[1] Transfer of Learned Opponent Models in Zero Sum Games

[2] Ismail Guennouni[1] & Maarten Speekenbrink[1]

[3] [1] Department of Experimental Psychology, University College London

[4] Author Note

[5] Correspondence concerning this article should be addressed to Ismail Guennouni,

[6] Department of Experimental Psychology, University College London, 26 Bedford Way,

[7] London WC1H 0AP, United Kingdom. E-mail: i.guennouni.17@ucl.ac.uk

8                                    Abstract

9    Enter abstract here. Each new line herein must be indented, like this line.

10         Human learning transfer abilities take advantage of important cognitive building blocks

11   such as an abstract representation of concepts underlying tasks and causal models of the

12   environment. One way to build abstract representations of the environment when the task

13   involves interactions with others is to build a model of the opponent that may inform what

14   actions they are likely to take next. In this study, we propose to explore opponent modelling

15   and its transfer with the use of computer agents possessing human-like theory of mind

16   abilities with limited degrees of iterated reasoning. Through two differentiated experiments,

17   we find that participants can deviate from Nash equilibrium play and learn to adapt to the

18   opponent strategy and exploit it. Moreover, we showed that participants transfer their

19   learning to new games and that the transfer is moderated by the level of sophistication of

20   the opponent. Computational modelling shows that it is likely that players start each game

21   using a model-based learning strategy that facilitates generalisation and opponent model

22   transfer, but then switch to behaviour that is consistent with a model-free learning strategy

23   in the latter stages of the interaction.

<sup>24</sup> Transfer of Learned Opponent Models in Zero Sum Games

<sup>25</sup> **Introduction**

<sup>26</sup>     Being able to transfer previously acquired knowledge to a new domain is one of the

<sup>27</sup> hallmarks of human intelligence. Humans are naturally endowed with the ability to extract

<sup>28</sup> relevant features from a situation, identify the presence of these features in a novel setting

<sup>29</sup> and use previously acquired knowledge to adapt to previously unseen challenges using

<sup>30</sup> acquired knowledge. More formally, @perkins1992transfer defines transfer of learning as the

<sup>31</sup> application of skills, knowledge, and/or attitudes that were learned in one situation to

<sup>32</sup> another learning situation. This typically human skill has so far eluded modern AI agents.

<sup>33</sup> Deep neural networks for instance can do very well on image recognition tasks and can even

<sup>34</sup> reach super-human performance levels on video and strategic board games. Yet they struggle

<sup>35</sup> to learn as fast or as efficiently as humans do, and more importantly they have a very

<sup>36</sup> limited ability to generalize and transfer knowledge to new domains. @Lake2017 argue that

<sup>37</sup> human learning transfer abilities take advantage of important cognitive building blocks such

<sup>38</sup> as an abstract representation of concepts underlying tasks and compositionally structured

<sup>39</sup> causal models of the environment.

<sup>40</sup>     One way to build abstract representations of the environment when the task involves

<sup>41</sup> interactions with others is to build a model of the person we are interacting with that may

<sup>42</sup> inform what actions they are likely to take next. Once we learn something about them, we

<sup>43</sup> can use this knowledge to inform how to best behave in novel situations. This may lead to

<sup>44</sup> very efficient generalization of knowledge, even to situations that are dissimilar to the history

<sup>45</sup> of interaction, assuming what we have learned about others is an abstract representation

<sup>46</sup> that is not too dependent on the environment of the initial interaction. There is evidence

<sup>47</sup> that people learn models of their opponents when they play repeated economic games

<sup>48</sup> [@stahl1995players], engage in bilateral negotiations [@baarslag2016learning], or simply try

<sup>49</sup> to exploit a non random player in chance games such as Rock-Paper-Scissors [@de2012higher].

In this paper, we are specifically interested in the way in which people build and use models of their opponent to facilitate learning transfer, when engaged in situations involving an interaction with strategic considerations. These situations arise frequently such as in negotiations, auctions, strategic planning and all other domains in which theory of mind [@premack1978does] abilities play a role in determining human behaviour. In order to explore learning transfer in these strategic settings, it is generally useful to study simple games as a model of more complex interactions. More specifically, we need a framework that allows the study of whether and how a player takes into consideration, over time, the impact of its current and future actions on the future actions of the opponent and the future cumulative rewards. Repeated games, in which players interact repeatedly with the same opponent and have the ability to learn about the opponent's strategies and preferences [@mertens1990repeated] are particularly adapted to the task of opponent modelling.

Early literature on learning transfer in games has mostly focused on measuring the proportion of people who play normatively optimal (Nash Equilibria) or salient actions (e.g Risk Dominance) in later games, having had experience with a similar game environment previously. For instance, @ho1998iterated measure transfer as the proportion of players who choose the Nash Equilibrium in later p-beauty contest games, after training on similar games. They find there is no evidence of immediate transfer (Nash equilibrium play in the first round of the new game) but positive structural learning transfer as shown by the faster convergence to equilibrium play by experienced vs non experienced players. @knez2000 test learning transfer in players exposed to two games with multiple equilibria sequentially and explore the ability of players to coordinate their actions to choose a particular equilibrium in subsequent games having reached it in prior ones. They distinguish between games that are similar in a purely descriptive way, meaning similar choice labels, identity of players, format and number of action choices; and games that are similar in a strategic sense, meaning similar payoffs from combination of actions, identical equilibrium properties or significant social characteristics of payoffs such as possibility of punishment, need for fairness and

cooperative vs competitive settings. They find that transfer of learning (successful coordination) occurs more readily in the presence of both descriptive and strategic similarity. If the games were only strategically similar, then the transfer was much weaker.

@Juvina2013 made a similar distinction between what they deemed surface and deep similarities and find that both contribute to positive learning transfer. However, they show that surface similarity is not necessary for deep transfer and can either aid or block this type of transfer depending on whether it leads to congruent or incongruent actions in later games. In a series of experiments using economic signalling games Cooper & Kagel [-@cooper2003lessons; -@Cooper2008] found that participants who have learned to play according to a Nash Equilibrium in one game can transfer this to subsequent games, even though the actions consistent with Nash Equilibrium in later games are different. They show that this transfer is driven by the emergence of sophisticated players who are able to represent the strategic implications of their actions and reason about the consequences of a different payoff structure on an opponent's actions.

Most studies fail to offer a formal explanation of transfer or a modelling framework that can explain the experimental observation of transfer between games and generalise it to extensive classes of games. A notable exception is the effort by @Haruvy2012 to specify a model of learning where players learn abstract rules that they can generalise and transfer across dissimilar games, rather than action choices that can only be used within the same game. Participants played ten games, presented in 4x4 normal form (matrix payoffs). Their results suggest that subjects do transfer learning over descriptively similar but strategically dissimilar games and that this learning transfer is significant. They also showed that players learn abstract aspects of the game that are then transferred to new settings. Their rule-learning model, based on @stahl1995players, was able to capture participants' dynamic behavior and shows that the propensity to select particular rules is perfectly transferred across games.

103      In exploring opponent modelling and learning transfer, most studies adopted two types

104 of opponents in the experimental setting. Either human participants were matched with

105 other participants or they played against a computer algorithm. Computer opponents were

106 generally programmed not to change their strategies over the course of the task, allowing

107 better experimental control. One of the commonalities in studies of how humans adapt to

108 computerised opponents is that they have mostly looked at the ability of players to detect

109 and exploit action-based learning rules. The strategies implemented by the computer

110 opponents had a style of play that was not "human-like" in the sense that humans are not

111 very good at playing specific mixed strategies with any precision of at detecting patterns

112 from long sequences of past play due to cognitive constraints. It is therefore important to

113 have agents that "play like humans", and one way of achieving that is to embed theses

114 agents with human-like theory of mind abilities. @simon1972theories explains that humans

115 have limited cognitive capacities and as such cannot be expected to solve computationally

116 complex problems such as finding Nash equilibria. Instead, they will try to "satisfice" by

117 choosing a strategy that is adequate in a simplified model of the environment, rather than an

118 optimal one. This concept finds its natural application in "level-k" theory, first adopted by

119 @stahl1995players. It posits that deviations from Nash equilibrium solutions in simple games

120 are explained by the fact that humans have a heterogeneous degree of strategic

121 sophistication. At the bottom of the ladder, level-0 players are non-strategic and play either

122 randomly or use a salient strategy in the game environment @arad201211. Level-1 players

123 are next up the ladder of strategic sophistication and will assume all their opponents belong

124 to the level-0 category and as such will best respond to them given this assumption.

125 Likewise, a level-2 player will choose actions that are the best response given the belief that

126 all opponents are exactly one level below, and so on.

127      In this study, we propose to explore opponent modelling and its transfer with the use

128 of computer agents possessing human-like theory of mind abilities with limited degrees of

129 iterated reasoning. The agents will either be a level-1 or level-2 player, mimicking human

theory of mind abilities and the limited recursion depth they exhibit [@goodie2012levels]. Our choice of using computer opponents instead of matching groups of participants makes it easier to disentangle the process of learning about the opponent from that of learning about the game structure and payoffs. When playing against other human opponents, players are learning about the game as well as trying to infer the potential dynamic strategy of the opponent simultaneously. Thus, it is harder to focus on an individual and how her strategies are changing and adapting to the opponent's play if we cannot experimentally control the behaviour of the opponent. The use of computer opponents to elicit learning behavior has been explored in the literature with encouraging results. For instance, @spiliopoulos2013strategic made humans play constant sum games against 3 computer opponents, designed to take advantage of known patterns in human play such as imperfect randomization and heuristics use. He found that human participants do adapt to the opponent they are facing. @shachat2004we made human participants face computer opponents playing various mixed strategies in a zero-sum asymmetric matching pennies game. They found that the players changed their strategies towards exploiting the deviations of the opponent from the Mixed Strategy Nash Equilibrium (MSNE), and that this exploitation was very likely if the deviation from the MSNE play was high. We measure transfer of learning about the opponent strategy between games with varying degrees of similarity. The first two games we use are structurally identical except for action labels. In one experiment, the third game is strategically similar to the first two but descriptively different, while in a second experiment, we introduce a third game that is dissimilar to the first two in terms of payoff matrix and strategic structure. In the first experiment, participants face the same opponent throughout the three games, and the opponents are randomised to be either level-1 or level-2 players. In the second experiment, participants faced both level-1 and level-2 opponent sequentially, with the order in which they are faced randomised across participants.

<sub>155</sub> **Experiment 1**

<sub>156</sub> **Methods**

<sub>157</sub> **Participants and Design.**   A total of 52 (28 female, 24 male) participants were

<sub>158</sub> recruited on the Prolific Academic platform. The mean age of participants was 31.2 years.

<sub>159</sub> Participants were paid a fixed fee of £2.5 plus a bonus dependent on their performance

<sub>160</sub> which averaged £1.06. The study used a 2 (computer opponent: level 1 or level 2) by 3

<sub>161</sub> (games: rock-paper-scissors, fire-water-grass, numbers) design, with repeated measures on

<sub>162</sub> the second factor.

<sub>163</sub> **Tasks.**   In the first experiment, the three games were rock-paper-scissors,

<sub>164</sub> fire-water-grass and the numbers game. A typical rock-paper-scissors game (hereafter RPS)

<sub>165</sub> is a 3x3 zero sum game, with a cyclical hierarchy between possible actions: rock blunts

<sub>166</sub> scissors, paper wraps rock, and scissors cut paper. If one player chooses an action which

<sub>167</sub> dominates their opponent's action, the player wins (receives a reward of 1) and the other

<sub>168</sub> player loses (receives a reward of -1). Otherwise it is a draw and both players receive a

<sub>169</sub> reward of 0. It has a unique MSNE consisting of randomly playing one of the three options

<sub>170</sub> each time.The second game is identical to Rock-Paper-Scissors in all but action labels. We

<sub>171</sub> call it Fire-Water-Grass (FWG): Fire burns grass, water extinguishes fire, and grass absorbs

<sub>172</sub> water. We are interested in exploring whether learning is transferred in a fundamentally

<sub>173</sub> similar game where the only difference is in the description of the choice actions. Finally, the

<sub>174</sub> numbers game is a generalization of rock-paper-scissors. In the variant we use, 2 participants

<sub>175</sub> concurrently pick a number between 1 and 5. To win in this game, a participant needs to

<sub>176</sub> pick a number exactly 1 higher than the number chosen by the opponent. For example, if a

<sub>177</sub> participant thinks their opponent will pick 3, they ought to choose 4 to win the round. To

<sub>178</sub> make the strategies cyclical as in RPS, the game stipulates that the lowest number (1) beats

<sub>179</sub> the highest number (5), so if the participant thinks the opponent will play 5, then the

<sub>180</sub> winning choice is to pick 1. This game has a structure similar to RPS in which every action

<sub>181</sub> is dominated by exactly one choice. All other possible combination of choices that are not

consecutive are considered ties. A win would add 1 point to the score of the player, while a loss deduces one point and a tie does not affect the score. Similar to RPS, the MSNE is to play each action with equal probability in a random way.

**Procedure.**   Participants played 3 games sequentially against the same computer opponent. The computer opponent either used a level-1 or level-2 strategy. Participants were informed they would play three different games against the same computer opponent. Each participant plays all three games consecutively and in the same order described above. Participants were told that the opponent cannot cheat and will choose its actions simultaneously without knowledge of the participant's choice. A total of 50 rounds of each game was played with the player's score displayed at the end of each game. The score was calculated as the number of wins minus the number of losses. Ties did not affect the score. In order to incentivise the participants to maximise the number of wins against the opponents, players were paid a bonus at the end of the experiment that was proportional to their final score. Each point worth £0.02. An example of the interface for the rock-paper-scissors game is provided in Figure 1.

**Results.**   Looking at the aggregate scores (See Figure 2 ), the RPS game had the lowest average score across participants (M = 0.289, SD = 0.348) followed by NUMBERS (M = 0.31, SD = 0.347) and finally the FWG game had the highest average score (M = 0.454, SD = 0.354). Aggregate average scores for each game were significantly different from 0 (hypothesised value of random play) using parametric one sample t-tests (RPS: $t(51) = 7.26$, $p < 0.001$ ; FWG: $t(51) = 10.04$ , $p < 0.001$ ; NUMBERS: $t(51) = 7.17$, $p < 0.001$). To analyse within and between game learning, we used a 2 (condition: level-1, level-2) by 3 (game: RPS, FWG, NUMBERS) by 2 (block: first half, second half) repeated measures ANOVA with the first factor varying between participants. There was a main effect of Game ( $F(2,100) = 8.54$, $\eta^2 = 0.05$, $p < 0.001$), showing that average scores varied significantly over the games. Post-hoc pairwise comparisons showed that performance in the FWG game was significantly higher than in the RPS game ( $t(100) = 3.78$, $p = 0.0008$ ), and the
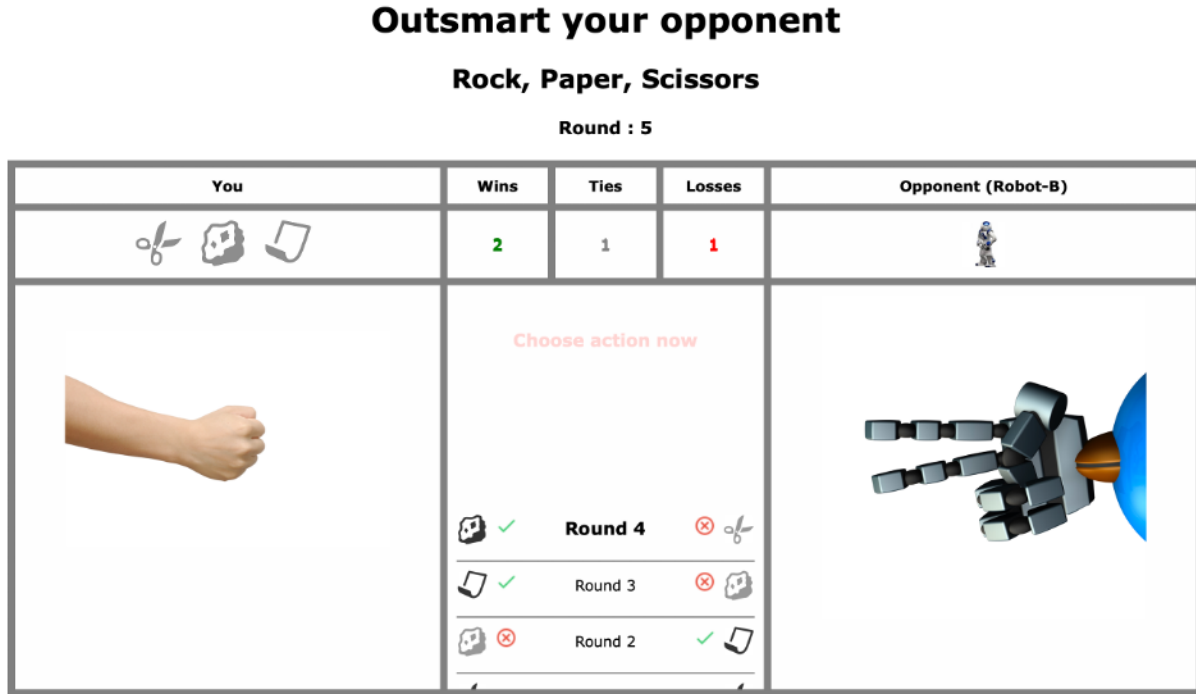
*Figure 1*. Screenshot of the feedback at the end of a round of Rock-Paper-Scissors

209  performance in the NUMBERS game was significantly lower than FWG game ( $t(100) =$

210  -3.32 , p = 0.0024). The score in RPS was not significantly different from the score in

211  NUMBERS ( $t(100) = 0.45$ , p = 0.65). The main effect of Block ( $F(1,50) = 22.51$ , p $<$

212  .001, $\eta^2 = 0.03$) shows that the average score in the first half of games (M = 0.29) was

213  significantly lower than in the second half of the games played (M = 0.40), which translates

214  to within-game learning. The main effect of Condition (F(1,50) = 5.44, p = .024, $\eta^2 = 0.05$)

215  indicates that scores were higher against the level-1 player (M = 0.43) than against the

216  level-2 player (M = 0.27). This indicates that it was harder for participants, on average, to

217  exploit the strategy of the more sophisticated opponent (level-2) compared to that of the

218  comparatively less sophisticated agent (level-1).

219       Finally, the analysis showed a significant block by game interaction ( $F(2,100) = 6.92$ ,

220  p = .002, $\eta^2 = 0.02$), indicating that within-game learning differed between the games.

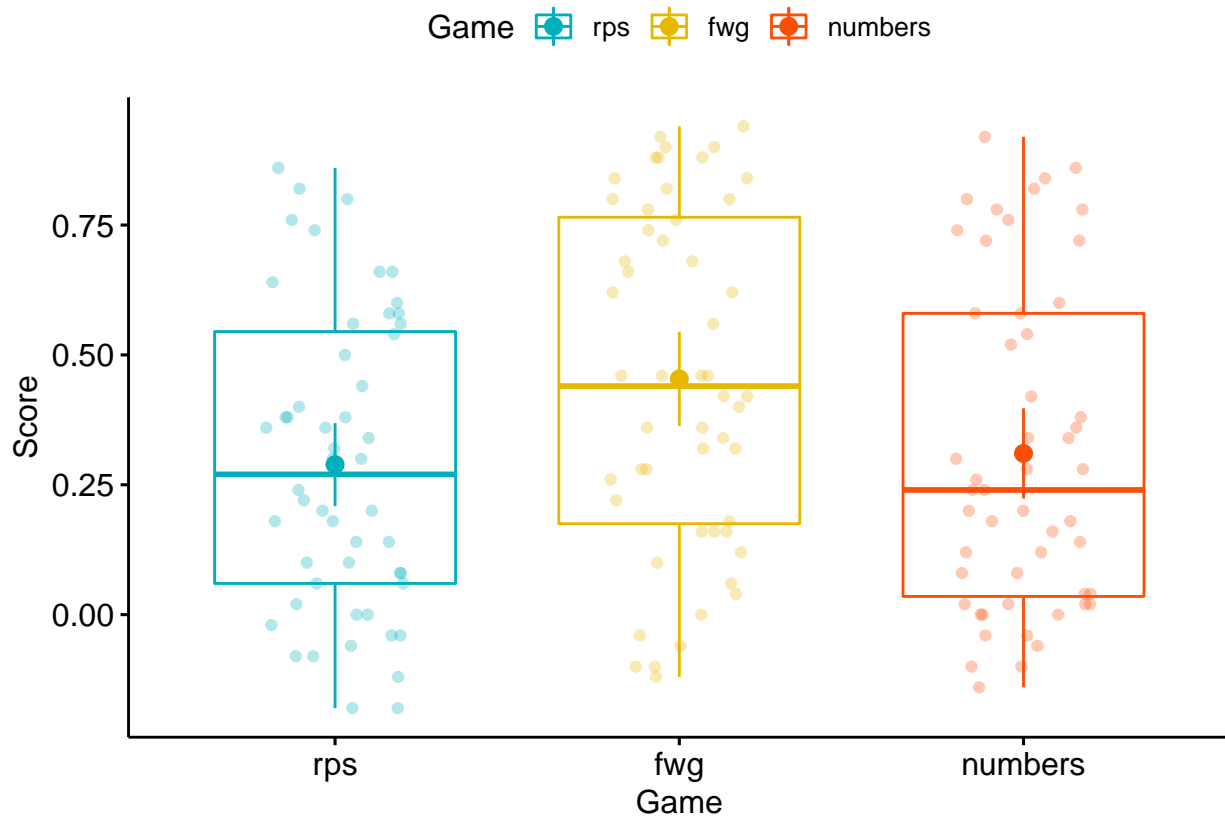221  Indeed, second half scores in RPS are significantly higher than first-half scores (t(150) =

*Figure 2*. Boxplot of scores per game across conditions

5.59, p <.0001), while there was no significant difference between block scores for the other two games. This is indicative of the significant within-game learning happening in the first game when players have no experience against the opponent, as opposed to much lower within game performance improvement in the latter games when participants have had some experience playing against the opponent and start with higher scores indicative of transfer. There was also a three-way interaction between condition, game, and block ( $F(2,100) = 3.88$ , p = .023, $\eta^2 = 0.01$), which indicates, as seen in Figure 3 that within-game learning changes across games also depend on the sophistication of the opponent. For instance, there is more within game learning in the third game against level-2 opponents, since the initial scores are lower than against level-1 opponent. The explanation for this will become clearer when we discuss the factors moderating learning transfer in the next section.
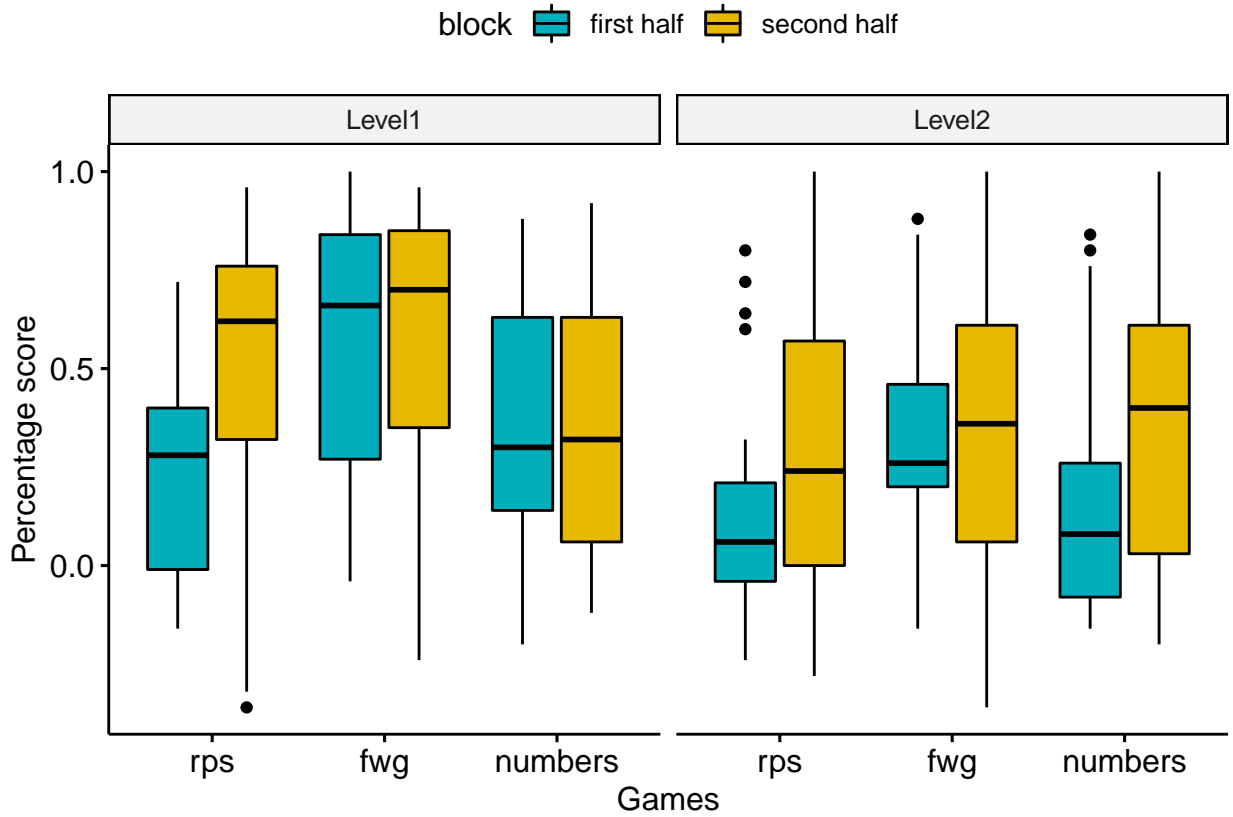
*Figure 3*. Average scores by game, block and condition

233    **Learning transfer.**   As a measure for learning transfer, we focus on participants'

234   scores in the first 5 rounds excluding the initial round (rounds 2-6). We exclude the very first

235   round as the computer opponent plays randomly here and there is no opportunity yet for the

236   human player to exploit their opponent's strategy. A group of players with no experience of

237   the game are expected to perform at chance level over the early rounds of a new game, as

238   was the case in RPS. Positive scores in the early rounds would therefore reflect generalization

239   of prior experience. For the FWG game, the score is significantly higher than 0 ( t(148.85) =

240   4.58 , p < 0.0001). This is also the case for the more dissimilar game : NUMBERS ( t(148.85)

241   = 3.00, p = 0.0092). For the RPS game, the average score is not significantly different from

242   0 as this is the first game and no learning is possible (t(148.85) = 1.04 , p = 0.89).

243    Next, we explore whether learning transfer is moderated by the type of opponent and
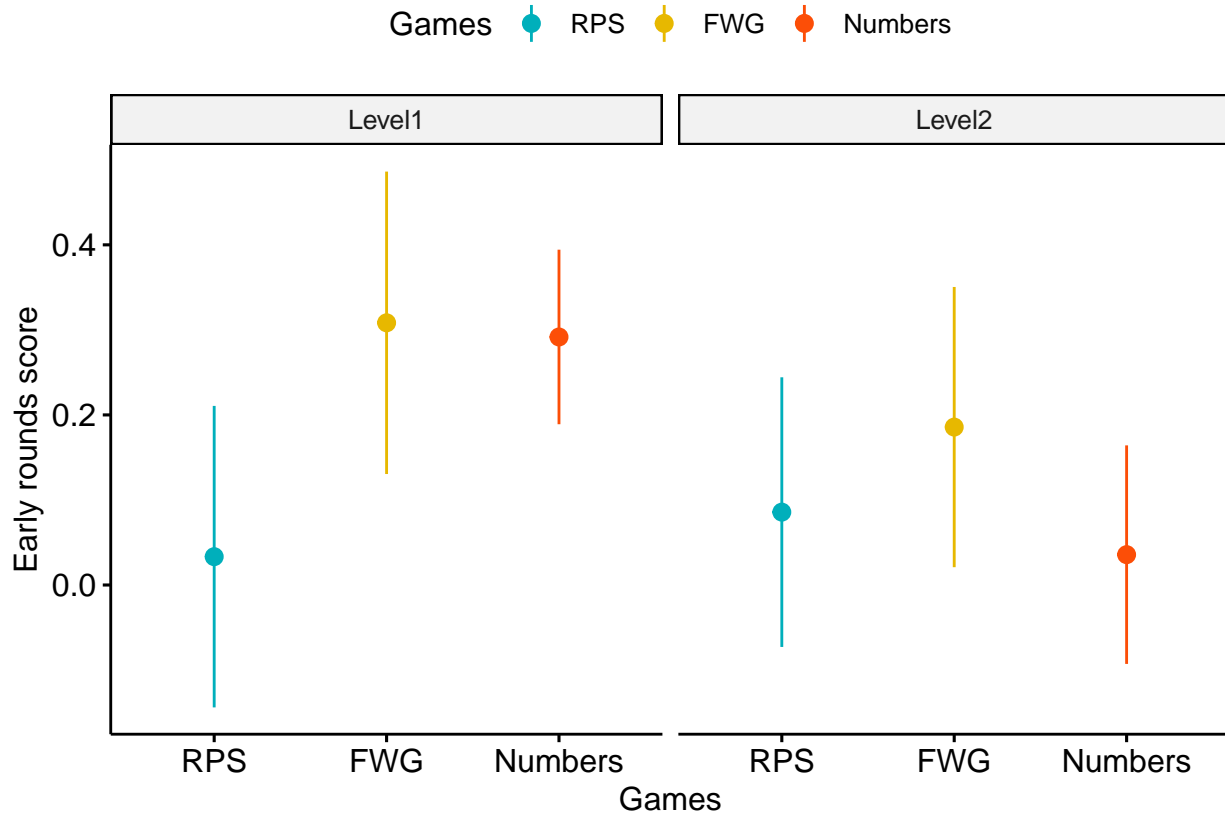
*Figure 4*. Mean and CI of early scores by game and opponent

game similarity. Figure 4 shows the mean scores for rounds 2-6 by game for both level-1 and level-2 facing players as well as the 95 percent confidence interval for the mean. Graphically we can see that the pattern is dissimilar between level-1 and level-2 players, and we suspect transfer to be positively associated with similarity and negatively with degree of sophistication of the agent. To explore this, we run statistical tests on early round scores by game and opponent against the null hypothesis of 0 (no transfer). For level-1 facing players, there is evidence of learning transfer from RPS to both FWG ( $t(150) = 3.96$, $p < 0.001$) and NUMBERS ($t(150) = 3.74$, $p < 0.001$) . For level-2 facing players, there is evidence for transfer from RPS to the similar game FWG, albeit scores are lower than for level-1 player ( $t(150) = 2.48$, $p = 0.01$) but not to the dissimilar game of NUMBERS.

<sub>254</sub>　　**Discussion Experiment 1.**　Our results when averaging across conditions (previous

<sub>255</sub>　section) showed that there was indeed evidence for transfer to the more dissimilar game

<sub>256</sub>　(NUMBERS). We can see from splitting the participants by opponent faced that this transfer

<sub>257</sub>　is exclusively driven by level-1 facing players, as average early round scores of level-2 facing

<sub>258</sub>　players are close to nil in the NUMBERS game. Therefore, both participants facing level-1

<sub>259</sub>　and level-2 agents can transfer learning to the similar game, but only those facing the less

<sub>260</sub>　sophisticated opponent are able to generalise to the less similar game.

<sub>261</sub>                                    **Second Experiment**

<sub>262</sub>            We ran a second experiment with various differentiated features to improve the

<sub>263</sub>    opportunity to measure learning transfer. Instead of making participants face either the

<sub>264</sub>    level-1 or level-2 player throughout, we made them face both opponent sequentially. Because

<sub>265</sub>    there were two distinct opponents, requiring potentially holding two opponent models in

<sub>266</sub>    memory, we also made it easier to recall the results of past rounds by providing participants

<sub>267</sub>    with the opportunity to see the history of the game since the beginning of each interaction.

<sub>268</sub>    Figure 1 shows an example of showing interaction history in the RPS game. Finally, we

<sub>269</sub>    changed the third game to a penalty shootout game, which has the same number of actions

<sub>270</sub>    as the first two. If we see evidence for differential play against opponents, it would show

<sub>271</sub>    participants adapting their strategies to the opponent they are facing, which is indicative of

<sub>272</sub>    opponent modelling.


<sub>273</sub>    **Methods**

<sub>274</sub>            **Participants & Design.**    A total of 48 participants (21 females, 28 males, 1

<sub>275</sub>    unknown) used the Prolific Academic platform to participate in the experiment. This was a

<sub>276</sub>    new set of participants unrelated to those taking part in Experiment 1. The average age of

<sub>277</sub>    players was 30.2 years, and the mean duration to complete the task was 39 minutes.

<sub>278</sub>    Participants were incentivised using a two-tier payment mechanism: a fixed fee of £2.5 for

<sub>279</sub>    completing the experiment plus a performance linked bonus that averaged £1.32.

<sub>280</sub>            **Tasks.**    The three games the participants played were Rock-Paper-Scissors,

<sub>281</sub>    Fire-Water-Grass, and the penalty shootout game. The first two games were identical to the

<sub>282</sub>    ones used in the first experiment. In the final game (shootout) the participants took the role

<sub>283</sub>    of the player shooting a football (soccer) penalty, with the AI opponent being the goalkeeper.

<sub>284</sub>    Players had the choice between three actions, like in the first two games: Shooting the

<sub>285</sub>    football to the left, right or centre of the goal. If the player shoots in a direction different

<sub>286</sub>    from that of where the goalkeeper dives, they win the round and the goalkeeper loses.

Otherwise, the goalkeeper catches the ball and the player loses the round. There is no possibility of ties in this game. Figure 5 shows a snapshot of play in the shootout game. What makes this game different however is that there are two ways to beat the opponent in each round: if we think the opponent is going to choose "'right"' in the next round, we can win by both choosing "'left"' and "'center"'. A level-1 player (thinks that his opponent will repeat his last action) has two ways to win the next round. As such, we have programmed the level-2 computer program to choose randomly between the two possibilities that a level-1 player may choose.
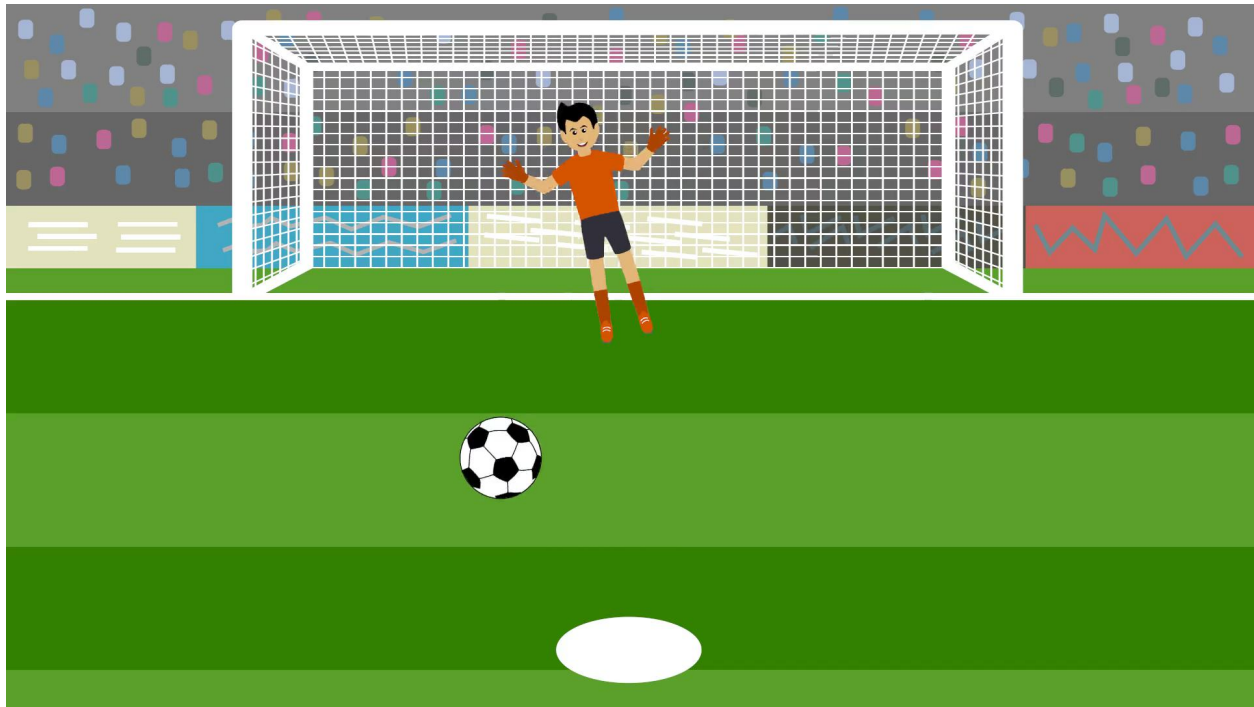


*Figure 5*. Screenshot of the shootout game

**Procedure.**    The participants played 3 games sequentially against both level-1 and level-2 computer opponents, rather than just one like in the first experiment. Like in the first experiment, the computer opponents retained the same strategy throughout the 3 games, however the participants faced each opponent twice in each game. Specifically, each game was divided into 4 stages numbered 1 to 4, consisting of 20, 20, 10, and 10 rounds respectively for a total of 60 rounds per game. Participants start by facing one of the

301 opponents in stage one, then face the other in stage two. This is repeated in the same order

302 in stages 3 and 4. Which opponent they faced first was counterbalanced. All participants

303 engage in the same three games (namely RPS, FWG and Shootout) in this exact order, and

304 were aware that the opponent was not able to know their choices beforehand but was

305 choosing simultaneously with the player. In order to encourage participants to think about

306 their next choice, a countdown timer of 3 seconds was introduced at the beginning of each

307 round. During those 3 seconds, participants could not choose an option and had to wait for

308 the timer to run out. A small delay that changed randomly (between 0.5 and 3 seconds) was

309 also introduced in the time it took the AI agent to give back their response, as a way of

310 simulating a real human opponent thinking time. After each round, the participants were

311 given detailed feedback about their opponent actions as well as whether they won or lost the

312 last round. Further information about the outcome of previous rounds was also visible on

313 the game screen below the feedback area. Throughout each stage, participants could scroll

314 down to recall the history of interaction. The number of wins, losses and ties were clearly

315 shown at the top of the screen for each game, and this scoreboard was reinitialised to zero at

316 the onset of a new stage game. As in the first experiment, all the games have a unique

317 MSNE consisting of randomising across actions. If participants follow this strategy, or simply

318 don't engage in learning how the opponent plays, they would score 0 on average against both

319 level-1 and level-2 players. Evidence of sustained wins would indicate that participants have

320 learned to exploit patterns in the opponent play.

**Results**

322      The RPS game had the lowest average score per round (M = 0.194, SD = 0.345)

323 followed by FWG (M = 0.27, SD = 0.394) and finally the Shootout game had an adjusted

324 average score in between the two (M = 0.289, SD = 0.326).[1] Using parametric t-tests on

---

[1] A higher score in shootout is expected as there are 2 out of three possible winning actions, compared to one

out of three in RPS and FWG. Indeed, a player not aiming to uncover the opponent's strategy and thus

325  adjusted scores, we reject the null hypothesis of random play in all three games (RPS: t(49)

326  = 6.26, $p < 0.0001$ ; FWG: t(49) = 7.25 , $p < 0.0001$ ; Shootout: t(49) = 13.61, $p < 0.0001$ ).

327  Using the average scores obtained by participants in each game and interaction, we explore

328  whether learning has occurred within and between games. We perform a two (condition:

329  level-1 first, level-2 first) by two ( opponent type: level-1 or level-2) by three (game: RPS,

330  FWG, Shootout) by two (interaction: first or second) repeated measures ANOVA with the

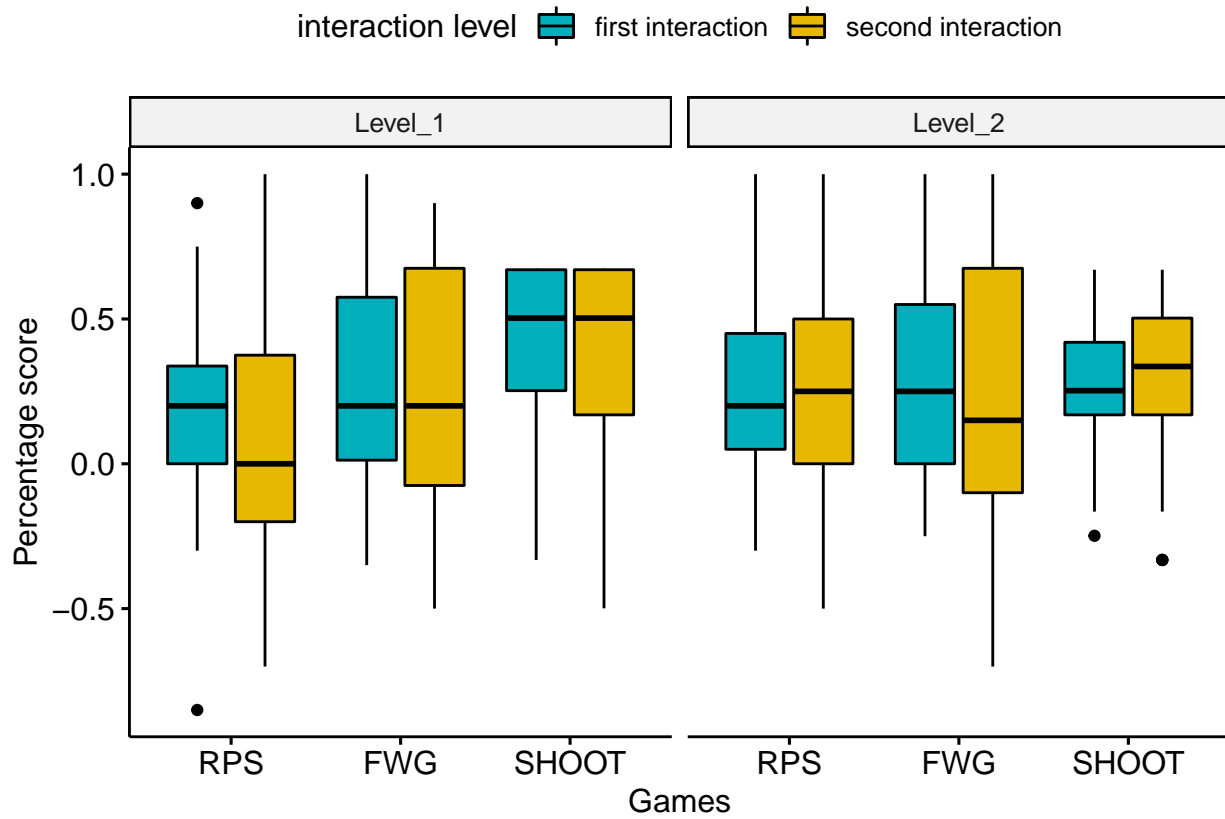331  first factor varying between participants.



*Figure 6*. Boxplot of scores per game and interaction by opponent type

332       There is evidence for a main effect of Game on scores (F(1.85,88.7) = 11.81, $\eta^2 = 0.04$,

choosing to play randomly should be expected to have on average score per round of 0 in both RPS and

FWG, and 0.33 in the Shootout game. To make the scores more comparable, and because we are interested

in player's performance that is not due to chance, we will adjust all scores in the shootout game by

subtracting the average score per round of a random strategy (0.33)

333  p < .0001). To explore these differences further, we look at post-hoc analyses for pairwise

334  comparisons between game scores (p-values adjusted using Holm method for multiple

335  comparisons). We find the performance in the games increases steadily throughout the

336  experiment, with FWG performance significantly higher than RPS (t(96) =2.53, p = 0.025),

337  and performance in the Shootout game also significantly higher than in FWG ( t(96) = 2.32,

338  p = 0.025 ). There was no main effect of either opponent type, the interaction factor( first or

339  second time opponent was faced) , or the condition factor (whether level-1 or level-2

340  opponent was faced first). There was however a significant interaction effect between Game

341  and opponent type ( $F(1.7, 81.82) = 5.31, \eta^2 = 0.02$, p = .01). Figure 6 shows boxplots of

342  game scores, averaged across participants, by game and opponent type. We also distinguish

343  between scores from the first time the players faced the opponent (first interaction) and the

344  second time they did (second interaction). We see that when facing level-1 agents, scores

345  increase steadily after each game, with FWG score significantly higher than RPS ( t(191) =

346  2.70, p = 0.03) and Shootout scores in turn significantly higher than FWG ( t(191) = 3.05, p

347  = 0.01). There was no significant difference between average scores on any two games when

348  facing level-2 agents however.

349  **Learning transfer.**   As a measure for learning transfer we will again compare scores

350  only on rounds 2-6 of each game, excluding the very first round where play is necessarily

351  random.

352  In Figure 7, we show the average score across participants and its 95 percent

353  confidence interval in rounds 2-6 of the first interaction with the opponent for each game.

354  These scores are also averaged across the levels of condition (meaning they are irrespective of

355  which opponent players faced first). For both the FWG and Shootout games, score in the

356  early rounds of the first interaction are significantly higher than 0 for both opponent types.

357  (Level-1 opponent: FWG: t(270) = 4.99, p < 0.0001; Shootout: t(270) = 6.66, p < 0.0001;

358  Level-2 opponent: FWG: t(270) = 4.40, p < 0.0001; Shootout: t(270) = 3.21, p =0.004 ).
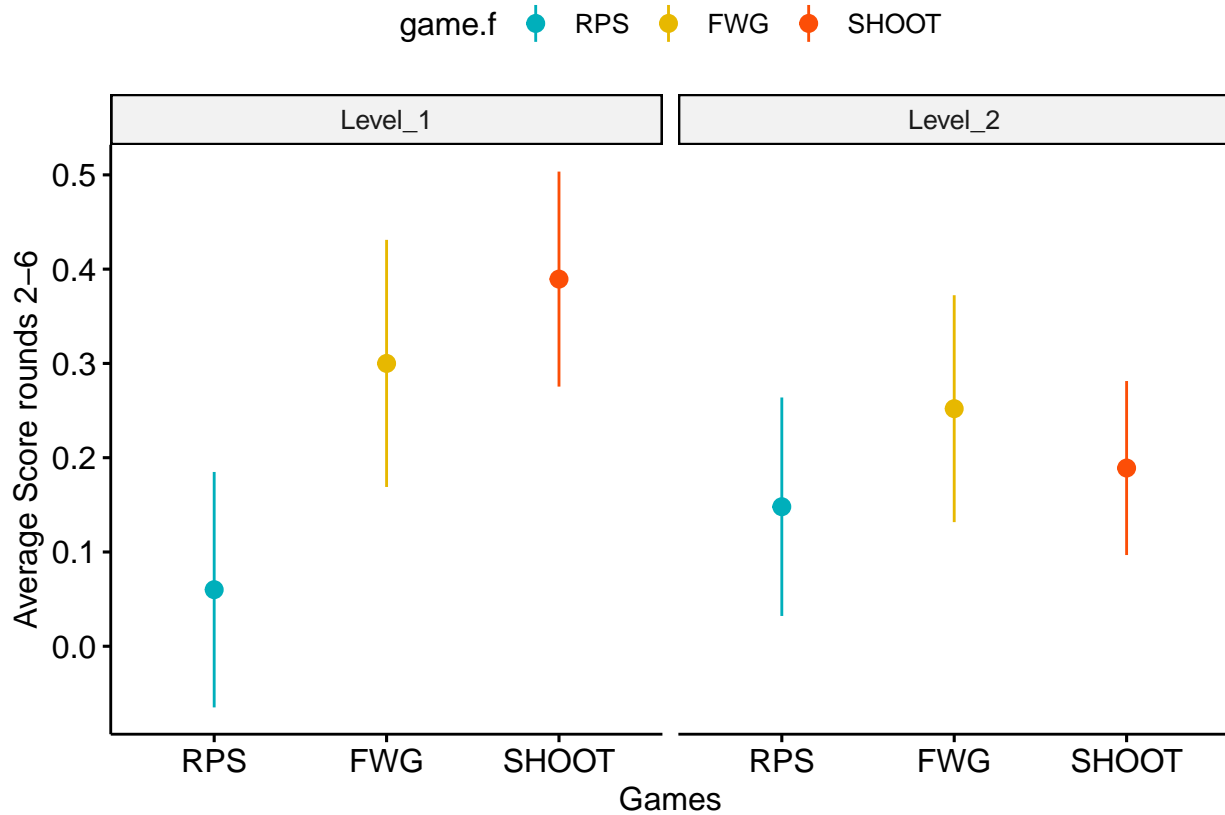
*Figure 7.* Average early round scores and confidence intervals by game and opponent for experiment 2

359   **Experiment 2 discussion:**   Looking at learning transfer by type of opponent faced,

360   we confirm the result from the first experiment that it is easier to transfer learning to the

361   more dissimilar game (Shootout) when facing a level 1 opponent. Indeed, while the early

362   scores of FWG for level-1 and level-2 facing players are not significantly different from each

363   other, the score of the players facing the level-1 opponent is indeed almost 0.2 point per

364   round higher than that of players facing level-2 opponents, and the difference is statistically

365   significant ( $t(144) = 2.45$ , $p = 0.01$). These early scores have also been adjusted to account

366   for the fact that the shootout game has higher average scores when playing randomly, and

367   therefore this difference is really due to better learning transfer and not due to chance.

<sub>368</sub>                                           **Computational modelling**

<sub>369</sub>          To gain more insight into how participants played the games against the computer

<sub>370</sub>   opponents, we estimated and compared different computational models of strategies the

<sub>371</sub>   players may have been using to learn how to beat the opponent. The baseline model assumes

<sub>372</sub>   play is random, and each potential action is chosen with equal probability. Note that this

<sub>373</sub>   corresponds to the Nash equilibrium strategy. In this section, we will go through the various

<sub>374</sub>   models we have used and explain how they update what they learn about the game or the

<sub>375</sub>   opponent

<sub>376</sub>   **Reinforcement Learning**

<sub>377</sub>          We include for comparison purposes a simple model-free reinforcement learning

<sub>378</sub>   algorithm, that reinforces actions that have led to positive rewards, and conversely lowers

<sub>379</sub>   the likelihood of choosing actions that led to a negative reward, irrespective of any state. We

<sub>380</sub>   will use a simple delta learning update rule:

$$V_{t+1}(a) = V_t(a) + \alpha * (R_t - V_t(a))$$

<sub>381</sub>          Where $V_t(a)$ is the value associated with action $a$ at time $t$, $\alpha$ is the learning rate and

<sub>382</sub>   $R_t$ the reward at time t. The probability of player $i$ choosing action $j$ at time $t + 1$ denoted

<sub>383</sub>   by $P_i^j(t + 1)$ is based on action values using a softmax choice rule:

$$P_i^j(t + 1) = \frac{e^{\lambda . V_t^j(a)}}{\sum_{k=1}^{m_i} e^{\lambda . V_t^k(a)}}$$

<sub>384</sub>          We extend this very simple model by adding a state space that consists of last round

<sub>385</sub>   human and agent play. This is akin to using a Q-learning algorithm [@watkins1992q]. The

<sub>386</sub>   update rule becomes:

$$Q^{new}(s_t, a_t) = Q(s_t, a_t) + \alpha * (R_t + \gamma * \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

Where $Q(s_t, a_t)$ is the value of taking action $a$ when in state $s$ at time $t$, $\gamma$ is the discount rate applied to future rewards. For instance, $Q(RS, P)$ denotes the value of taking action "Paper" this round if the player's last action was "Rock" and the opponent played "Scissors". This is a much richer model allowing the players to compute the values of actions conditional on past play.

**EWA Models.** Next, we use a self-tuning Experience Weighted attraction model [@ho2004economics]. EWA models particularity is that they nest two seemingly different approaches, namely reinforcement learning and belief learning. Belief based models are based on the assumption that players keep track of the frequency of past plays and best respond to that. In contrast, reinforcement learning does not take into account beliefs about other players, but is such that an action followed by a positive reward is more likely to be repeated than an action followed by a negative reward. The self-tuning EWA model has been shown to perform better than both these nested models in multiple repeated games and has the advantage of having only one free parameter, the inverse temperature in the softmax choice function.

Let's define some notation in order to write the update rule of the self-tuning EWA model. For player i, there are $m_i$ strategies, denoted $s_i^j$ (i.e player i's strategy number j). Strategies actually played by i in period t, are denoted $s_i(t)$, while the opponent's strategy at time t is denoted $s_{-i}(t)$. After playing $s_i^j$ at time $t$, player i pay-off is denoted $\pi_i(s_i^j, s_{-i}(t))$, and the actual pay-off the player received is $\pi_i(t)$.

The EWA model is based on updating "Attractions" for each action over time. For instance, the attraction of strategy $j$ to player $i$ at time $t$ is written $A_i^j(t)$. Future action choice probabilities are based on these attractions using the softmax playing rule:

$$P_i^j(t+1) = \frac{e^{\lambda.A_i^j(t)}}{\sum_{k=1}^{m_i} e^{\lambda.A_i^k(t)}}$$

410     The attractions are updated over every time step $t$ using the following update rule :

$$A_i^j(t) = \frac{\phi.N(t-1).A_i^j(t-1) + [\delta + (1-\delta).I(s_i^j, s_i(t))].\pi_i(s_i^j, s_{-i}(t))}{\phi.N(t-1) + 1}$$

411     Here, I(x,y) is the indicator function equal to 1 if $x = y$ and 0 otherwise. A simple way

412     to think about this update rule is that attractions are multiplied by a parameter that

413     represents experience $(N(t-1))$ which is itself decaying by a weight $\phi$. The result is added

414     to either the pay-off received (when the indicator function is 0), or to $\delta$ times the foregone

415     pay-off (when indicator function is 1). We can see that setting $\delta = 0$ leads to reinforcement

416     of past actions, while positive and high delta parameters make the update rule take into

417     account foregone pay-offs, which is similar to weighted fictitious play [@cheung1994learning].

418     While the assumption in expanding the update rule above is that $\phi$ and $\delta$ are free

419     parameters [@camerer1997experience], the self-tuning aspect of the model comes from the

420     fact that these are now self-tuned using the formulas expanded in [@ho2004economics].

421     **ToM models.**   In this set of models, we assume that participants have a belief that

422     the opponent is a level-k agent, with uniform probability of the level k, and use evidence of

423     past play to update their beliefs in a Bayesian way about the true value of k. We use values

424     of k in $0, 1, 2$. These models also assume the opponent can deviate from these level-k

425     strategies and play randomly with probability $\theta$, a parameter to be fit. We distinguish

426     between multiple ToM models based on their ability to keep what was learned about the

427     opponent in memory and hence facilitate transfer. In a No-Between-Transfer (NBT) model,

428     players have no memory of what was learned about the opponent and start every new game

429     assuming each level-k has equal probability. In the context of Experiment 2 where players

430     face both opponents, this model assumes that participants transfer learning within the same

431  game, from the first to the second interaction with the opponent, but are not able to transfer

432  that learning to new games (So within but no between transfer)

433      Conversely, In a Between-Transfer model (BT), players are assumed to keep in memory

434  what was learned about the type of opponent faced (vector of probabilities of level-k) and

435  use that at the beginning of each new game. In the context of experiment 2, we still assume

436  that if between transfer is present, then within game transfer is also present (from first to

437  second interaction).

438      We fit the above two models to experiment 1 data. In experiment two, on top of these

439  two models, we fit another model in which all stages of the game and all new games start

440  with a fresh uniform probability of level-k opponent (NT), so no within or between opponent

441  model learning transfer.

442      **Estimation and model comparison.**  In both experiments, all models were fit to

443  each participant data, with optimal parameters being estimated using maximum likelihood.

444  Using information criteria based Bayesian model comparison (BIC), the best fitting model

445  for each participant was chosen and we compared the number of participants whose behavior

446  was best explained by each model.

447  **Experiment 1 modelling :**

448      Figure 8 shows the results for experiment 1: we can see that while some participant's

449  learning behavior was either random or explained by some of the base models, a significant

450  number of participants in experiment one had learning most consistent with Q-learning with

451  states defined by last round play.
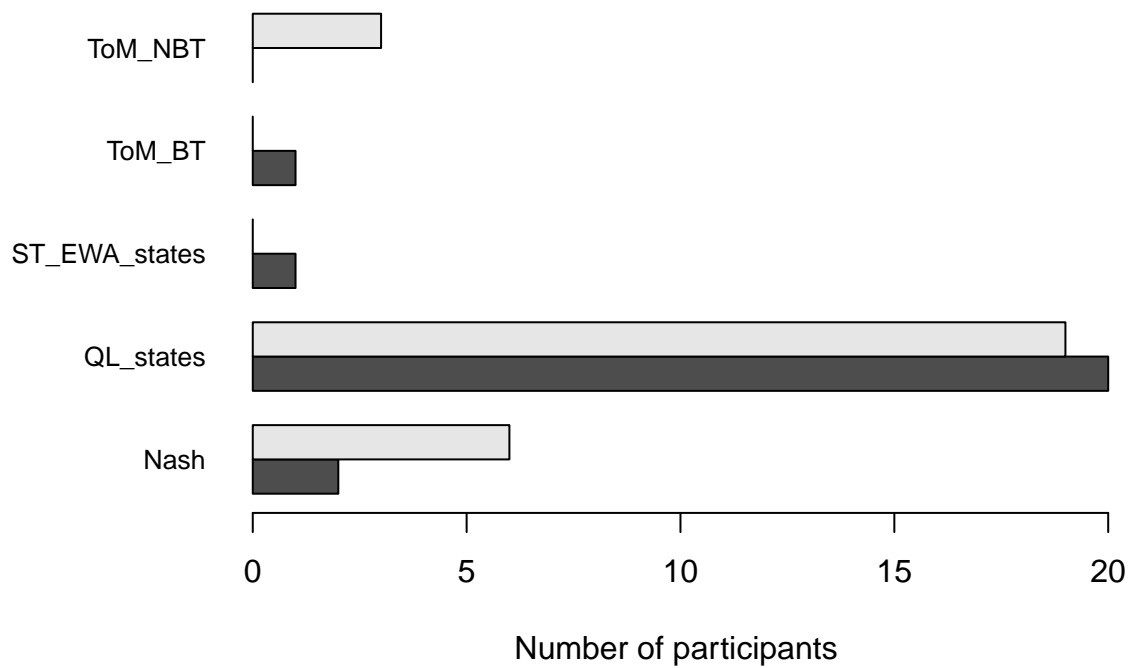
452      Table  1 shows the model BIC weights as well as the number of participants best fit by

453  each model.

454      Next we compared the performance of players whose actions are best fit by each of our

Table 1

*Experiment one Average BIC weights and number of participants best fit by model*

|                   | Nash | ToM_BT | ToM_NBT | QL_states | STEWA_States |
|-------------------|------|--------|---------|-----------|--------------|
| Model BIC weights | 0.10 | 0.07   | 0.05    | 0.73      | 0.05         |
| Count best fit    | 8.00 | 1.00   | 3.00    | 39.00     | 1.00         |



*Figure 8*. Experiment 1 - Histogram of best fitting computational models by condition

hypothesized models. Figure 9 shows the average cumulative performance of players across games, for participants grouped by which model best fits their behavior in experiment 1. We can see that participants whose actions are most consistent with learning a ToM opponent model in a Bayesian way had the best overall performance (without transfer), followed by Q-learning conditional on last round play. EWA, QL and random players had, understandably the lowest performance.
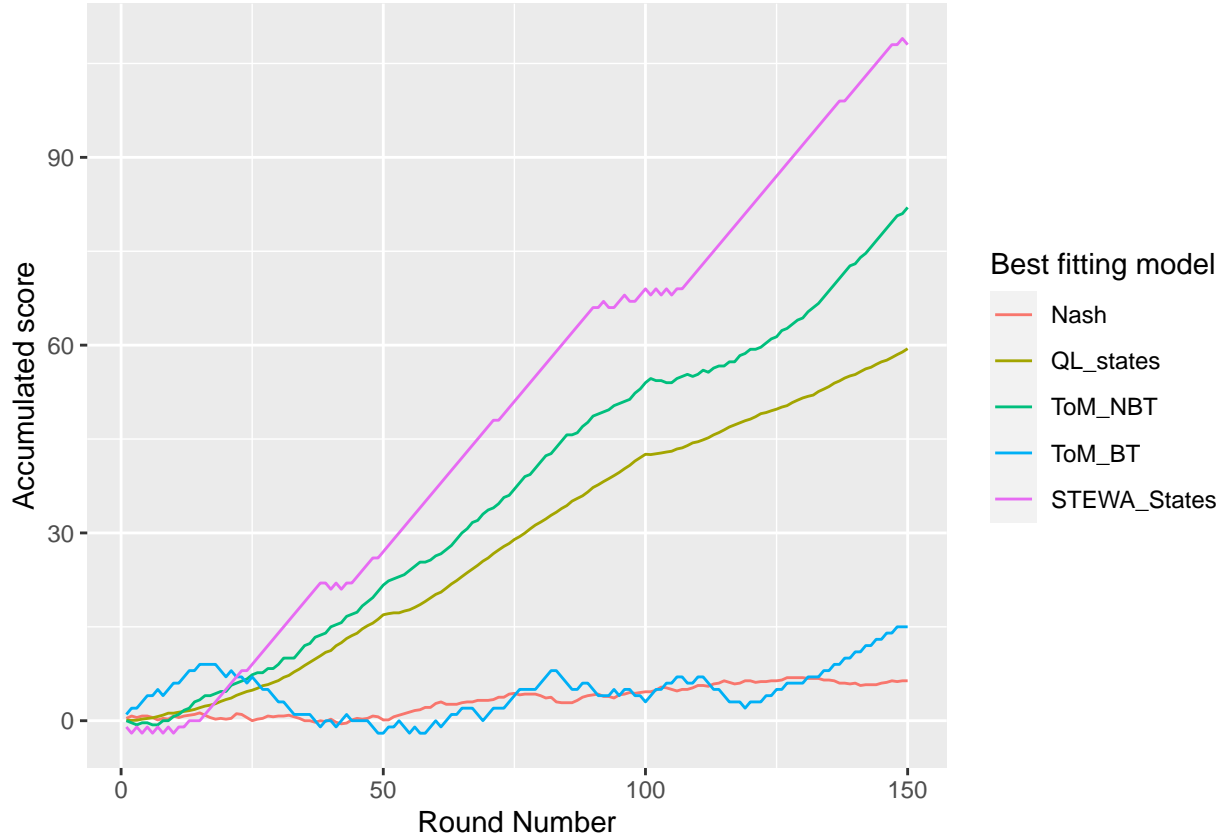


*Figure 9.* Experiment 1 - Average cumulative scores of participants by best fitting model

**Experiment 2**

```
##                  Nash TOM BT TOM NBT TOM NT QLS NT QLS Tr STEWA NT STEWA Tr
## BIC weights      0.13    0.02    0.06    0.06    0.18    0.49      0.04      0.03
## Count best fit 9.00    0.00    2.00    3.00    8.00  27.00      1.00      0.00
```

In experiment 2, we can see from Figure 10 that Q-learning with the aforementioned
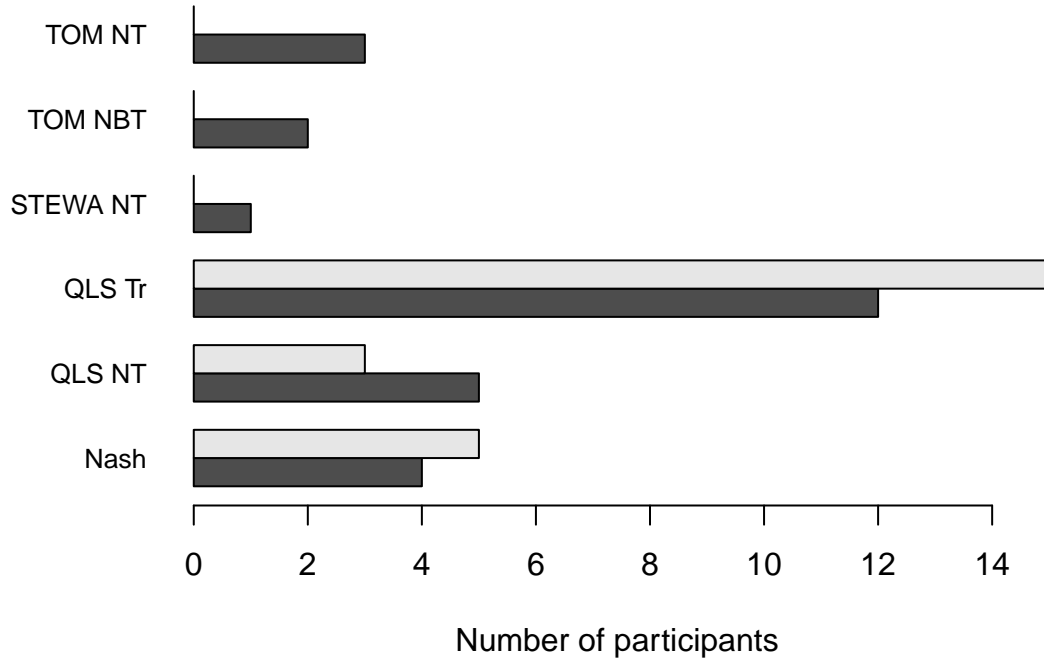
*Figure 10*. Experiment 2 Histogram of best fitting computational models by condition

state-space was again more successful than the Bayesian models in fitting player's action

choices. In experiment 2 when participants faced both level-1 and level-2 agents sequentially,

the Bayesian models (with or without transfer) did not fit players observed data as well.

This is also reflected in BIC weights in Table 2

Plotting cumulative scores by best model for experiment 2, we see very similar results

looking at Figure 11, in that participants whose behavior was best fit by a ToM model of

learning the opponent strategy had the highest cumulative performance. Out of these ToM

models, the one in which there is within-game but no between-game transfer (NBT) had the

best cumulative performance (although it only fit 2 participants best), followed by a model

in which both within and between transfer of opponent models is allowed (BT). The next

best model from a performance perspective was a Q-learning model with states and within

Table 2

*Experiment 2 - Average BIC weights and number of participants best fit by model*

|                | Nash | TOM BT | TOM NBT | TOM NT | QLS NT | QLS Tr | STEWA NT | STEWA Tr |
|----------------|------|--------|---------|--------|--------|--------|----------|----------|
| BIC weights    | 0.13 | 0.02   | 0.06    | 0.06   | 0.18   | 0.49   | 0.04     | 0.03     |
| Count best fit | 9.00 | 0.00   | 2.00    | 3.00   | 8.00   | 27.00  | 1.00     | 0.00     |

game transfer, followed by ToM models where players reset opponent models at each stage of each game (NT). As expected, random play was at the bottom of cumulative performance.

**Using Hidden Markov Model to explore strategy switching**

The computational modelling indicates that most players are best fit by Q-learning type models with states defined by last round play. This is at odds with the findings from the section regarding learning transfer: If indeed most participants use Q-learning with states to choose their actions, then they should not be able to transfer learning to the early rounds of the new game. In order to understand better what is going on, we plot the likelihood by trial for each game and each of the three strategies: Q-learning with states, and Theory of Mind models with and without the possibility of across game transfer.

We start with experiment 1 data. Figure 12 shows that in the later games, the likelihood for the ToM models is higher in the initial rounds in which learning transfer is measured, but that over time, the likelihood of Q-learning model becomes more important and exceeds that of ToM models.

Likewise, in experiment 2, we want to understand the dynamic of strategy choice by plotting the likelihood by trial for each strategy, using the optimal parameters found when fitting the model. Figure 13 shows that, as in experiment 1, ToM models had higher likelihood in the early stages of the second (most similar) game, however the likelihood of Q-learning with states models increases steadily to be the highest in the later stages of all games. In the third and more dissimilar game, we get a result that is different from
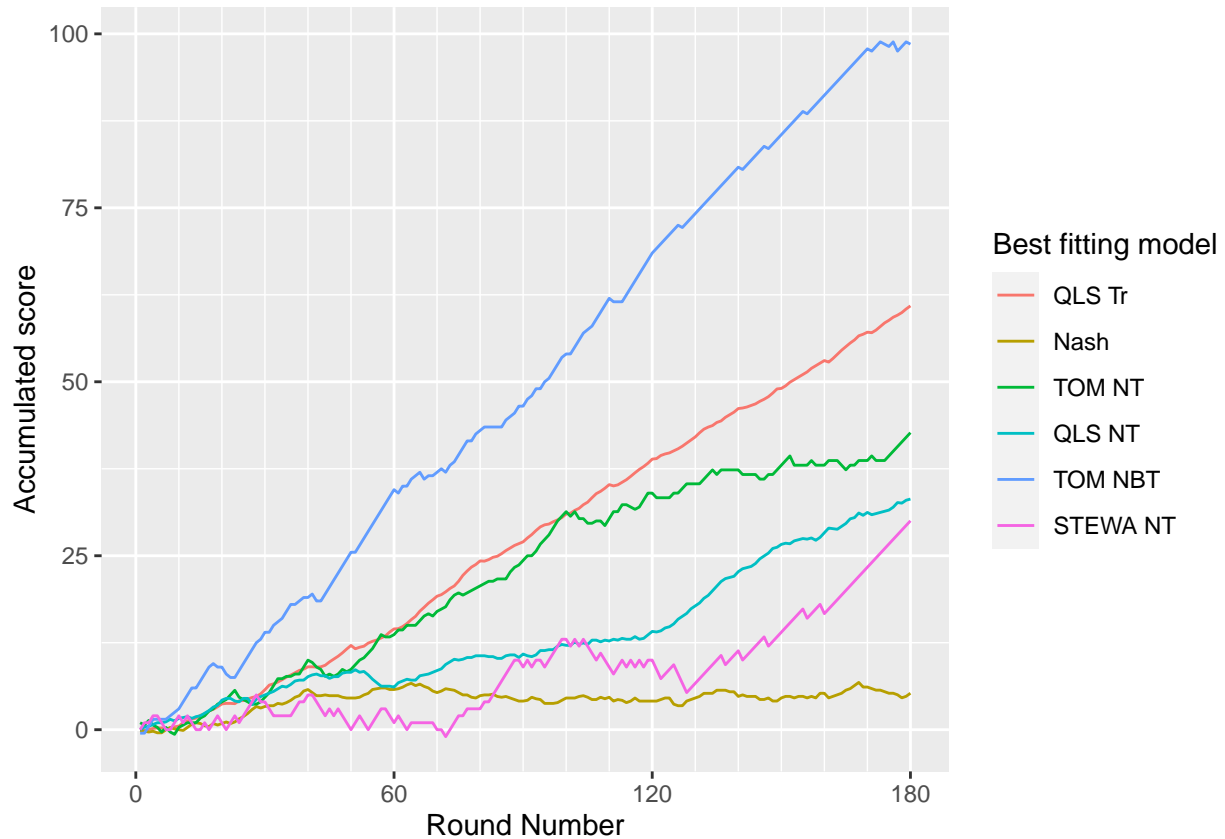
*Figure 11*. Experiment 2 Average cumulative scores of participants by best fitting model

experiment 1. In this instance, the likelihoods of the ToM models stay constant and close to their initial values.

The fact that the likelihoods of the main strategies considered cross over in both experiments could be interpreted as indicative of participants switching between strategies as the games progressed. Indeed, in both experiments, following our results, it seems that in the earlier stages of the latter games, the ToM based strategies fitted observed action choices better than Q-learning based ones, with a reversal of the roles in later stages.

In order to test for the existence of strategy switching in participants' play, we fit Hidden Markov Models in which the latent states are the 3 strategies used (Q-learning with state space consisting of previous round play, ToM based model with opponent model
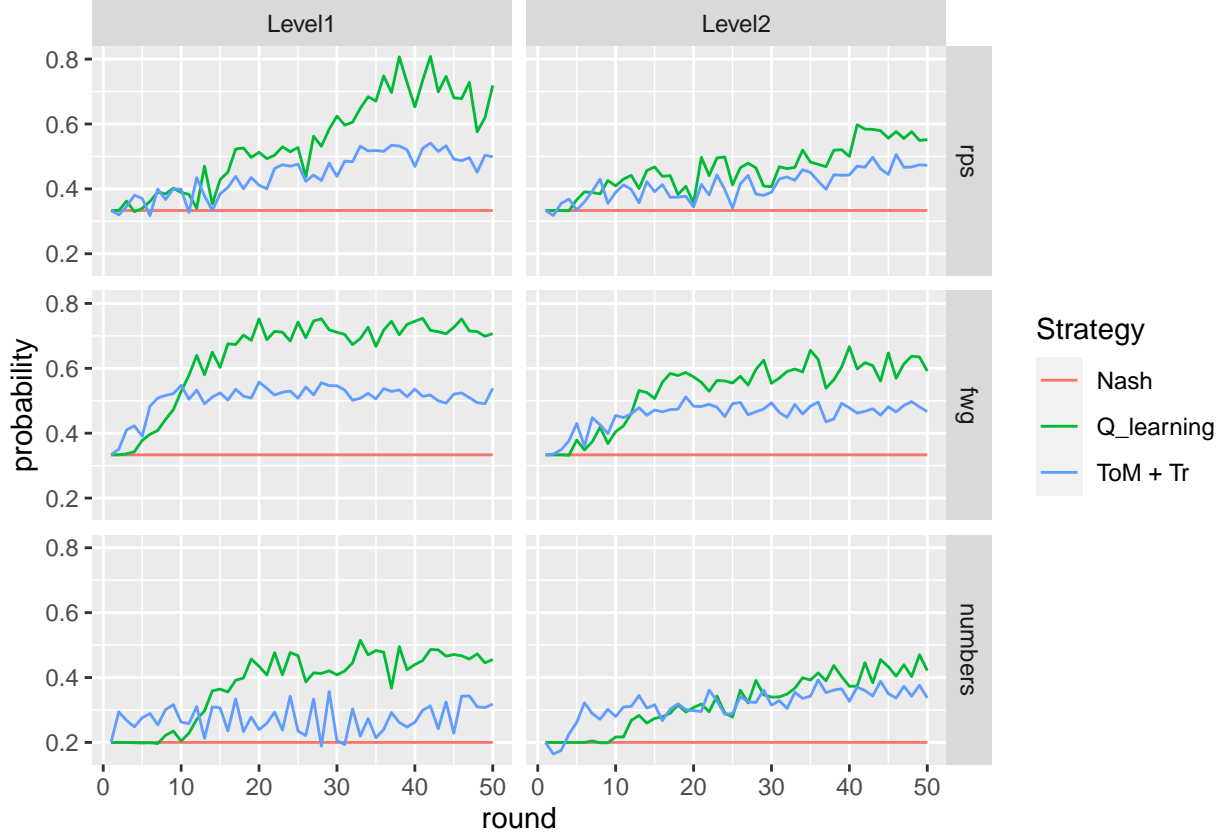
*Figure 12*. Experiment 1 Likelihood by trial by game and opponent faced

507  transfer, and a base model consisting of random play consistent with a Nash equilibrium

508  strategy). Hidden Markov models are useful tools to explore structure in observed time

509  series. They are named as such because of two properties: First, they make the assumption

510  that any observable action at time t results from a process whose state at time t , named $S_t$

511  is "hidden" from the observer. Second, it also assumes that this hidden process has a Markov

512  property, meaning that given state $S_{t-1}$, the value of $S_t$ is independent of all states occurring

513  before time $t - 1$. We also assume that $S_t$ has a discrete probability distribution in that it

514  take one of K discrete values. The model is therefore specified by initial probabilities of

515  being in each state $1, 2, ..., K$ and transition probabilities for moving from state $i$ to state $j$.

516  These probabilities are fit using observed actions generated from these hidden states.

517       To investigate the possibility of strategy switching, we fit two different hidden Markov
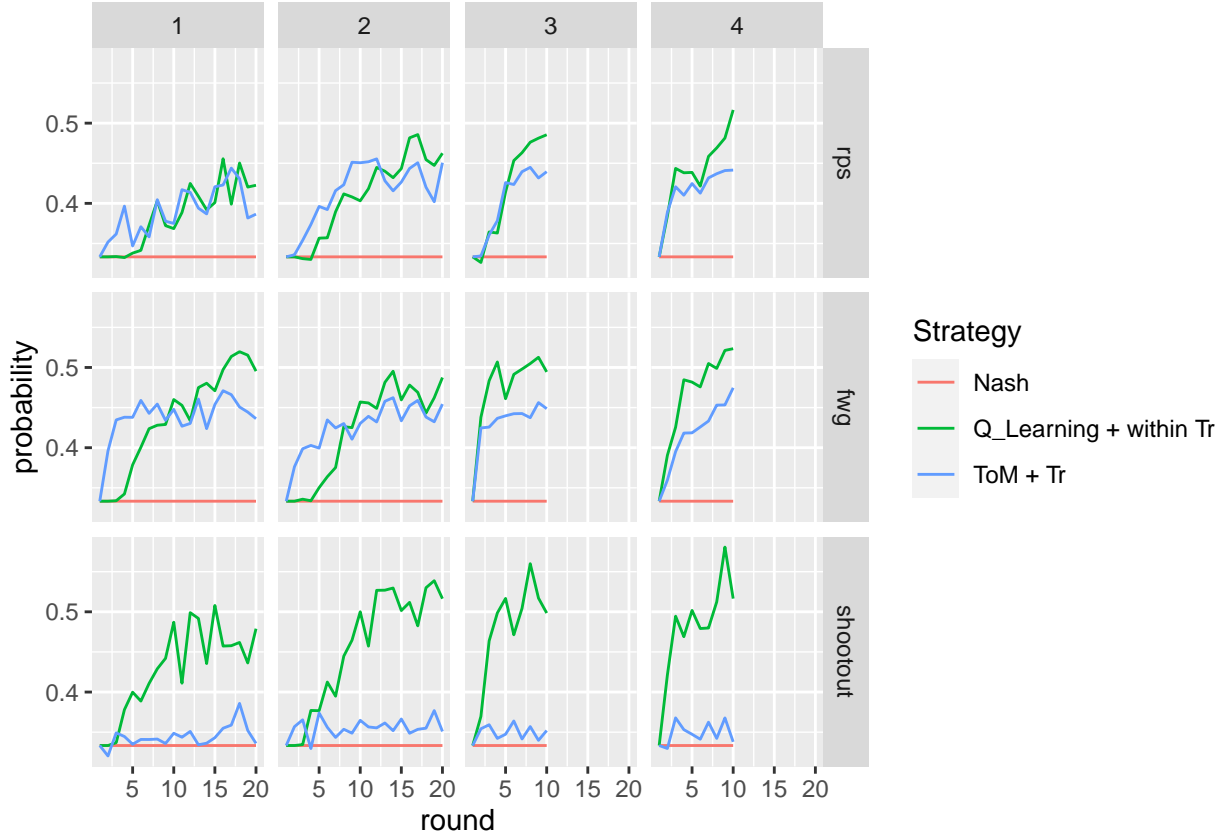
*Figure 13*. Experiment2 likelihood by trial by game and opponent faced

models with the depmixS4 R package. In the first model, we allow for a non-nil probability of players transitioning from one state (strategy) to another. In the second model, we assume that such switching does not happen, and as such assume implicitly that when players start with a particular strategy, they continue using it throughout the experiment. We then compare the likelihoods of each HMM model using a likelihood ratio test.

**Experiment 1:** In experiment 1, the HMM model with switching fits significantly better than the non-switching one ($p < .001$). This is further statistical evidence in favour of the hypothesis that participants switch between strategies. In order to understand at which stage of the games the switching might happen, and whether there are any differences between games and type of opponents faced, we plot in Figure 14 the average (across participants) posterior probabilities of each state (strategy), as a function of trial and
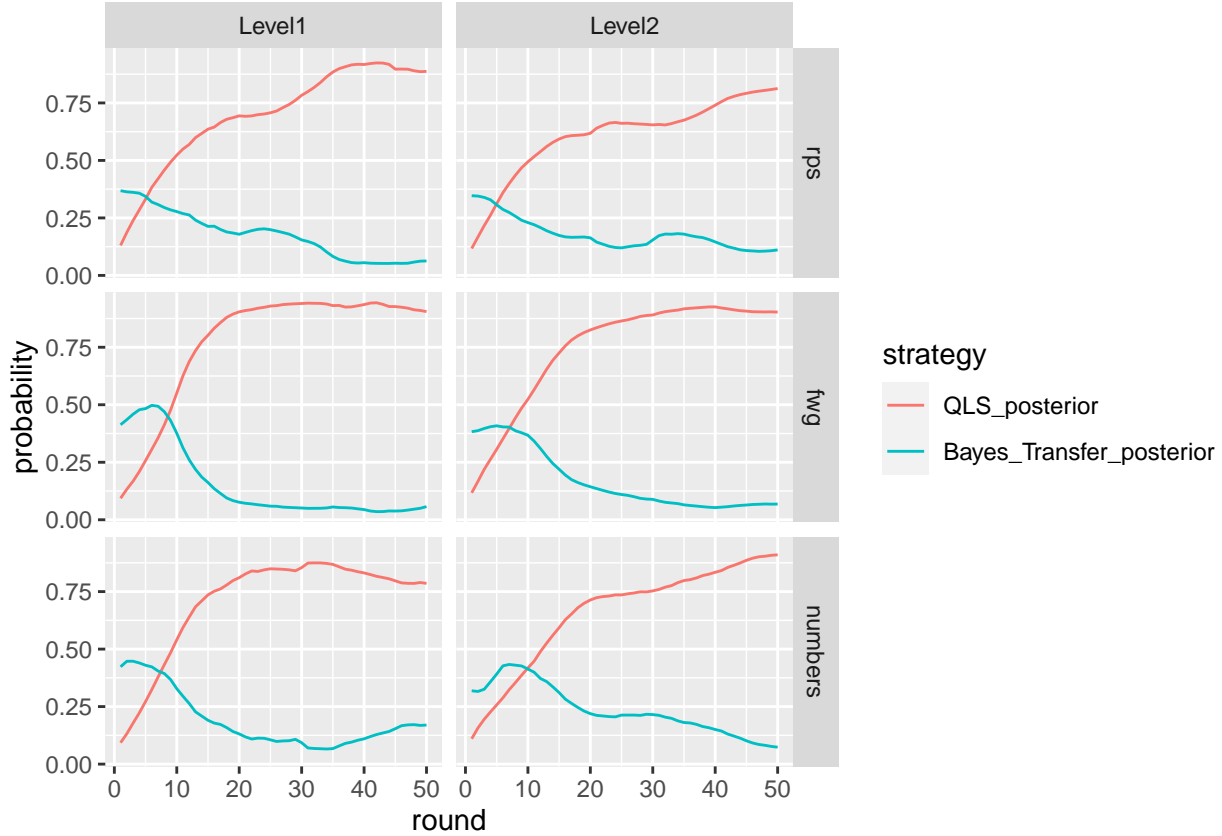
*Figure 14*. Experiment1 posterior probability of strategies by game and opponent faced

<sup>529</sup> opponent faced. The posterior probability is the probability that an observation comes from

<sup>530</sup> a component distribution a posteriori, i.e. given the value of the observation. In the first

<sup>531</sup> experiment, we can see from the plots of fwg and numbers games for level-1 opponent that

<sup>532</sup> although the likelihoods are very close, the posterior probability of the Bayesian model with

<sup>533</sup> transfer is slightly higher than that of the QLS model in the very early rounds, but decreases

<sup>534</sup> rapidly while the posterior probability of the QL-learning with states models keeps

<sup>535</sup> increasing.

<sup>536</sup>     **Experiment 2:**   The switching model in experiment two has also significantly higher

<sup>537</sup> likelihood ($p = 0.00$). On top of indications from looking at the likelihood by trial graphs,

<sup>538</sup> we have therefore further evidence that participants did indeed switch their strategies as the

<sup>539</sup> games progressed. The posterior probability plot in Figure 15 shows switching much more
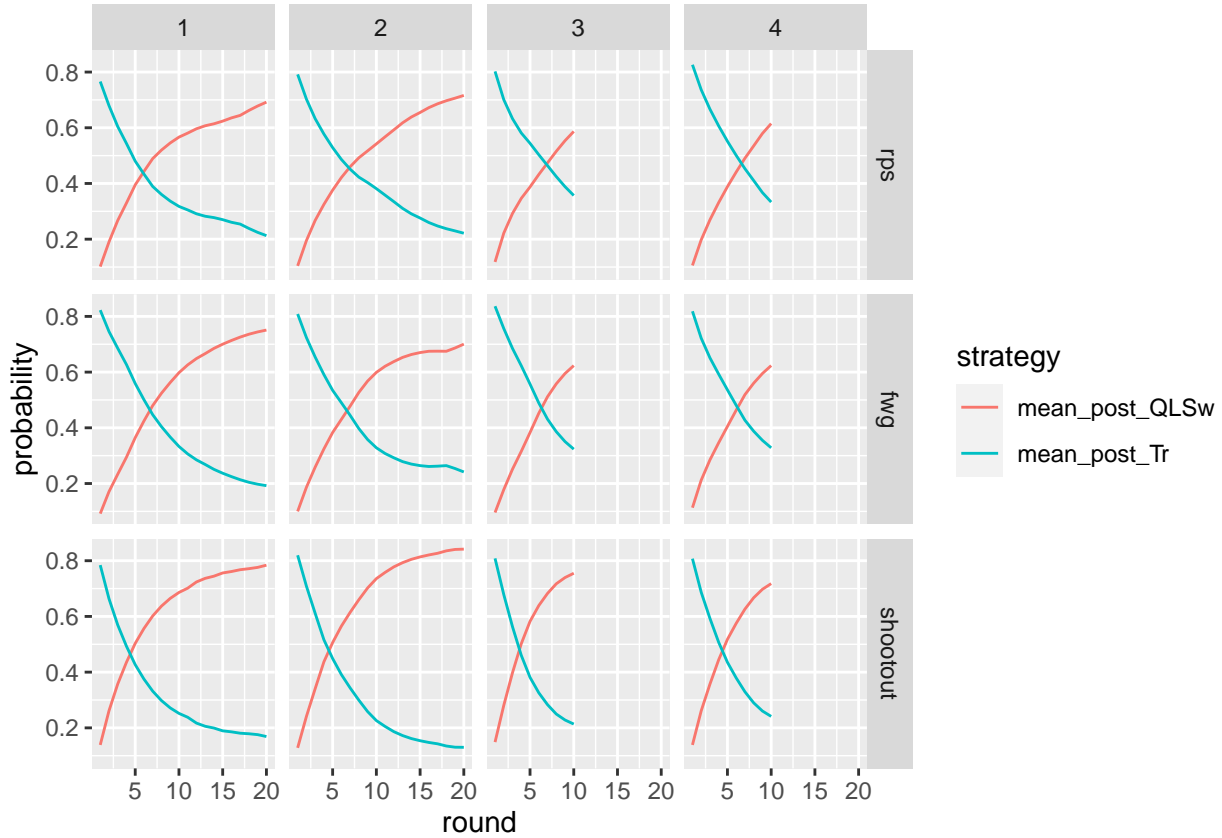
*Figure 15*. Experiment2 likelihood by trial by game and opponent faced

⁵⁴⁰ clearly across games and stages. The switching also seems to happen very early on at the

⁵⁴¹ beginning of each game and stage, and is also consistently in the same direction: The

⁵⁴² probability of Bayesian models with transfer being initially high, then decreasing rapidly

⁵⁴³ while the posterior probability of QL-Learning with states and within transfer learning

⁵⁴⁴ increases rapidly.

⁵⁴⁵     Therefore, HMM modelling shows clear evidence in favour of strategy switching by

⁵⁴⁶ participants, specifically after a few rounds of play. The strategy switching is consistently

⁵⁴⁷ from ToM models towards Q-learning with states models in both experiments.

<sup>548</sup>                                          **Discussion**

<sup>549</sup>        In this study, we investigated human learning transfer across games by making human

<sup>550</sup> participants play against computer agents with rule-based level-k strategies. We were

<sup>551</sup> interested in exploring whether participants learn about the strategy of their opponent and

<sup>552</sup> transfer such knowledge between games, and whether this is modulated by the similarity

<sup>553</sup> between games and the sophistication of the agent.

<sup>554</sup>        The results of our first online experiment show that the majority of participants learn

<sup>555</sup> to adapt to the opponent strategy over multiple interactions and generalise this learning to

<sup>556</sup> the similar game. We used results on very early rounds for measuring transfer as they are

<sup>557</sup> unlikely to be tainted by any within game learning. Using this approach, we showed that

<sup>558</sup> transfer to the more dissimilar game was modulated by the degree of sophistication of the

<sup>559</sup> agent, with evidence for transfer when players face the less sophisticated agent but not the

<sup>560</sup> more sophisticated one.

<sup>561</sup>        In the second online experiment, there were many more opportunities to test transfer

<sup>562</sup> than before: 2 opportunities to transfer opponent models within each game, and a total of

<sup>563</sup> three games, which means 6 opportunities to test transfer. The results on learning transfer

<sup>564</sup> confirmed prior findings from the first experiment. While there was no evidence of higher

<sup>565</sup> scores across interactions within the same game (likely due to the lower number of rounds

<sup>566</sup> per interaction and the higher cognitive load of facing two opponents rather than one), we

<sup>567</sup> found evidence for learning transfer across games as early round scores analysis confirmed.

<sup>568</sup> We also found that learning transfer is modulated by the type of opponent faced. When the

<sup>569</sup> players faced the level-1 opponent, they were able to transfer learning. However, when they

<sup>570</sup> faced the level-2 opponent, there was weaker evidence for transfer. The lack of transfer when

<sup>571</sup> facing the more sophisticated opponent might be due to the difficulty of learning that

<sup>572</sup> opponent strategy to start with. A player cannot transfer what they have not learnt and as

<sup>573</sup> such, since it might be harder to learn the strategy of the level-2 opponent, this in turn

might translate into weaker evidence for transfer.

Coming back to learning transfer, we observed evidence that participants start off new games with prior knowledge as their scores are significantly higher than chance, confirmed both by early stage analysis as well as rounds 2-6 scores analysis. The question we ask ourselves therefore is: What exactly did the players learn in RPS that allowed them to beat the opponent in FWG and Shootout? what did the players learn specifically about the opponent strategy and what form did this learning take?

We will proceed by considering multiple potential answers to this question. First, maybe players simply learn spatial heuristics that allow them to perform better than chance. An example is a spatial heuristic that consists of choosing "weapons" in a particular order. For instance, it is possible to keep winning against a level-k opponent by choosing actions in a particular spatial order such as cycling through them from left to right. This was one of the weaknesses in the design of experiment one, as it was indeed possible using very simple spatial sequences to beat the opponent on most rounds. We took this into account in designing experiment two by randomly shuffling the spatial order of action choices in each round. Still, the learning and conclusions were similar, so this could not explain both learning and transfer in experiments one and two.

A second possible hypothesis for learning the opponent's strategy is the use of simple rules based on last round play (for instance, I play scissors whenever opponent played rock in last round, or whenever the last round play was rock/scissors, I should play paper in this round, etc. . . ). Our Q-learning with states as prior-round play model is a good proxy for this type of strategies. While this approach certainly seemed to be the best fit for some player's behavior, it is unsatisfactory in explaining some of the learning transfer evidence we showed. Indeed, learning the best action in a particular state is not transferable to a new game since the state space is different and there is no single mapping between the state spaces of the initial and latter games. These rules would therefore need to be learned anew in the latter

600 game which is inconsistent with above chance performance in very early rounds.

601 Likewise, assuming that players learn a complete model of the environment (for

602 instance the transition probabilities from last round play to new play) might explain learning

603 within games but is equally unable to account for early games transfer of learning as such

604 models, besides being cognitively very expensive to learn, would require many rounds of

605 practice. Another issue with these hypotheses is that they are not consistent with significant

606 score differences between those facing level-2 and level-1 opponents. More specifically, if we

607 assume that participants were using some type of associative learning or relying on spatial

608 heuristics, then their scores should not depend on the degree of strategic sophistication of

609 the opponent since their approaches would render this variable irrelevant. To be sure, if a

610 participant learns to pick say "scissors" whenever the opponent last picked "rock", then the

611 degree of strategic sophistication of the opponent (its level k) should not impact this

612 learning, and we would expect in this case there would be no difference between scores when

613 facing level-1 and level-2 opponents, which is not the case here. The fact that the degree of

614 sophistication of the opponent matters points to the importance of opponent modelling to

615 successful transfer of learning.

616 We are left with two possible explanations: First, it is possible that the players have

617 uncovered a heuristic that allows them to beat the opponent without explicitly modelling

618 their strategy, and is robust to transfer. Indeed, because of the cyclicality in action choices

619 (e.g : Rock beats Scissors beats Paper beats Rock), it is possible to beat level-2 opponents

620 most of the time by following a simple rule: Play in the next round whatever the opponent

621 played in the last round. This is a rule that wins and is also robust to transfer as it does not

622 depend on action labels and even works in the dissimilar game.

623 The second explanation of learning transfer is that it is driven by a group of

624 participants that are able to build a mental representation of what the strategy of the

625 opponent is. A successful mental representation would take the perspective of the opponent

or endow it with intentionality in order to detect its strategy when the opponent is playing

based on a level-k reasoning model. For instance, the player may think "My opponent is

always trying to be one step ahead of me, therefore, I will be one step ahead of where it

thinks I will be". This mental representation would facilitate the use of theory of mind

abilities and thus enable the players to learn opponent strategies when they are based on

human-like reasoning models such as level-k or cognitive hierarchy. This type of learning

would be deemed "explicit" in the psychology literature as a process through which

knowledge consists of cognitive representations of concepts and rules, as well as the

relationship between them. It involves the evaluation of explicit hypotheses and results in

better problem-solving skills (Mandler, 2004). Since it is less context dependent, this type of

learning is generalizable to new situations, akin to the more general framework of rule-based

learning explored by Stahl (2000, 2003).

Our second experimental design allows us to test whether the first explanation holds.

Since there is a simple transferable heuristic that works against level-2 players, and since as

far as we know, there are no similar ones against level-1 players, if indeed participants were

using this, they would perform better and transfer learning more easily when facing level-2

opponents. Because level-2 opponents use a higher level of strategic reasoning, they should

in fact be harder to play against and in the absence of such a heuristic, performance and

learning transfer should be worse.

Our results show that in fact, it was harder to transfer learning when facing level-2

opponents, both comparing first interactions across games and using early rounds analysis.

Based on our assumptions, we conclude therefore that the most likely explanation is that

participants who are able to beat the opponent and transfer learning are likely to be

explicitly modelling the opponent strategy using level-k reasoning, compared to using simple

learning rules they uncovered during the course of learning.

Our computational modelling allowed us to delve deeper into what might be driving

652 the observed learning transfer. Initial modelling of observations using all available data

653 seems to indicate that the most likely model was a Q-learning type model. However, as we

654 argued above, that would be inconsistent with the evidence for learning transfer. Breaking

655 down likelihoods by trial and fitting a hidden markov model to the data with states being

656 the various strategies that participants are assumed to be using, we showed evidence for

657 within game switching of strategies. Participants start the early rounds of a new game acting

658 in a way consistent with a Bayesian Theory of mind level, which would be accurate and

659 generalisable but computationally expensive. However, as trials continue, participants seem

660 to switch to a habitual type of learning (QL-models).

661 Why is this switching happening? We believe that participants show flexibility in their

662 use of learning strategies. When a new game is started that is similar to a previously played

663 game with the same opponent, participants need a way to transfer prior knowledge of the

664 opponent and apply it to the new game in order to best respond. Adopting a Bayesian

665 model based on ToM achieves the goal of transferring the opponent model and thus coming

666 up with best responses in the early trials. However, Bayesian ToM models are

667 computationally expensive and require higher order thinking (I think that you think that I

668 think. . . ). As such, as the games progresses, they may become burdensome and the higher

669 amount of historical interaction in the new game allows participants to have enough data to

670 start using the cognitively cheaper model-free learning strategies such as Q-learning. The

671 preference for less computationally demanding strategies is well established [@Kool_2011].

672 Moreover, the ability to flexibly switch is also consistent with evidence from the literature on

673 learning strategies in humans, showing that they indeed shift between model-based and

674 model-free learning when the environment requires it [@Simon_Daw_11].

## Conclusion

Our online experiments results are consistent with behavioural game theory findings, in that human players can deviate from Nash equilibrium play and learn to adapt to the opponent strategy and exploit it when the opponent itself is deviating from Nash equilibrium. Moreover, we showed that participants transfer their learning to new games with varying degrees of similarity. The transfer is also moderated by the level of sophistication of the opponent, with participants showing more success in learning and transferring against opponents adopting a less sophisticated strategy.

Having said that, there remains a high degree of heterogeneity between players. There is a high positive association between players who learn to beat the sophisticated and less sophisticated opponents, indicating that some players are more able to detect the patterns in opponent play and learn how to exploit them. Moreover, the computational modelling shows that it is likely that players start each game using a model-based learning strategy that facilitates generalisation and opponent model transfer, but then switch to behaviour that is consistent with a model-free learning strategy as the experiment goes on. This is likely driven by a trade-off between computational complexity and accuracy between model based and model free strategies.

692 # References