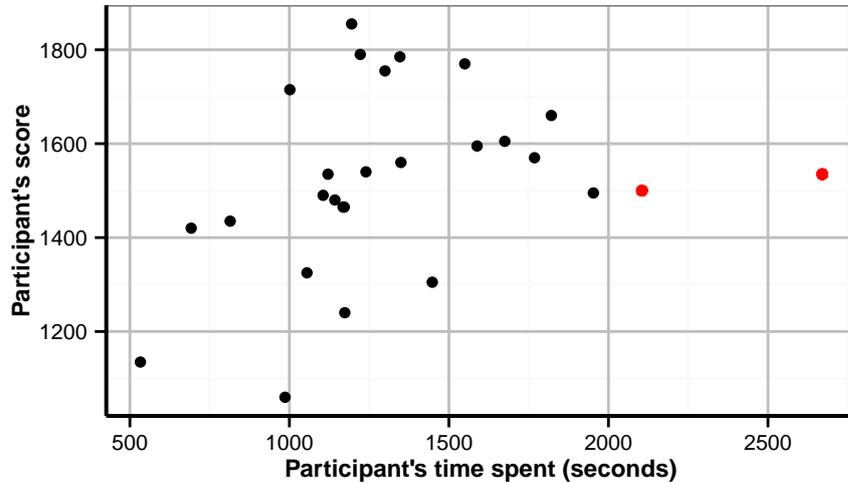


Figure 4.9: Scatter plot of the participants' total time spent on the colored trails experiment and their score. The two red dots depict outliers.

only had one minute per round to come up with a plan, it was very unlikely that they used third-order theory of mind.

The agent, Alex, used an evaluation algorithm to decide what to offer. This system was also used to evaluate what order of theory of mind a participant most likely used. The agent compared every offer of the participant with what it would expect a  $ToM_0$ ,  $ToM_1$  and  $ToM_2$  agent would offer. The closer the participant's offer matched one of those offers, the more confident the agent became in the fact that the participant was actually using that order of theory of mind (and therefore the confidence in the use of the other orders decreased). So at every step, the agent had three confidences per participant: a confidence that the participant used  $ToM_0$ , a confidence that the participant used  $ToM_1$ , and a confidence that the participant used  $ToM_2$ . The confidences were adjusted at every offer and were not reset during the colored trails experiment. The agent was initially set to be most confident that the participant used  $ToM_2$ . The average confidences per theory of mind order are shown per block of the experiment in Figure 4.10.

A Friedman test showed that the confidence levels for all three orders of theory of mind significantly differed:  $\chi^2(2) = 119.14$ ,  $p < .001$ . Post hoc analysis with Wilcoxon signed-rank tests was conducted with a Bonferroni correction applied, resulting in a significance level of .017.

- **Confidences in  $ToM_1$  &  $ToM_2$ :**  $Z = -2.59$ ,  $p = .0096$
- **Confidences in  $ToM_0$  &  $ToM_2$ :**  $Z = -7.81$ ,  $p < .001$
- **Confidences in  $ToM_0$  &  $ToM_1$ :**  $Z = -7.82$ ,  $p < .001$

The confidence in a  $ToM_2$  participant was significantly different from the confidence in a  $ToM_1$  participant as well as from the confidence in a  $ToM_0$  participant. The confidence in a  $ToM_1$  participant was also significantly different from

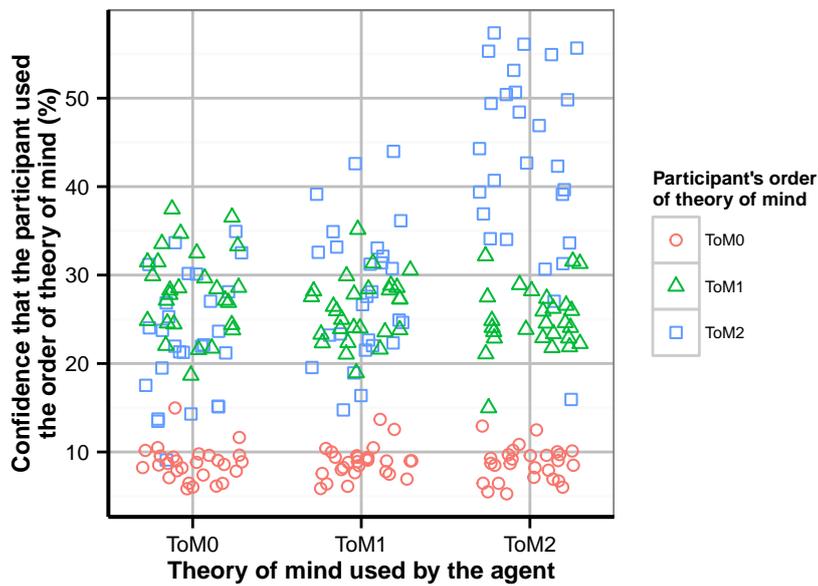


Figure 4.10: Scatter plot of the confidences that a participant most likely used a certain order of theory of mind, split per block in the colored trails experiment (opponent who used  $ToM_0$ ,  $ToM_1$  or  $ToM_2$ ). Note that the order of the blocks was randomized. Every data point is the mean confidence in one of the orders of theory of mind of one participant for a particular block. The confidences do not add up to 100%, the remainder is the uncertainty of the system.

the confidence in a  $ToM_0$  participant. The evaluation system thus managed to make a clear distinction between its confidences in all three orders of theory of mind. The confidence in a  $ToM_0$  participant was lowest, as can be seen in Figure 4.10, which indicates that the participants used higher orders theory of mind while playing colored trails.

The confidences for the three orders of theory of mind were compared between the participants from the control and test group. There was no significant difference between the two groups for the  $ToM_2$  confidence,  $ToM_1$  confidence, nor the  $ToM_0$  confidence. This was (again) opposed to what was expected.

- **Confidence in  $ToM_2$ :**  $t(76.35) = -0.053$ ,  $p = .96$
- **Confidence in  $ToM_1$ :**  $t(77.89) = -0.98$ ,  $p = .33$
- **Confidence in  $ToM_0$ :**  $t(77.63) = -0.53$ ,  $p = .59$

It was tested whether the confidences in the order of theory of mind used by the participants was different between the three blocks (one block against a  $ToM_0$  agent, one against a  $ToM_1$  agent and one against a  $ToM_2$  agent). A Friedman test showed that the confidences did not significantly differ between the blocks ( $\chi^2(2) = 5.65$ ,  $p = .059$ ).

To test whether the separate confidences per order of theory of mind did differ per block, three more Friedman tests were performed.

- **Confidence in  $ToM_2$ :**  $\chi^2(2) = 24.89$ ,  $p < .001$
- **Confidence in  $ToM_1$ :**  $\chi^2(2) = 2.67$ ,  $p = .26$
- **Confidence in  $ToM_0$ :**  $\chi^2(2) = 0.52$ ,  $p = .77$

The test was not significant for the confidence in  $ToM_0$  over the three blocks, nor for the confidence in  $ToM_1$ . For the confidence in  $ToM_2$  however, the confidences did significantly differ per block.

To see where the differences between the blocks lay for the  $ToM_2$  confidence, post hoc analysis with Wilcoxon signed-rank tests was conducted with a Bonferroni correction applied, resulting in a significance level of .017.

- **Blocks  $ToM_0$  &  $ToM_1$ :**  $Z = -2.25$ ,  $p = .025$
- **Blocks  $ToM_0$  &  $ToM_2$ :**  $Z = -5.08$ ,  $p < .001$
- **Blocks  $ToM_1$  &  $ToM_2$ :**  $Z = -4.24$ ,  $p < .001$

The  $ToM_2$  confidences were not significantly different between the blocks with a  $ToM_0$  and  $ToM_1$  opponent. There was a significant difference between the  $ToM_2$  confidences for the  $ToM_0$  and  $ToM_2$  opponent. The difference between the  $ToM_2$  confidences for the  $ToM_1$  and  $ToM_2$  opponent was significant as well. In both cases, the confidence in  $ToM_2$  was higher in the block with the  $ToM_2$  opponent. This indicates that the participants most likely used more second-order theory of mind when they were (unknowingly) facing a  $ToM_2$  agent.

There was no correlation between the overall score on the colored trails experiment and the average confidence in the participant using  $ToM_2$  or  $ToM_1$ . The correlation between the overall score on the colored trails experiment and the average confidence in the participant using  $ToM_0$ , on the other hand, was significant: the higher the achieved score on the colored trails experiment, the higher the confidence in the use of  $ToM_0$ .

- **Confidence in  $ToM_2$ :**  $r(25) = .17, p = .39$
- **Confidence in  $ToM_1$ :**  $r(25) = -.044, p = .83$
- **Confidence in  $ToM_0$ :**  $r(25) = .47, p = .013$

Since the confidences that the participants might have used  $ToM_0$  were very low, it is hard to draw conclusions from this result.

Since more reasoning steps are needed for higher orders of theory of mind, it was tested whether there was a difference in reaction times between the three blocks. The total time per game was used, which could maximally be 180 seconds (three rounds, each with a time limit of one minute). A Friedman test showed that there was a significant difference:  $\chi^2(2) = 9.56, p = .0084$ . Again, post hoc analysis with Wilcoxon signed-rank tests was conducted with a Bonferroni correction applied, resulting in a significance level of .017.

- **Blocks  $ToM_0$  &  $ToM_1$ :**  $Z = -2.40, p = .016$
- **Blocks  $ToM_0$  &  $ToM_2$ :**  $Z = -0.62, p = .54$
- **Blocks  $ToM_1$  &  $ToM_2$ :**  $Z = -3.16, p = .0016$

The reaction times against the  $ToM_0$  agent did not differ from the reaction times against the  $ToM_2$  agent. The reaction times against the  $ToM_1$  opponent did significantly differ from the reaction times against the  $ToM_0$  agent and the  $ToM_2$  agent. The participants needed more time against the  $ToM_1$  agent ( $M = 59.96$  s,  $SD = 29.48$  s) than against the  $ToM_0$  agent ( $M = 54.54$  s,  $SD = 28.57$  s) and  $ToM_2$  agent ( $M = 53.07$  s,  $SD = 24.88$  s). This might be the result of the manipulating behavior of the  $ToM_1$  agent, which could have hindered the participants. It is surprising to see that no more time was needed in the  $ToM_2$  block compared to the  $ToM_0$  block, even though the participants most likely used more  $ToM_2$  themselves in the  $ToM_2$  block, which would have resulted in more reasoning steps.

### 4.3.3 Questionnaire: the Participant's Experience

The participants filled out a questionnaire about the colored trails experiment. The first question was whether they found this part difficult. There was no significant difference between the reported difficulty by the control and test group ( $U = 72.00, Z = -1.08, p = .28$ ) and there was no relation between the reported difficulty and the achieved overall score on the colored trails experiment ( $r_{\text{tau}}(25) = -.30, p = .055$ ). The results are shown in Figure 4.11. The answers show that the colored trails experiment was a challenge for the participants, but not one which was too difficult.

An inspection of the scatter plot in Figure 4.11 indicates that there might be an outlier (data point at bottom left). This was confirmed by Grubb's test for one outlier ( $G = 2.78, p = .036$ ). With the outlier removed, the reported difficulties from the control and test group were still not significantly different ( $U = 72.00, Z = -1.01, p = .47$ ). However, without the outlier, the relation between reported difficulty and the achieved score on the colored trails experiment was significant:  $r_{\text{tau}}(25) = -.45, p = .0061$ . Participants with higher scores tended to report lower difficulty.

The following question was about how they reasoned about the other player ('Alex'). Did they reason about what Alex would do? The results are shown in

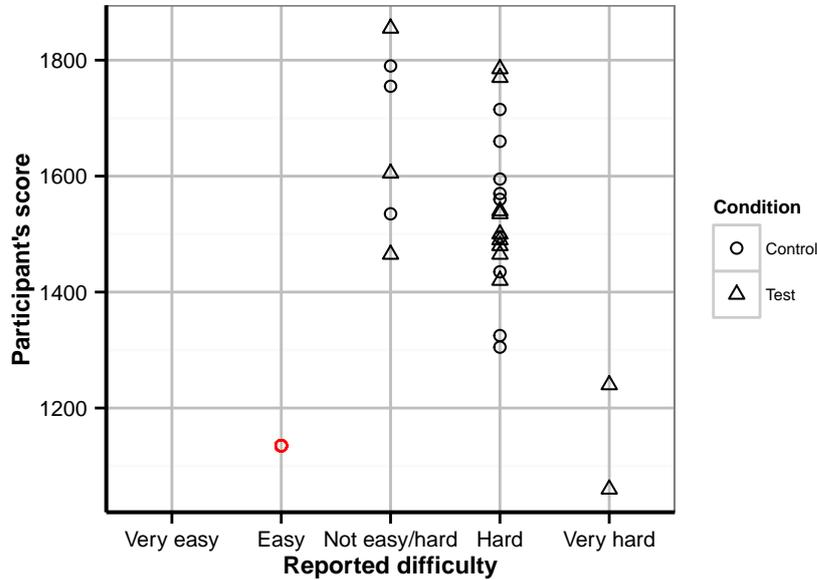


Figure 4.11: Scatter plot of the overall scores on colored trails and how difficult the participants reported the game was. The red data point is an outlier.

Figure 4.12, which shows that most participants did reason about Alex to some extent. There was no significant difference in the reported reasoning about Alex between the control and test group ( $U = 87.50$ ,  $Z = -0.15$ ,  $p = .88$ ). The relation between the reported reasoning about Alex and the achieved overall scores on colored trails was not significant either ( $r_{\tau}(25) = .23$ ,  $p = .14$ ). There was also no significant relation between the reported amount of reasoning about Alex and the agent's confidence in the participant's use of either of the three orders of theory of mind.

- **Confidence in  $ToM_2$ :**  $r_{\tau}(25) = -.16$ ,  $p = .29$
- **Confidence in  $ToM_1$ :**  $r_{\tau}(25) = -.13$ ,  $p = .40$
- **Confidence in  $ToM_0$ :**  $r_{\tau}(25) = .22$ ,  $p = .14$

There was another question about the reasoning of the participants: did they reason about the goal location of Alex? The results are shown in Figure 4.13, which shows the most of the participants reasoned about Alex' goal location to some extent, but less than about what Alex would do (Figure 4.12). There was no significant relation between the reported amount of reasoning about Alex' goal and the group a participant was in ( $U = 110.50$ ,  $Z = -0.96$ ,  $p = .34$ ), nor between this reported reasoning and the participant's overall score on colored trails ( $r_{\tau}(25) = .076$ ,  $p = .61$ ). The relation between the reported amount of reasoning about Alex' goal and the confidence in the used order of theory of mind was not significant either.

- **Confidence in  $ToM_2$ :**  $r_{\tau}(25) = -.25$ ,  $p = .10$

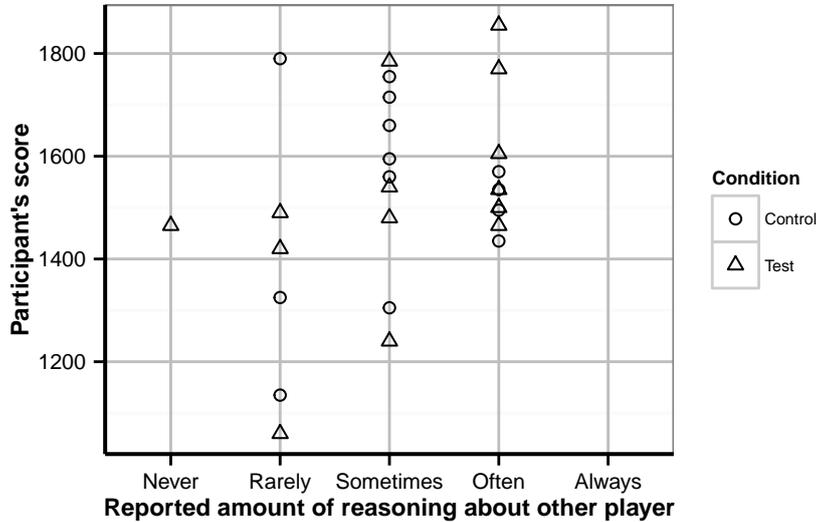


Figure 4.12: Scatter plot of the overall scores on colored trails and the reported amount of reasoning by the participants about what the other player ('Alex') would do.

- **Confidence in  $ToM_1$ :**  $r_{\text{tau}}(25) = -.043$ ,  $p = .78$
- **Confidence in  $ToM_0$ :**  $r_{\text{tau}}(25) = .18$ ,  $p = .25$

In the comments the participants gave on the previous two questions, many participants stated that they looked at which chips Alex was consistently giving away / keeping in a colored trails game. Some people also indicated that they looked at the possible paths Alex could create with the chips he asked, for example to be able to 'propose' an alternative route. One person said that it would not be necessary to reason about Alex' goal and a few others stated that they did not have enough time for that. One participant indicated that s/he took into consideration that after the third round it was likely that the agent would withdraw. Overall the results from the questionnaire support the idea that the participants used higher orders of theory of mind.

As stated before, the colored trails experiment was divided in three parts: one block against a  $ToM_0$  agent, one block against a  $ToM_1$  agent and one block against a  $ToM_2$  agent. This was not pointed out to the participants and therefore they were asked in the questionnaire whether they thought Alex changed his strategy during the colored trails experiment and if so how often. Unfortunately, the participants' comments on this question revealed that this question was not expressed specifically enough: part of the participants thought they were asked about strategy change during a single colored trails game. However, the comments of some participants showed that they did interpret the question correctly. Appendix D.1 shows those comments. Since these data were not reliable, no tests were performed on it.

The final question was whether the participants thought that playing marble

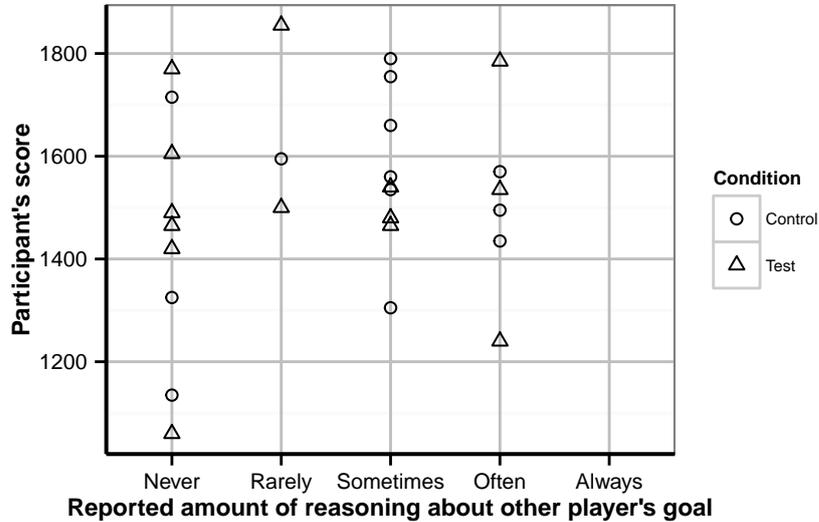


Figure 4.13: Scatter plot of the overall scores on colored trails and the reported amount of reasoning by the participants about the other player's ('Alex') goal location.

drop first, influenced how they played colored trails. There were only five (out of twenty-seven) participants who reported that marble drop had influenced them, which is consistent with the missing differences between the control and test group. Three of the participants who did report an influence were in the test group, the other two in the control group. The achieved scores of those participants varied from 1465 to 1790 points ( $M = 1554$  points,  $SD = 134.88$  points).

Three of those participants (all from the test group) commented that it had helped because of the practice in / awareness of thinking about the actions of the opponent. Another participant, from the control group, stated that it had helped him/her because the marble drop games stressed that the score of the opponent was not important, which made him/her less vindictive during the colored trails games. The last participant stated that it helped with thinking about what the goal of the opponent could be and basing choices on that. This is an indication that marble drop might support people in playing colored trails.

#### 4.3.4 Influence of Marble Drop

For determining whether there was a relation between the results on marble drop and on colored trails, only the data from the test group was used. Since all participants in the control group obtained a 100% accuracy on the marble drop games they played (zero-order level), it was not meaningful to use those results.

There was no correlation between the overall score on marble drop and the

score on colored trails ( $r(12) = .19, p = .50$ ). There was no correlation between the overall score on marble drop and the average confidence in the participant using any of the three orders of theory of mind.

- **Confidence in  $ToM_2$ :**  $r(12) = -.31, p = .28$
- **Confidence in  $ToM_1$ :**  $r(12) = -.18, p = .53$
- **Confidence in  $ToM_0$ :**  $r(12) = .23, p = .42$

These results and the lack of differences between the control and test group on other parts of the data indicate that the marble drop experiment did not influence the results of the colored trails experiment.

## 4.4 Interpersonal Reactivity Index

From the personality questionnaire, the scores for four scales were calculated per participant: the perspective taking scale (PT), the fantasy scale (FS), the empathic concern scale (EC) and the personal distress scale (PD). The results are shown in Figure 4.14.

According to [33], the average scores on all four scales differ between men and women. For our participants this difference was only found for EC ( $Z = -2.28, p = .023$ ). Therefore, tests on PT, FS and PD were conducted on all participants instead of on men and women separately.

Correlation tests for each of the four scales and the score on the colored trails experiment did not yield any significant results. Correlation tests between the confidence of the participant using  $ToM_0$ ,  $ToM_1$  or  $ToM_2$  during the colored trails experiment and all of the four scales did not yield any significant results for most tests, except for the confidence in the use of  $ToM_2$  and PD ( $r(25) = .49, p = .010$ ) and for the confidence in the use of  $ToM_1$  and FS ( $r(25) = .50, p = .0081$ ). In both cases a higher score on PD/FS tended to result in a higher confidence. Overall, however, personality traits do not seem to influence the performance on colored trails.

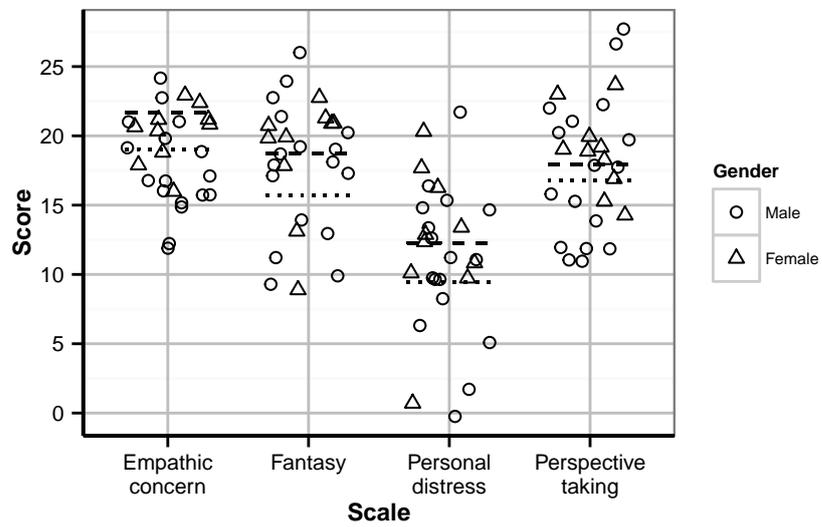


Figure 4.14: Scores of the participants on the Interpersonal Reactivity Index questionnaire, per scale. The dashed line depicts the average scores for women, the dotted line the average scores for men, as found in [33].

## Chapter 5

# Discussion

In this thesis we investigated whether people use theory of mind in the colored trails negotiation game. The question was: if they use theory of mind, which order of theory of mind would they use? And will this be influenced by the order of theory of mind used by the opponent player? Furthermore we investigated whether training of the use of theory of mind via the marble drop game would help in the use of theory of mind in the colored trails game. Finally it was tested whether personality aspects of the participants influenced their results on colored trails.

### 5.1 Marble Drop

With an overall accuracy of 99% on the zero-order marble drop games and 93% on the first-order marble drop games, the participants performed well. However, on the second-order marble drop games the overall accuracy was only 67%, which is low compared to previous research [29]. The questionnaire about the marble drop experiment showed that most participants did report to have reasoned about the opponent, even the ones that had a low accuracy. This shows that the participants realized that they had to take into account the opponent player, but apparently they did not exactly know how to do that yet. This might have been caused by the small number of second-order games they played. Both in [29] and in the current study, the participants played four zero-order, eight first-order and eight second-order level marble drop games. For the current study, this was the end of the marble drop part. In [29] this was just the training phase: the participants then continued with another 32 second-order level marble drop games, of which the accuracy was 90%-94%. Apparently, participants need more training before they master the second-order level marble drop games.

There were two games in particular on which the participants scored poorly. The pay-off structures of those games are shown in Figure 4.3. These were the only games where the black marble had to end up in the second bin from the left. Why did the participants score so poorly on these pay-off structures? In the next paragraphs the participant is the blue player and the computer the

orange player.

Obviously, when the participants would have used  $ToM_2$  (correctly), they would have ended up at the correct bin. So apparently the participants were not using  $ToM_2$  or, if they did, they made a mistake.

If they would have been using  $ToM_0$ , they would have started with opening the right trapdoor, since the darkest blue marble was in the second bin from the left. This would have resulted in the correct answer, because the computer would have opened his left trapdoor. So most of the participants apparently were not using  $ToM_0$  either.

If one would use  $ToM_1$  for these games, one would think that the opponent would open his right trapdoor in order to get closer to the darkest orange marble reachable at that point. The blue marble in the two right most bins are both lighter than the blue marble in the left most bin. Therefore a  $ToM_1$  player would decide to open the left trapdoor at the first decision point. This would result in the wrong answer. So this makes it likely that participants were using  $ToM_1$  to answer these questions. If the low accuracies were caused by the use of  $ToM_1$ , then the next question is: why did they perform better on other games where the use of  $ToM_1$  would have led to the wrong answer as well? So this is not a sufficient explanation of the poor accuracy.

The two games with the lowest accuracies were the only games where the correct answer was to let the black marble end up in the second bin from the left and this is what might explain the poor accuracy. For marble drop, if common knowledge of rationality is assumed, using *backwards induction* leads to the correct answer [36]. With backwards induction one starts at the final set of trapdoors, comparing the marbles in each bin, working your way up. In [37], eye-fixations were recorded of participants playing second-order level marble drop games. Their data showed that participants used *forward reasoning plus backtracking* more often than backward induction. In forward reasoning plus backtracking one starts reasoning at the first set of trapdoors, to decide in which direction the highest pay-off lies. This is followed by backward reasoning to check whether that bin is reachable.

This means that when someone uses forward reasoning with backtracking, many reasoning steps are needed when the highest attainable pay-off is in the second bin. The more reasoning steps needed, the higher the chance that one will make a mistake. This is a plausible explanation for the low accuracies on these two games. Unfortunately, in our experiment, no eye-fixation data was collected nor reaction times, which also tell something about the used strategy by the participants [38], to confirm this.

## 5.2 Training

The marble drop experiment was part of our research so we could see whether training with marble drop would stimulate the use of theory of mind in colored trails. The control group only played zero-order level games, the test group zero-, first- and second-order level games. No evidence was found for an influence of marble drop on the scores on colored trails nor on the most likely used order of theory of mind by the participants in the colored trails experiment.

A possible explanation for this is the low accuracy on the second-order marble drop games (67%). If the training phase with marble drop would be extended, the accuracies would probably get higher, as was found in [29]. The reason for this, as is described in [39], is probably that the participants started with simple strategies which they adjusted when necessary. To adjust the  $ToM_1$  strategy to a  $ToM_2$  strategy might take more time than the eight games the participants played. Further research will have to show if playing more second-order level marble drop games will support a learning effect for colored trails.

## 5.3 Colored Trails

### 5.3.1 Scores

All participants scored far better than the absolute minimum score possible and also better than the ‘withdrawal score’ (the score the participants would have gotten if they would have withdrawn from every negotiation). This shows that the participants understood the game quite well and were not just ‘fooling around’. However, there was also room for improvement, because no one came close to the Pareto front. Be that as it may, the colored trails game is quite complex and the participants had no experience with it, so it was not expected that the participants would play optimally.

The participants’ scores per game did, in general, increase when the experiment progressed, so it seems as if they learned during the experiment. It is, however, not that easy to jump to this conclusion. Every game was different and sometimes it was harder to get a high score (when participant and agent both needed the same chips) than in other games. Furthermore one could never exactly know what the agent would do. This made it difficult to make a perfect prediction, especially since the agent switched its behavior twice by changing the order of theory of mind it used. Thirdly, one was always dependent on the behavior of the opponent player. These are all aspects which influence the learning effect. In order to confirm that a learning effect occurs in the colored trails experiment, those aspects should be controlled for.

The score per colored trails game was influenced by the starting player. Both the participants and the agents scored better on games in which the agent started the negotiation. There was one exception: in the block against the  $ToM_2$  agent, the participants scored better when they started themselves. The influence of the starting player presumably came from the initial offer this player made [40]. Both the participants and the agents generally ascribed high scores to themselves with their first offers. The participants did that just as well as the  $ToM_0$  and  $ToM_1$  agents did, but the  $ToM_2$  outclassed them on this. The  $ToM_0$  and  $ToM_1$  agents ascribed higher scores to the participants with their first offers than vice versa. The  $ToM_2$  agents and the participants ascribed similar scores to each other with their first offers.

So the results from the blocks with the  $ToM_0$  and  $ToM_1$  agents were similar: the final scores were better when the agent started the negotiation; with their first offers both the participants and agents ascribed similar (high) scores to themselves; and the agents’ first offers ascribed higher scores to the participants

than vice versa. How can this be explained? The last point might be the result of the difference between computer agents and human beings: agents have much more computing power and can ‘reason’ more optimally. This might have helped them in coming up with a better, ‘fairer’, distribution than the participants were able to (in the first round of a negotiation).

In the block against the  $ToM_2$  agents, the participants ascribed similar scores to the agents with their first offers as the agents ascribed to them. It seems that the participants compensated for the superior computing power of the agents. A difference between the blocks with the  $ToM_0$  and  $ToM_1$  agents on the one hand and the block with the  $ToM_2$  agent on the other hand is that in the latter the participants used more  $ToM_2$  (see also Section 5.3.2). This might have helped them in creating better initial offers for the opponent player.

In the  $ToM_2$  block, the final score of the participant was better when the participant started. Why was there no benefit from the agent in this situation? How can a higher score from the initial offer result in a higher score at the end of a negotiation in the first place? This might be a result of the anchoring principle [41]: people tend to use the first piece of information as a reference point or ‘anchor’. They make new estimates based on the anchor, adjusting away from it. Different starting points yield different estimates and might thus result in different scores at the end of a negotiation.

In this case the first offer is actually not the first piece of information, the initial distribution is. Nonetheless, the first offer is an important piece of information: it is the first information about how much the other player is willing to give to you. Therefore it can still be seen as an anchor. The anchors that both the agents and participants created for their opponent player were different: the anchor the participant received from the agent was higher than the anchor the participant gave to the agent, except in the  $ToM_2$  block where they were similar. A ‘higher’ reference point for the participant might have resulted in a more ambitious participant, which could have led to higher scores at the end of the negotiation. In the  $ToM_2$  block both players ascribed similar scores to each other with their initial offers; this might explain the absence of the ‘extra motivation’ for the participant from the agent, resulting in better scores for a player when that player started the negotiation. Another explanation is that the games in the  $ToM_2$  block might not have been as suitable for creating the results in block  $ToM_0$  and  $ToM_1$  due to differences in the available chips and their distribution.

In the colored trails part, the participants unknowingly played against three types of agents ( $ToM_0$ ,  $ToM_1$ ,  $ToM_2$ ). The participant’s score on a colored trails game was influenced by the order of theory of mind that the agent used. It was lower when the participant played against the  $ToM_1$  agent. This might be due to the nature of the  $ToM_1$  agent. This type of agent attributes meaning to the offers of the participant. Based on those offers, it tries to determine the most likely goal location of the participant. To decide what the best offer is, it tries to evaluate its offer from the perspective of the participant. This way, it can manipulate the opponent: the agent can offer something which it thinks the opponent will not accept, but will change the opponent’s beliefs in such a way that the opponent will come up with an even better offer. This type of agent is

thus quite devious and not very generous, which might explain the low scores in this block.

In spite of this, the agent's scores were also lowest when it was using  $ToM_1$ . The social welfare in this block, measured as the sum of both players' scores, was therefore low: the cooperation between the players was not very good in this block. A possible explanation for the agent's low scores is that it assumed that the opponent player was using  $ToM_0$ . The results however, show that it is more likely that the participants were using either  $ToM_1$  or  $ToM_2$  during this block. So the  $ToM_1$  agent's predictions about the participant might have been wrong quite often, which could explain its low scores during this block.

The  $ToM_2$  agent is a more efficient version of the  $ToM_1$  agent and can also try to manipulate the other player. Nevertheless, the participants' scores were much better against this type of agent. The  $ToM_2$  agent assumes that the other player will try to interpret its offers and can use this to signal, for example, its goal location. This can speed up the negotiation process, leaving the participant with more time for other processes such as creating a good offer. This might have compensated for the manipulations, resulting in higher scores for the participants. The scores for the agent were also higher in the  $ToM_2$  block, compared to the results from the  $ToM_1$  block, and therefore the social welfare was also higher. The  $ToM_2$  agent might have scored better than the  $ToM_1$  agent because its predictions were more accurate (it was more likely that the participant used  $ToM_1$  than  $ToM_0$ ).

The participant's scores in the block against the  $ToM_0$  agent were similar to the scores against the  $ToM_2$  agent and the social welfare was similar as well, but only in the  $ToM_0$  block scores higher than 100 points for the participant occurred. The  $ToM_0$  agent bases its beliefs on the behavior of the participant. For example: when the participant rejects an offer, the agent believes that the participant will also reject offers where the participant gets even less chips. Therefore this type of agent is more willing to create generous offers for the participant, as long as the alternative (the initial distribution, via withdrawal) is worse. This can sometimes result in very high scores for the participant.

The data also showed that high scores for one player were usually obtained when the other player accepted the offer. Apparently the players had to create good offers themselves to get high scores. It was not very likely that the other player would create an offer that would help one reach the goal, although this did occur. This might indicate that both players had problems with figuring out where the opponent's goal location was. Another explanation is that a player did not want to give away all chips the other player needed, because this would lead to a low score for the player itself.

The more time a participant spent on the colored trails experiment, the higher his/her overall score tended to be. Taking more time to evaluate the situation seemed to have helped the participants to come up with better proposals, which is what one would expect.

### 5.3.2 Used Orders of Theory of Mind by the Participant

The evaluation model of the agent showed that it is most likely that the participants were using  $ToM_1$  or  $ToM_2$ , as opposed to  $ToM_0$ . This finding was supported by the comments of the participants, where many people stated that they looked at which chips the agent was consistently keeping / giving away during a negotiation, which was for them an indication of what the agent wanted. The use of  $ToM_0$  and  $ToM_1$  did not differ between the three blocks (of which the existence was unbeknown to the participants), but the use of  $ToM_2$  did: participants used more  $ToM_2$  when they played against the  $ToM_2$  agent.

It is interesting to see that an opponent who uses a higher order of theory of mind stimulated the use of higher orders of theory of mind in the participants. If a participant would play optimally, then that is also what one would expect. Therefore this is another indication that the participants were actually using theory of mind.

There are, on the other hand, also other strategies that the participants might have used. The agent's confidences of the participants using either  $ToM_0$ ,  $ToM_1$  or  $ToM_2$  did not add up to 100%, which means that there is uncertainty in the system. Examples of other strategies which the participants could have used are: 1) only look at possible distributions which lead to your goal tile and try those (repeatedly); 2) try to determine what the odds are that the opponent player needs a chip with a certain texture based on how often that texture is present on the board.

Besides the other strategies, part of the uncertainty can also be explained by theory of mind itself. As reference points, the agent used the *optimal* offers from a  $ToM_0$  /  $ToM_1$  /  $ToM_2$  player. Because of the time limit and the complexity of the game, it is not likely that the participants would always come up with the optimal offer for the order of theory of mind they were using. This may have led to wrongful decreases in the confidence that the participant was using that order of theory of mind. Therefore the results on the used orders of theory of mind by the participants are an indication and not an absolute certainty.

Another reason why more research is needed to confirm our results is the following. In the block with the  $ToM_2$  agent the confidences that the participants were using  $ToM_2$  was highest. This could indicate that the participants most likely mainly used  $ToM_2$  and less  $ToM_1$  in this block. Another explanation, however, could be that the offers from the  $ToM_2$  agent were of such a form that they elicited offers from the participants which were easier to evaluate by the evaluation system than the participants' offers in the other two blocks. Therefore the confidences in  $ToM_2$  could be higher in this block than in the other blocks, even if, possibly, the participants used  $ToM_2$  in all of them.

A significant relation was found between the scores on the colored trails experiment and the confidence that participants used  $ToM_0$ . It is not clear what this means. The likeliness that a participant actually used  $ToM_0$  in any of the three blocks was low. Further research is needed to interpret this result and to confirm that it is not just a coincidence.

Since more reasoning steps are needed for using higher orders of theory of mind, one would expect that the participants needed more time against the  $ToM_2$  agent, because the participants most likely used more  $ToM_2$  in that block

themselves. Surprisingly, the participants needed the same amount of time against the  $ToM_2$  agent as against the  $ToM_0$  agent. Against the  $ToM_1$  agent, they even needed more time. A possible explanation for the fact that no longer reaction times were found in the block with the  $ToM_2$  agent, is that this type of agent tried to create offers in such a way that the agent's goal location was easy to interpret for the participant. This might have reduced the time needed to form good offers by the participant, compensating for the the extra time needed for the extra reasoning steps. That the participants needed more time against the  $ToM_1$  agent, might be because this type of agent proposed distributions which were not good for the participants (low scores in that block). It might have taken the participants extra time to overcome this.

### 5.3.3 Comparison with Previous Research

The current study is based on the study in [12], where the three types of computer agents used in the current study played colored trails against each other. In [12] it was shown that two  $ToM_0$  players do not negotiate very well with each other. The results of the current study show that it is not very likely that the participants used  $ToM_0$ , therefore no comparisons between these results can be made. In the blocks with the  $ToM_0$  and  $ToM_1$  agent, the confidence that the participants used  $ToM_1$  is similar to the confidence that the participant used  $ToM_2$ . Therefore no good comparisons could be made with these data either, since one would not know whether one would be comparing the agent with a  $ToM_1$  or  $ToM_2$  player. In the block with the  $ToM_2$  agent however, a clear distinction in the confidences was found: it was most likely that the participants used  $ToM_2$ . In [12] two  $ToM_2$  computer agents could negotiate with each other well. The same was found in our study: both the agents and the participants scored well in this block, as expected.

The social welfare results in [12] regarding two  $ToM_2$  agents showed that the social welfare was high in this setting. In the current study, the highest social welfare was found in the blocks with the  $ToM_2$  agent and  $ToM_0$  agent. Since the scoring system used in [12] was different than the scoring system used in the current study, it cannot be compared whether the level of social welfare is similar. Since it was not clear to tell which order of theory of mind the participants used in the  $ToM_0$  agent block, no further comparisons regarding the social welfare in this block could be made.

## 5.4 Interpersonal Reactivity Index

In the colored trails game, and with theory of mind in general, it is important to be able to view things from the perspective of someone else. The interpersonal reactivity index (IRI) tests how high someone's tendency is to adopt the point-of-view of someone else. The expected relation with perspective taking, which is necessary for the use of theory of mind, and the results on colored trails was not found. The only significant results were a higher confidence in the use of  $ToM_2$  for a higher score on the personal distress scale and a higher confidence in the use of  $ToM_1$  for a higher score on the fantasy scale.

Regarding the personal distress scale: the time limit that was present in the experiment could have put participants under pressure. There is, however, no clear reason why this would increase the use of  $ToM_2$ , especially since reasoning in this order of theory of mind takes more time than the other two orders. Regarding the fantasy scale: people who can more easily transpose themselves imaginatively into fictitious characters might be better at engaging in the colored trails game and with the agent. This might have supported the (correct) use of  $ToM_1$ . The last part of the IRI was a test for ‘empathic concern’: to what extent does someone have feelings like sympathy and concern for (unfortunate) others. A person who scores high on this score will probably be less greedy in the colored trails game, resulting in lower scores. No proof was found for this.

## 5.5 Research Questions

In Chapter 1, one main research question and three subquestions were formulated. Based on the results found in the current study, those questions will be answered.

### 5.5.1 Subquestion 1: The Influence of the Opponent’s Order of Theory of Mind

The first subquestion was formulated as follows:

*How is the use of theory of mind in the colored trails negotiation game influenced by the order of theory of mind the opponent uses?*

People use more second-order theory of mind in the colored trails negotiation game when the opponent uses second-order theory of mind as well. In other situations people use either first- or second-order theory of mind.

### 5.5.2 Subquestion 2: The Influence of Training

The second subquestion runs as follows:

*What is the influence of training with the marble drop game on the use of theory of mind in the colored trails negotiation game?*

There is no training effect from playing marble drop when at the end of the training phase, people have not yet mastered the use of second-order of theory of mind.

### 5.5.3 Subquestion 3: The Influence of Personality Traits regarding Empathy

The third subquestion was:

*What is the influence of personality traits regarding empathy on the performance on the colored trails negotiation game?*

People who reported that they are inclined to take into account someone else's perspective in daily life, are not better at using (higher orders of) theory of mind in the colored trails negotiation game than others.

#### 5.5.4 Main Question: Participant's Use of Theory of Mind

The main question this thesis tried to answer, is:

*Do people use theory of mind when playing a negotiation game, for which the use of second order theory of mind has proven to be useful, against a computer agent and if so, what order of theory of mind do they use?*

Yes, people use theory of mind in order to play the colored trails negotiation game. They mostly use first- and second-order theory of mind, but are inclined to predominantly use second-order theory of mind when they are (unknowingly) facing an opponent who uses second-order theory of mind.

## 5.6 Future Research

### 5.6.1 Training or Transfer

As mentioned at the start of this chapter, no effect was found of the marble drop experiment on the colored trails experiment. Related to this is the method we used. In the current study, one half of the participants only played zero-order level marble drop games, the other half also played first- and second-order level marble drop games before continuing to the colored trails experiment. No differences were found between the two groups in the colored trails experiment. This seems to indicate that the participants did not learn from the marble drop experiment, but this is not necessarily true.

To explain this, we will first introduce a new concept. Learning is an important part of our lives, since we need to adapt to new situations, products and procedures all the time. Starting from scratch in every new situation would be very inefficient, so using knowledge from other situations would be helpful in mastering something new. Being able to use 'old' knowledge in new situations is called *transfer* of learning.

The concept of transfer was first mentioned already more than a century ago in an article by Woodworth and Thorndike [42]. They suggested that practicing one task would only benefit another task if the tasks share identical elements. This is the *Theory of Identical Elements*.

Based on the Theory of Identical Elements, one would expect that more transfer occurs when more elements are shared between tasks. The similarity between tasks is called transfer *distance* or transfer *similarity* [43]. How to determine this distance?

One framework that could be used describes the cognitive representation of a task in terms of the user's Goals, Operations, Methods, and Selection rules. It is called the *GOMS model* and was created by Card, Moran, and Newell in 1983 (as cited by [44, 45, 46]).

Based on this, the *cognitive complexity theory* (CCT) was developed [47]. In this theory, the items described in a GOMS analysis are represented as a set of production rules. Those production rules, the production system, are used to model user behavior. Each rule has a condition and an action, which will be carried out when the condition is satisfied. The training time depends on the number of unique rules that one has to learn to be able to perform a task. Transfer between tasks depends on the number of shared production rules (the 'elements' from the Theory of Identical Elements) relative to the total number of productions.

A similar more recent system which goes one step deeper is ACTransfer [48]. In this system, the overlap between tasks can be determined by comparing the so called *primitive information processing elements* (PRIMs) of those tasks (the basic elements of cognitive skills) and how these are combined into productions. The PRIMs themselves are also production rules, but they lack a control component. This means that they do not control in which order they are carried out. When productions (consisting of several PRIMs) are created in one task which can also be used in another task, transfer can occur.

All above-mentioned systems are based on the idea that, to encourage transfer to occur, it is important to look at the underlying structures of tasks. If tasks only look similar superficially, there is little chance for transfer to occur.

The training effect that we hoped to see in the current experiment, was actually a form of transfer. Therefore the above-mentioned points might explain the missing training effect. Maybe marble drop and colored trails are not similar in such a way that the knowledge gained in the marble drop experiment can transfer to the colored trails experiment. Even though the use of theory of mind is useful in both, making the underlying structure similar, it is obvious that the superficial appearances of both games are very different and only colored trails is a mixed-motive game, having a cooperation element besides the competitive element. Despite these differences, 'far' transfer could still occur [49].

Marble drop is a game of complete and perfect information which is surveyable: even in the second-order theory of mind form, it is possible to deduce the winning strategy (for both players). With perfect reasoning one can know exactly what will happen, when one knows that the opponent reasons perfectly as well. In colored trails there is incomplete information so, at least at the start of the negotiation, even if one reasons perfectly, it is not possible to know for sure what the other player will do. Therefore it is harder to apply theory of mind and besides that, there are also other strategies that can be used (e.g. what are the chances that the other player needs a chip with this texture?). Maybe general problem solving was the most difficult part of colored trails, instead of the use of theory of mind. A plausible explanation for the fact that no transfer was found is therefore that the games are not similar enough in the underlying structure, but further research using the above-mentioned methods will have to show whether that is the case or not. Another explanation might

be that the difficulty levels of both games were too different, which has proven to be of influence on transfer as well [50].

Another reason for the lack of training effect, i.e. transfer, could be the low accuracy on the second-order level marble drop games. Maybe the participants did not learn enough from the marble drop experiment. In our case the training phase on second-order theory of mind consisted of only eight questions. This might have helped the participants to remember and activate their knowledge on second-order theory of mind, but it was most likely too short to learn a new skill. To test for this, a pre-test post-test design should be used in future research.

Even when future research shows that the participants will learn more from the marble drop experiment when more games are played and marble drop and colored trails are similar enough for transfer to occur, this is no guarantee for better performance on colored trails. There is evidence that when one becomes an expert in one situation, transfer decreases [51]. The participants might create a strategy for the second-order level marble drop games which is specific for that game. For example, if the participants use backward induction on these games, they might forget that they are actually reasoning about the intentions of the opponent player. Since colored trails is a much more complex game than marble drop, this strategy will not work for colored trails. So being an expert on marble drop might even diminish the transfer. Therefore, future research will have to be conducted carefully.

### 5.6.2 Adjusting the Colored Trails Set-up

Our data showed a positive relation between the time spent on the colored trails experiment and the overall score on this part. Furthermore, some participants stated in their comments that they did not have enough time to reason about the opponent. An interesting adjustment for future work would therefore be to extend the time limit to, for example, three minutes. This might not just increase the performance, but it could also increase the use of second-order theory of mind. With more time, people could more easily find out what all possible paths are and thus better deduce the opponent's goal location. In the current experiment, participants reported that they did not very often try to find out where the opponent's goal location was. It is unclear whether this was due to the short time limit, so an extended time limit will not necessarily increase the search for the opponent's goal.

### 5.6.3 Colored Trails as Negotiation Practice

The current study showed that playing colored trails against a computer agent who used second-order theory of mind, increased the use of second-order theory of mind in the participants. Playing colored trails could therefore be a useful tool to practice the use of higher orders of theory of mind. New experiments are needed to see if playing colored trails against a second-order theory of mind computer agent is beneficial for the use of theory of mind and if it will increase performance in negotiations with incomplete information. A useful adaptation could be to include conversations with a well-informed person on theory of mind

and negotiations. This person could indicate, at the end of each game, how theory of mind should have been applied and how it would have been beneficial for the negotiation. This could make the participants even more aware of the power of theory of mind.

## Chapter 6

# Bibliography

- [1] Premack, D., Woodruff, G.: Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences* **1**(4) (1978) 515–526
- [2] Perner, J.: Theory of mind. In Bennett, M., ed.: *Developmental Psychology: Achievements and Prospects*. Psychology Press (1999) 205–230
- [3] Wimmer, H., Perner, J.: Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception. *Cognition* **13**(1) (1983) 103–128
- [4] Perner, J., Wimmer, H.: “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) (1985) 437–471
- [5] Wellman, H.M., Cross, D., Watson, J.: Meta-analysis of theory-of-mind development: The truth about false belief. *Child Development* **72**(3) (2001) 655–684
- [6] Onishi, K.H., Baillargeon, R.: Do 15-month-old infants understand false beliefs? *Science* **308**(5719) (2005) 255–258
- [7] Miller, S.A.: Children’s understanding of second-order mental states. *Psychological Bulletin* **135**(5) (2009) 749–773
- [8] Hedden, T., Zhang, J.: What do you think I think you think?: Strategic reasoning in matrix games. *Cognition* **85**(1) (2002) 1–36
- [9] Verbrugge, R.: Logic and social cognition: The facts matter, and so do computational models. *Journal of Philosophical Logic* **38**(6) (2009) 649–680
- [10] De Weerd, H., Verbrugge, R., Verheij, B.: How much does it help to know what she knows you know? An agent-based simulation study. *Artificial Intelligence* **199-200** (2013) 67–92
- [11] De Weerd, H., Verbrugge, R., Verheij, B.: Agent-based models for higher-order theory of mind. In Kamiński, B., Koloch, G., eds.: *Advances in Social*

- Simulation, Proceedings of the 9th Conference of the European Social Simulation Association. *Advances in Intelligent Systems and Computing* 229. (2013) 213–224
- [12] De Weerd, H., Verbrugge, R., Verheij, B.: Higher-order theory of mind in negotiations under incomplete information. In: *Proceedings of the 16th International Conference on Principles and Practice of Multi-Agent Systems, Lecture Notes in Artificial Intelligence* 8291. (2013) 101–116
- [13] Baron-Cohen, S.: The autistic child’s theory of mind: A case of specific developmental delay. *Journal of Child Psychology and Psychiatry* **30**(2) (1989) 285–297
- [14] Happé, F.: An advanced test of theory of mind: Understanding of story characters’ thoughts and feelings by able autistic, mentally handicapped, and normal children and adults. *Journal of Autism and Developmental Disorders* **24**(2) (1994) 129–154
- [15] Purdy, J.E., Domjan, M.: Tactics in theory of mind research. *Behavioral and Brain Sciences* **21** (1998) 129–130
- [16] Penn, D.C., Povinelli, D.J.: On the lack of evidence that non-human animals possess anything remotely resembling a ‘theory of mind’. *Philosophical Transactions of the Royal Society B: Biological Sciences* **362**(1480) (2007) 731–744
- [17] Oxoby, R.J., McLeish, K.N.: Sequential decision and strategy vector methods in ultimatum bargaining: Evidence on the strength of other-regarding behavior. *Economics Letters* **84**(3) (2004) 399–405
- [18] Fehr, E., Schmidt, K.M.: A theory of fairness, competition, and cooperation. *The Quarterly Journal of Economics* **114**(3) (1999) 817–868
- [19] Camerer, C.F., Ho, T.H., Chong, J.K.: A cognitive hierarchy model of games. *The Quarterly Journal of Economics* **119**(3) (2004) 861–898
- [20] Peled, N., Gal, Y., Kraus, S.: A study of computational and human strategies in revelation games. *Journal of Autonomous Agents and Multi-Agent Systems* (2014) 1–25
- [21] Haim, G., Gal, Y., Gelfand, M., Kraus, S.: A cultural sensitive agent for human-computer negotiation. In: *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems - Volume 1. AAMAS ’12, Richland, SC, International Foundation for Autonomous Agents and Multiagent Systems* (2012) 451–458
- [22] Oshrat, Y., Lin, R., Kraus, S.: Facing the challenge of human-agent negotiations via effective general opponent modeling. In: *Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems - Volume 1. AAMAS ’09, Richland, SC, International Foundation for Autonomous Agents and Multiagent Systems* (2009) 377–384

- [23] Rosenfeld, A., Zuckerman, I., Halevi, E.S., Drein, O., Kraus, S.: NegoChat: A chat-based negotiation agent. In: Proceedings of the 2014 International Conference on Autonomous Agents and Multi-agent Systems. AAMAS '14, Richland, SC, International Foundation for Autonomous Agents and Multiagent Systems (2014) 525–532
- [24] Ficici, S.G., Pfeffer, A.: Modeling how humans reason about others with partial information. In: Proceedings of the 7th International Joint Conference on Autonomous Agents and Multiagent Systems - Volume 1. AAMAS '08, Richland, SC, International Foundation for Autonomous Agents and Multiagent Systems (2008) 315–322
- [25] Lueg, C., Pfeifer, R.: Cognition, situatedness, and situated design. In: Second International Conference on Cognitive Technology, IEEE Computer Society (1997) 124–135
- [26] Gal, Y., Grosz, B., Pfeffer, A., Shieber, S., Allain, A.: The influence of task contexts on the decision-making of humans and computers. In Kokinov, B., Richardson, D., Roth-Berghofer, T., Vieu, L., eds.: Modeling and Using Context. Volume 4635 of Lecture Notes in Computer Science. Springer Berlin Heidelberg (2007) 206–219
- [27] Grosz, B.J., Kraus, S., Talman, S., Stossel, B., Havlin, M.: The influence of social dependencies on decision-making: initial investigations with a new game. In: Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems - Volume 2. AAMAS '04, Washington, DC, USA, IEEE Computer Society (2004) 782–789
- [28] Gal, Y., Grosz, B., Kraus, S., Pfeffer, A., Shieber, S.: Agent decision-making in open mixed networks. *Artificial Intelligence* **174**(18) (2010) 1460–1480
- [29] Meijering, B., Van Maanen, L., Van Rijn, H., Verbrugge, R.: The facilitative effect of context on second-order social reasoning. In Catrambone, R., Ohlsson, S., eds.: Cognitive Science Society, Cognitive Science Society (2010) 1423–1429
- [30] Nadler, J., Thompson, L., Van Boven, L.: Learning negotiation skills: Four models of knowledge creation and transfer. *Management Science* **49**(4) (2003) 529–540
- [31] Loewenstein, J., Thompson, L., Gentner, D.: Analogical encoding facilitates knowledge transfer in negotiation. *Psychonomic Bulletin & Review* **6**(4) (1999) 586–597
- [32] Gentner, D., Loewenstein, J., Thompson, L.: Learning and transfer: A general role for analogical encoding. *Journal of Educational Psychology* **95**(2) (2003) 393–408
- [33] Davis, M.H.: A multidimensional approach to individual differences in empathy. *JSAS Catalog of Selected Documents in Psychology* **10** (1980) 85–103

- [34] Völlm, B.A., Taylor, A.N.W., Richardson, P., Corcoran, R., Stirling, J., McKie, S., Deakin, J.F.W., Elliott, R.: Neuronal correlates of theory of mind and empathy: A functional magnetic resonance imaging study in a nonverbal task. *NeuroImage* **29**(1) (2006) 90–98
- [35] Schulte-Rüther, M., Markowitsch, H., Fink, G., Piefke, M.: Mirror neuron and theory of mind mechanisms involved in face-to-face interactions: A functional magnetic resonance imaging approach to empathy. *Journal of Cognitive Neuroscience* **19**(8) (2007) 1354–1372
- [36] Aumann, R.J.: Backward induction and common knowledge of rationality. *Games and Economic Behavior* **8**(1) (1995) 6–19
- [37] Meijering, B., Van Rijn, H., Taatgen, N.A., Verbrugge, R.: What eye movements can tell about theory of mind in a strategic game. *PLoS ONE* **7**(9) (2012) e45961+
- [38] Bergwerff, G., Meijering, B., Szymanik, J., Verbrugge, R., Wierda, S.M.: Computational and algorithmic models of strategies in turn-based games. In: *Proceedings of the 36th Annual Conference of the Cognitive Science Society*. (2014)
- [39] Meijering, B., Taatgen, N.A., Van Rijn, H., Verbrugge, R.: Modeling inference of mental states: As simple as possible, as complex as necessary. *Interaction Studies* (in press)
- [40] Van Poucke, D., Buelens, M.: Predicting the outcome of a two-party price negotiation: Contribution of reservation price, aspiration price and opening offer. *Journal of Economic Psychology* **23**(1) (2002) 67–76
- [41] Tversky, A., Kahneman, D.: Judgment under uncertainty: Heuristics and biases. *Science* **185**(4157) (1974) 1124–1131
- [42] Woodworth, R.S., Thorndike, E.L.: The influence of improvement in one mental function upon the efficiency of other functions. (I). *Psychological Review* **8**(3) (1901) 247–261
- [43] Butterfield, E., Nelson, G.: Theory and practice of teaching for transfer. *Educational Technology Research and Development* **37**(3) (1989) 5–38
- [44] Polson, P.G., Bovair, S., Kieras, D.: Transfer between text editors. In: *Proceedings of the SIGCHI/GI Conference on Human Factors in Computing Systems and Graphics Interface*. CHI '87, New York, NY, USA, ACM (1987) 27–32
- [45] Singley, M.K., Anderson, J.R.: The transfer of text-editing skill. *International Journal of Man-Machine Studies* **22**(4) (1985) 403–423
- [46] Ziegler, J.E., Hoppe, H.U., Fahnrich, K.P.: Learning and transfer for text and graphics editing with a direct manipulation interface. *SIGCHI Bulletin* **17**(4) (1986) 72–77

- [47] Polson, P.G., Kieras, D.E.: A quantitative model of the learning and performance of text editing knowledge. *SIGCHI Bulletin* **16**(4) (1985) 207–212
- [48] Taatgen, N.A.: The nature and transfer of cognitive skills. *Psychological Review* **120**(3) (2013) 439–471
- [49] Barnett, S.M., Ceci, S.J.: When and where do we apply what we learn?: A taxonomy for far transfer. *Psychological Bulletin* **128**(4) (2002) 612–637
- [50] Wierda, S., Arslan, B.: Modeling theory of mind in ACTransfer. In Szymanik, J., Verbrugge, R., eds.: *Proceedings of the Second Workshop Reasoning About Other Minds: Logical and Cognitive Perspectives*, co-located with *Advances in Modal Logic*. (2014) 40–44
- [51] Taatgen, N.A.: Diminishing return in transfer: A PRIM model of the Frensch (1991) Arithmetic Experiment. In: *Proceedings of the 12th International Conference on Cognitive Modelling*. (2013) 29–34



# Appendix A

## Marble Drop

The numbers used in the Tables A.2-A.4 (e.g. ‘marble 3’) correspond to the numbers in Table A.1.

### A.1 Colors

Table A.1: The colors used in marble drop. The shades are on a scale of 1-4, where 1 is the lightest shade and 4 the darkest.

Color	HTML-code
Blue 1	e1e9ed
Blue 2	a2c2f4
Blue 3	5380d1
Blue 4	1111d5
Orange 1	fff4e6
Orange 2	ffd7a6
Orange 3	ffaf4d
Orange 4	ff8c00

### A.2 Structures

Table A.2: The settings used in the zero-order level of marble drop. ‘marble 1’ indicates a light marble, ‘marble 4’ a dark one.

Bin 1 participant	Bin 1 computer	Bin 2 participant	Bin 2 computer	Optimal reachable bin participant
marble 1	marble 1	marble 4	marble 4	bin 2
marble 1	marble 4	marble 4	marble 1	bin 2
marble 4	marble 1	marble 1	marble 4	bin 1
marble 4	marble 4	marble 1	marble 1	bin 1

Table A.3: The settings used in the first-order level of marble drop. ‘marble 2’ indicates the lightest marble, ‘marble 4’ the darkest one.

<b>Bin 1</b>	<b>Bin 1</b>	<b>Bin 2</b>	<b>Bin 2</b>	<b>Bin 3</b>	<b>Bin 3</b>	<b>Optimal</b>
<b>participant</b>	<b>computer</b>	<b>participant</b>	<b>computer</b>	<b>participant</b>	<b>computer</b>	<b>reachabe</b>
						<b>bin participant</b>
marble 3	marble 2	marble 2	marble 4	marble 4	marble 3	bin 1
marble 3	marble 4	marble 2	marble 3	marble 4	marble 2	bin 1
marble 3	marble 3	marble 4	marble 2	marble 2	marble 4	bin 1
marble 3	marble 2	marble 4	marble 3	marble 2	marble 4	bin 1
marble 3	marble 4	marble 4	marble 2	marble 2	marble 3	bin 1
marble 3	marble 2	marble 4	marble 4	marble 2	marble 3	bin 2
marble 3	marble 4	marble 4	marble 3	marble 2	marble 3	bin 2
marble 3	marble 4	marble 4	marble 3	marble 2	marble 2	bin 2
marble 3	marble 4	marble 2	marble 2	marble 4	marble 3	bin 3

Table A.4: The settings used in the second-order level of marble drop. ‘marble 1’ indicates the lightest marble, ‘marble 4’ the darkest one. (‘Part.’ is participant, ‘comp.’ is computer.)

<b>Bin 1</b>	<b>Bin 1</b>	<b>Bin 2</b>	<b>Bin 2</b>	<b>Bin 3</b>	<b>Bin 3</b>	<b>Bin 4</b>	<b>Bin 4</b>	<b>Bin 4</b>	<b>Bin 4</b>	<b>Optimal reachable</b>
<b>part.</b>	<b>comp.</b>	<b>bin participant</b>								
marble 3	marble 2	marble 1	marble 3	marble 4	marble 1	marble 2	marble 4	marble 1	marble 4	bin 1
marble 3	marble 4	marble 1	marble 2	marble 2	marble 3	marble 4	marble 1	marble 2	marble 1	bin 1
marble 3	marble 3	marble 4	marble 2	marble 1	marble 4	marble 2	marble 1	marble 2	marble 1	bin 2
marble 3	marble 1	marble 2	marble 3	marble 4	marble 2	marble 1	marble 4	marble 1	marble 4	bin 1
marble 3	marble 4	marble 4	marble 2	marble 1	marble 3	marble 2	marble 1	marble 2	marble 1	bin 2
marble 3	marble 4	marble 2	marble 2	marble 4	marble 3	marble 1	marble 3	marble 1	marble 1	bin 3
marble 3	marble 1	marble 2	marble 3	marble 1	marble 2	marble 4	marble 2	marble 4	marble 4	bin 4
marble 3	marble 4	marble 1	marble 2	marble 2	marble 1	marble 2	marble 1	marble 4	marble 3	bin 4



## Appendix B

# Colored Trails

### B.1 Scenarios

Below, in Tables B.1-B.3, are the boards, goals and chip distributions of all games used in the colored trails experiment. The tiles of the board are depicted as follows: '0' is a diagonally striped tile, '1' an empty tile, '2' a horizontally striped tile, '3' a dotted tile, and 'S' the start tile. Below each board the accompanying goal locations and chips distributions are given for both players. The first number, in *italic*, is the goal position. The numbers correspond to the tiles shown in Figure B.1. The second number is the number of diagonally striped chips that player started with, the third number the number of empty chips, the fourth number the number of horizontally striped chips, and the fifth number the number of dotted chips.

0	1		2	3
4				5
		<b>S</b>		
6				7
8	9		10	11

Figure B.1: Every number depicts a possible goal tile in the colored trails game.

Table B.1: Scenarios where the computer agent used  $ToM_0$ . See the text for an explanation.

	Starting player agent					Starting player participant				
<b>Row 1</b>	2	1	1	2	1	0	0	2	1	1
<b>Row 2</b>	1	3	2	2	2	1	2	2	0	1
<b>Row 3</b>	3	3	S	1	1	0	1	S	2	0
<b>Row 4</b>	0	1	3	0	2	3	3	0	0	2
<b>Row 5</b>	3	0	3	1	2	3	0	3	2	2
<b>Agent</b>	10	2	1	1	0	3	1	1	2	0
<b>Participant</b>	2	2	0	1	1	3	0	2	1	1
<b>Row 1</b>	2	3	2	0	2	2	0	3	3	2
<b>Row 2</b>	1	3	0	1	1	2	2	2	0	2
<b>Row 3</b>	3	2	S	3	1	0	2	S	0	1
<b>Row 4</b>	3	2	0	3	3	3	3	3	0	1
<b>Row 5</b>	1	0	1	0	3	0	2	0	1	2
<b>Agent</b>	6	1	1	2	0	10	1	2	1	0
<b>Participant</b>	1	0	2	0	2	1	0	0	2	2
<b>Row 1</b>	3	3	1	1	0	0	1	0	1	1
<b>Row 2</b>	1	1	0	3	2	2	0	1	3	3
<b>Row 3</b>	0	3	0	1	2	0	3	3	1	3
<b>Row 4</b>	3	2	1	2	0	3	3	1	1	1
<b>Row 5</b>	1	0	0	3	3	0	1	1	3	3
<b>Agent</b>	5	0	0	1	3	9	1	0	2	1
<b>Participant</b>	6	1	2	0	1	0	1	1	2	0
<b>Row 1</b>	0	3	0	2	2	2	2	2	0	1
<b>Row 2</b>	2	3	3	1	2	2	0	2	0	0
<b>Row 3</b>	3	3	0	0	2	1	3	2	3	2
<b>Row 4</b>	1	2	0	2	3	3	2	0	0	0
<b>Row 5</b>	0	2	3	0	3	0	2	2	3	3
<b>Agent</b>	4	1	1	1	1	3	1	0	1	2
<b>Participant</b>	2	1	0	3	0	10	0	1	3	0

Table B.2: Scenarios where the computer agent used  $ToM_1$ . See the text for an explanation.

	Starting player agent					Starting player participant				
Row 1	3	2	0	1	1	0	2	3	0	2
Row 2	1	3	3	0	3	3	3	3	3	3
Row 3	1	1	0	2	1	3	1	1	1	3
Row 4	0	2	2	0	3	3	0	1	0	2
Row 5	1	0	2	1	2	0	2	2	1	1
Agent	2	2	0	1	1	8	1	1	1	1
Participant	8	1	1	1	1	5	1	0	2	1
Row 1	1	1	0	1	0	3	3	0	0	2
Row 2	0	0	2	2	3	2	3	0	0	2
Row 3	3	3	0	0	0	3	2	1	1	0
Row 4	0	1	1	1	1	1	3	0	1	2
Row 5	0	3	3	1	0	2	0	3	1	2
Agent	6	0	1	1	2	4	1	1	2	0
Participant	3	1	1	1	1	11	2	2	0	0
Row 1	2	2	1	1	2	2	2	1	2	1
Row 2	0	3	2	3	1	2	3	3	3	2
Row 3	0	2	0	0	1	3	0	2	3	1
Row 4	2	0	2	2	3	2	1	0	3	2
Row 5	2	1	3	3	0	2	3	2	3	1
Agent	6	0	1	2	1	0	1	1	1	1
Participant	8	1	1	1	1	3	2	0	1	1
Row 1	2	0	1	0	1	0	1	1	1	0
Row 2	1	3	1	3	0	2	0	2	0	0
Row 3	1	0	3	2	3	3	0	1	0	1
Row 4	3	1	3	1	1	1	3	2	3	3
Row 5	0	1	3	0	0	3	0	0	3	1
Agent	0	1	1	2	0	8	3	1	0	0
Participant	4	0	1	1	2	8	1	1	1	1

Table B.3: Scenarios where the computer agent used  $ToM_2$ . See the text for an explanation.

	Starting player agent					Starting player participant				
<b>Row 1</b>	0	0	0	2	2	1	0	3	1	0
<b>Row 2</b>	2	2	1	3	2	1	1	1	2	1
<b>Row 3</b>	3	2	2	1	2	0	2	1	2	2
<b>Row 4</b>	0	2	0	3	3	2	1	1	3	1
<b>Row 5</b>	3	3	1	1	2	0	0	3	0	2
<b>Agent</b>	6	1	2	1	0	4	2	1	0	1
<b>Participant</b>	1	0	1	2	1	0	1	1	2	0
<b>Row 1</b>	1	3	1	3	3	1	3	3	1	1
<b>Row 2</b>	3	1	3	2	0	1	1	0	3	2
<b>Row 3</b>	3	1	3	1	2	2	0	0	2	3
<b>Row 4</b>	3	0	2	2	1	2	3	3	2	3
<b>Row 5</b>	1	1	2	2	3	0	2	0	2	0
<b>Agent</b>	3	0	2	2	0	11	1	1	0	2
<b>Participant</b>	4	2	0	1	1	1	0	1	0	3
<b>Row 1</b>	2	0	3	0	2	3	0	2	0	2
<b>Row 2</b>	0	1	0	3	2	1	1	1	3	1
<b>Row 3</b>	2	3	0	3	3	1	2	3	0	1
<b>Row 4</b>	1	2	2	0	1	3	1	3	2	1
<b>Row 5</b>	3	2	0	0	1	3	2	2	0	2
<b>Agent</b>	3	2	1	1	0	8	1	1	1	1
<b>Participant</b>	0	0	1	1	2	5	0	0	3	1
<b>Row 1</b>	1	1	1	3	0	2	2	3	0	3
<b>Row 2</b>	2	3	1	2	1	0	1	2	1	3
<b>Row 3</b>	1	3	0	2	1	1	2	0	3	2
<b>Row 4</b>	0	1	3	1	1	1	3	0	0	1
<b>Row 5</b>	2	3	0	0	2	1	1	2	2	2
<b>Agent</b>	4	0	2	0	2	6	2	1	1	0
<b>Participant</b>	1	2	0	2	0	9	1	0	1	2

## B.2 Control questions

At the end of the instructions on colored trails, the participants had to answer several questions, which are shown in Table B.4. The original questions were in Dutch.

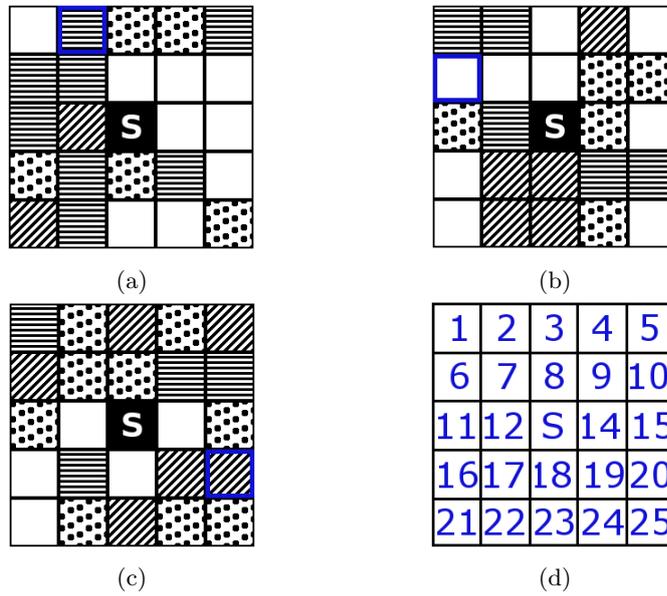


Figure B.2: Three boards used to test the participants knowledge and a board with numbers as references to the tiles.

Table B.4: Control questions on colored trails, translated from Dutch.

Question	Answer
How much time do you have per round (in seconds)?	60 seconds
After how many rounds does the negotiation stop automatically?	6 rounds
<i>Which statements are true?</i>	
If I reach the goal, I get 50 bonus points	True
If I do not reach the goal, I get 10 points deduction	False, 10 points deduction per missing step
If I do not use a chip, I get 5 bonus points for that chip	True
The goal of the opponent is marked	False
<i>Which tile can you reach in the situations in Figure B.2?</i>	
<i>(You only need to indicate the tile closest to the goal, use the numbers in Figure B.2.)</i>	
First board	Tile 7
Second board	Tile 7
Third board	Tile 19

# Appendix C

## Questionnaires: questions

### C.1 Questions about Marble Drop

Every question is followed by the answer options, which are presented in italic. The original questions and answer options were in Dutch.

- Did you find this part to be difficult?
- *Very easy / Easy / Not easy nor hard / Hard / Very hard*
- Did you reason about what the opponent would do?
- *Never / Rarely / Sometimes / Often / Always*
- If so, how?
- *Open answer*
- Do you have general remarks on marble drop?
- *Open answer*

### C.2 Questions about Colored Trails

Every question is followed by the answer options, which are presented in italic. The original questions and answer options were in Dutch.

- Did you find this part to be difficult?
- *Very easy / Easy / Not easy nor hard / Hard / Very hard*
- Did you reason about what Alex would do?
- *Never / Rarely / Sometimes / Often / Always*
- If so, how?
- *Open answer*
- Did you reason about where Alex' his goal was?
- *Never / Rarely / Sometimes / Often / Always*
- If so, how?
- *Open answer*
- Did you have the idea that Alex changed his strategy during the game?
- *Yes / No*
- If so, how often?
- *1 / 2 / 3 / 4 / 5 / 6 / 7 / 8 / 9 / 10*
- If so, how did you notice this?
- *Open answer*

- Do you have the idea that playing marble drop influenced how you played colored trails?
- *Yes / No*
- If so, how?
- *Open answer*
- Do you have general remarks on colored trails?
- *Open answer*

### C.3 Interpersonal Reactivity Index

The interpersonal reactivity index, from [33]. The answer scale was as follows:

**A:** Does not describe me very well

**B:** Does not describe me well

**C:** Neutral

**D:** Describes me well

**E:** Describes me very well

Behind each question a category is mentioned. ‘PT’ is the Perspective Taking scale, ‘FS’ the Fantasy scale, ‘EC’ the Emphatic Concern scale, and ‘PD’ the Personal Distress scale. The minus sign indicates which scoring systems should be use:

**Without minus sign:** A=0, B=1, C=2, D=3, E=4

**With minus sign:** A=4, B=3, C=2, D=1, E=0

The questions:

- I daydream and fantasize, with some regularity, about things that might happen to me. (FS)
- I often have tender, concerned feelings for people less fortunate than me. (EC)
- I sometimes find it difficult to see things from the ‘other guy’s’ point of view. (PT-)
- Sometimes I don’t feel very sorry for other people when they are having problems. (EC-)
- I really get involved with the feelings of the characters in a novel. (FS)
- In emergency situations, I feel apprehensive and ill-at-ease. (PD)
- I am usually objective when I watch a movie or play, and I don’t often get completely caught up in it. (FS-)
- I try to look at everybody’s side of a disagreement before I make a decision. (PT)
- When I see someone being taken advantage of, I feel kind of protective toward them. (EC)
- I sometimes feel helpless when I am in the middle of a very emotional situation. (PD)
- I sometimes try to understand my friends better by imagining how things look from their perspective. (PT)
- Becoming extremely involved in a good book or movie is somewhat rare for me. (FS-)
- When I see someone get hurt, I tend to remain calm. (PD-)
- Other people’s misfortunes do not usually disturb me a great deal. (EC-)

- If I'm sure I'm right about something, I don't waste much time listening to other people's arguments. (PT-)
- After seeing a play or movie, I have felt as though I were one of the characters. (FS)
- Being in a tense emotional situation scares me. (PD)
- When I see someone being treated unfairly, I sometimes don't feel very much pity or them. (EC-)
- I am usually pretty effective in dealing with emergencies. (PD-)
- I am often quite touched by things that I see happen. (EC)
- I believe that there are two sides to every question and try to look at them both. (PT)
- I would describe myself as a pretty soft-hearted person.(EC)
- When I watch a good movie, I can very easily put myself in the place of a leading character. (FS)
- I tend to lose control during emergencies. (PD)
- When I'm upset at someone, I usually try to "put myself in his shoes" for a while. (PT)
- When I am reading an interesting story or novel, I imagine how I would feel if the events in the story were happening to me. (FS)
- When I see someone who badly needs help in an emergency, I go to pieces. (PD)
- Before criticizing somebody, I try to imagine how I would feel if I were in their place. (PT)



## Appendix D

# Questionnaires: participants' answers

### D.1 Colored Trails: Strategy Change by Opponent

Part of the answers on the question whether the opponent ('Alex') changed his strategy during the colored trails experiment. The other answers are omitted due to misunderstandings of the question. The original answers were in Dutch. Every participant who is stated here, reported that Alex changed his strategy twice. At the end of each answer the order in which the three blocks ( $ToM_0/ToM_1/ToM_2$  opponent) were presented is indicated (the participants were not informed about these blocks).

- At the start Alex was quite stubborn and did not give in. During the middle part he seemed to be more willing to make concessions, maybe so he would not get extra chips, but would reach his goal. At the end he was stubborn again, trying not to let me win. ( $ToM_1, ToM_0, ToM_2$ )
- Later on in the game he seemed to accept my offers more easily, settling for something lower himself. I offered him something, he came up with a counter-proposal, but I offered my first proposal again and then he did accept it. Later on he did this less often. ( $ToM_0, ToM_2, ToM_1$ )
- At a certain moment it seemed like Alex was using his offer strategically. For example, proposing something very unreasonable at the end of a game so I had to choose between something really bad and the initial chip distribution. A little while before that it seemed like he made an extra 'step' in his strategy. ( $ToM_1, ToM_0, ToM_2$ )
- At the start, Alex was very compliant. Until the middle part, where it was somewhat more sturdy. At the end it was less sturdy again. ( $ToM_0, ToM_2, ToM_1$ )
- Every now and then Alex 'was hating' and he demanded all the chips. This did not take very long. Furthermore I imagine that he had a temper and sometimes made kind offers (five chips for me and three for himself)

and sometimes his offers were more aggressive (the other way around).  
However, this could also be all in my head. ( $ToM_0$ ,  $ToM_1$ ,  $ToM_2$ )

