# Freedom from habits: the capacity for autonomous behaviour

**Keiji Ota[1], Lucie Charles[1], Patrick Haggard[1]**

[1] Institute of Cognitive Neuroscience, University College London, London, United Kingdom

Corresponding author

Keiji Ota

**Keywords**

Cognitive control / Free action / Habits / Volition / Competitive pressure

# Abstract

The capacity for autonomous behaviour is key to human intelligence, and fundamental to modern social life. However, experimental investigations of the cognitive bases of human autonomy are challenging, because experimental paradigms typically constrain behaviour using controlled contexts, and elicit behaviour by external triggers. In contrast, the sources of human autonomy and freedom are assumed to be endogenous. Here we propose a new theoretical construct of adaptive autonomy, meaning the capacity to make behavioural choices that are free from constraints of both immediate triggers and habitual responding. Participants played a competitive game in which they had to choose the right time to act, in the face of an opponent who punished (in separate blocks) either choice biases, habitual sequences of action timing across trials, or habitual responses to the effects of reinforcement. Adaptive autonomy with respect to each habit was measured by the ability to maintain performance against the opponent even when the corresponding habit was punished. We found that participants were able, under pressure from their opponent, to become free from habitual choices of when to act, but were not able to free themselves from win-stay, lose-shift patterns of reinforcement, even when these resulted in punishment. These results propose a new testing ground of autonomous behaviour as a flexible adaptation of more or less habitual behaviours that co-exist with different classes of external constraint.

# Introduction

Animal behaviour depends upon both exogenous, environmental factors and endogenous factors. The endogenous factors can be conceptualized as a dimension extending from stereotypical behavioural patterns like habits, and flexible, intelligent actions. The latter are thought to play a special role in human autonomy and volition. Adaptability and variation of behavioural choices allows humans to adapt to environmental challenges and find novel solutions. The capacity of human autonomy has been extensively examined in experimental tasks which encourage participants to act freely (Brass & Haggard, 2007; Fleming et al., 2009; Jahanshahi et al., 1995; Libet et al., 1983) or to act randomly (Baddeley et al., 1998; Baddeley, 1966; Jahanshahi et al., 2000) by explicitly telling them to do so. Such instructions in voluntary-action studies invite participants to behave in a way that reflects their understanding of volition and freedom (Haggard, 2008). Outside the laboratory, in contrast, people readily switch between more stereotyped and more autonomous behaviours without explicit instruction, as a function of multiple situational and internal factors.

As such, it remains unclear the extent to which people can express their behavioural autonomy through volitional actions. Competitive games might be a good testing ground for this question, for two reasons. First, many competitive games require people to initiate an action endogenously, rather than in response to an external stimulus. For example, in the 'rock, paper, scissors' game, each participant selects an action without first seeing the action of their opponent. This stimulus-independence is considered a necessary condition for volition (Jenkins et al., 2000). Second, volitional actions are often contrasted with habitual, or routine actions (Haggard, 2019). Competitive games offer a convenient way to manipulate the extent to which any individual action is or is not habitual. For example, if a player behaves habitually in a competitive game, their opponent will be able to predict their upcoming choice, and adjust their strategy accordingly. Therefore, an agent playing a competitive game should avoid habitual or exploitative behaviours, and must innovate in order to avoid being predicted. Non-human primates indeed respond to competitive pressure by initiating exploratory behaviour (Barraclough et al., 2004; Lee et al., 2004; Lee et al., 2005).

The present study therefore investigates human volition and autonomy in a competitive game task in which participants could not react to their competitor's current move (stimulus independence), and additionally could receive reward only when they avoided the competitor's prediction (habit independence). Further, we used several virtual competitor algorithms, each one designed to punish a particular kind of habit. Whereas psychologists have often thought of

65   habits as personality traits, we reasoned that a person may stop behaving habitually when a

66   competitor begins to predict, exploit and punish their habitual beahviour. Thus, the *change* in

67   habitual behaviour under competitive pressure offers a quantitative measure of individual

68   autonomy with respect to habits.

69        We conceptualised three distinct "habitual families"[1]. The first was automatic response

70   selection (Dolan & Dayan, 2013; Du et al., 2022; Robbins & Costa, 2017). We will refer to this

71   as *standard choice habits*. Consider the simple task of generating one of three digits in each

72   turn (see Figure 1A). Agent X may prefer to choose "1", for whatever reason, while agent Y

73   may be less biased. The statistical similarity between an individual's observed choice pattern

74   and a random pattern can be measured (depicted by the right arrow in Figure 1A). Suppose

75   now that a competitor punishes X for repeating one choice within a game scenario. If X can

76   break her habit, she should now choose the two other digits more often (shared area in Figure

77   1A). Agent Y may be less adaptive and stick to his original choice pattern. This adaptive

78   capacity may reflect the autonomy each agent has over their choice habits (left arrow in Figure

79   1A). We call this quantity *adaptive autonomy*.

80        The second habit family we considered is *transition habits*. This refers to action chains

81   or routines (Lashley, 1951; Robbins & Costa, 2017; Rosenbaum et al., 2007). In our example,

82   integer counting ("1, 2, 3") is such a habit (agent X in Figure 1B). The only way to completely

83   avoid such habits is to generate each choice independently from the previous trial. Yet studies

84   of random number generation show people find this difficult (Baddeley et al., 1998; Baddeley,

85   1966; Bar-Hillel & Wagenaar, 1991). The shaded area and left arrow in Figure 1B illustrate

86   potential behavioural adaptation to punishing transition habits.
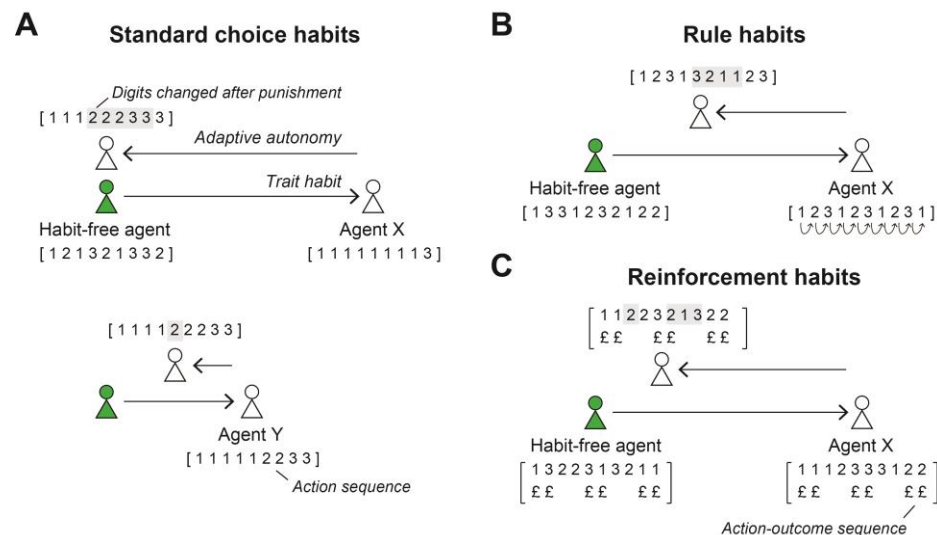
87        Lastly, we address how people respond to action successes and failures, by

---

[1] The concept of a habit has been discussed over centuries by philosophers (Barandiaran & Di Paolo, 2014). Its theoretical framework has been formalised by psychologists and neuroscientists' works, but it is still controversial (Dolan & Dayan, 2013; Du et al., 2022; Robbins & Costa, 2017). A psychologist's view of a habit is that it reflects the formation of stimulus-response associations (Wood & Runger, 2016). A traditional testing ground for habits is a reward devaluation paradigm. Here a reward previously assigned to a stimulus-response association is devalued (Robbins & Costa, 2017). If the response is still automatically evoked by the stimulus, its behaviour is said to be a habit rather than goal-directed. In the present study, we do not aim to examine whether or how habits are formed. Rather, we investigate how people become liberated from their endogenous habitual patterns. The three habit families we conceptualised are considered different expressions of habits. From the perspective of a stimulus-response association, standard choice habits are any obvious responses in any task. Rule habits are those responses evoked by the last response. While reinforcement habits are the responses elicited by the last outcome/feedback. The precise mechanism by which each habit forms is unimportant. Rather, the concept of three habit families refers to statistical patterns that are not random nor independent (see below) and that may be continually shaped by ongoing experience. Thus, volition in our paradigm requires the regulation of salient/habitual patterns and the exploration of new behavioural patterns.

88    considering *reinforcement habits*. In Figure 1C, agent X show typical win-stay lose-shift
89    behaviour in a digit generation task. In contrast, a habit-free agent generates each choice
90    independently from whether the previous outcome was rewarded or not. The vast majority of
91    studies in reinforcement learning assume that a 'win-stay lose-shift' strategy is natural, or even
92    unavoidable (Worthy et al., 2013). Here we test whether people can unlearn this familiar
93    reinforcement habit when it is punished by a competitor. Adaptive autonomy would mean that
94    an agent would be able to break the association between their next action and the previous
95    outcome (see the potential change in Fig. 1C).

96         In this experiment, we designed a structured series of competitors in a game scenario,
97    in order to selectively punish these specific habits, and measure individuals' capacity for
98    adaptive autonomy, as the change in behaviour when a specific habit family was punished.
99    Using this framework, we tested the capacities or limits of human autonomy for three "habitual
100   families" of choice, transition and reinforcement. To explore whether adaptive autonomy
101   reflects a general capacity, or rather is specific to a particular habit family, we explored
102   correlations across individuals in our adaptive autonomy measures for each habit. Finally, we
103   modelled the learning process by which people generated a new action in order to avoid the
104   competitor.



**Figure 1.** Three habitual families in a hypothetical experiment. An agent is asked to generate
one digit from three in each turn. A sequence of generated digits is shown in a square bracket
from left to right. The precise task is not important. **A**. Agent X habitually selects the digit "1"
while agent Y occasionally selects the other two digits. A habit-free agent should select each of
the 3 digits randomly. A right arrow indicates a pattern similarity from the habit-free agent to

111  agent X or Y. **B.** Agent X has a trait rule-based transition habit, in which they count up from the

112  last digit. The habit-free agent selects the digit independently from the previous digit. **C.** In

113  certain situations where an agent is rewarded, the pattern that is dependent on the reward

114  assigned captures a trait reinforcement habit. In all habitual families, we measure the extent to

115  which selection patterns change when agents are punished for habitual action. Shaded areas

116  represent digits changed after the punishment and hypothetical agents move their location

117  closer to the habit-free agent. A left arrow illustrates adaptive autonomy, a change in

118  behavioural patterns with respect to habits.

119

## 120 Results

### 121 Experimental task

122  Participants were asked to decide when to press a key that caused some food to be delivered

123  to a storage location. They were competing with a virtual competitor, represented as a flock of

124  birds (Fig. 2A). The birds tried to catch the food during the delivery process, by deciding when

125  to fly out of a tree and across the field. The participant's task was to deliver the food without it

126  being caught by the birds. We programmed the birds to predict the time of the participant's next

127  action based on the history of their reaction times. Based on this prediction, the birds made a

128  choice of when to fly on each trial. They flew at a time that was designed to intercept the food

129  thrown by a participant within one of three intervals: 1) early throw (0–1.5 sec), 2) middle throw

130  (1.5–3.0 sec) or 3) late throw (3.0–4.5 sec). The participants could win a trial by pressing the

131  key during one of two intervals that the birds did not select. The intervals were not explicitly

132  demarcated for the participant, who experienced a continuum of potential action times in each

133  trial. There was no time for participants to perform the task reactively because the birds could

134  travel much faster than the food. If the participants simply waited for a moment when no birds

135  flew and then threw, the birds could suddenly appear and intercept the food. Therefore, the

136  participants were asked to predict when the birds would appear and avoid them. This feature

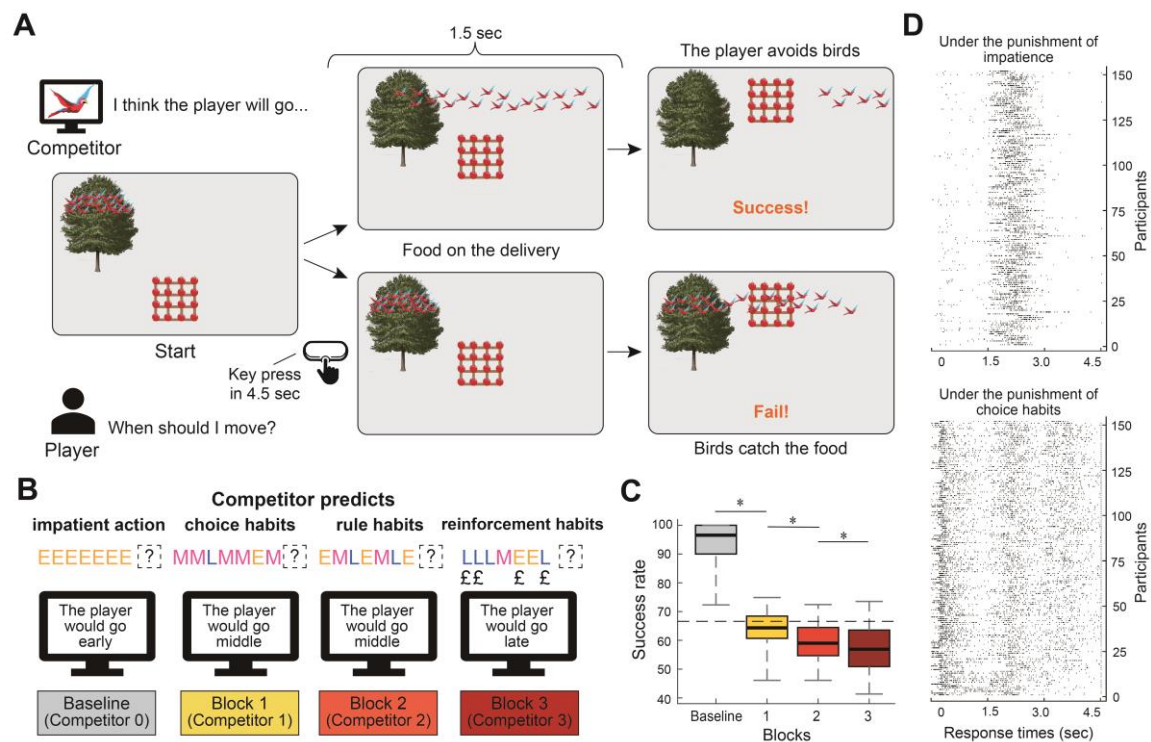137  means that our participant's actions were stimulus-independent.

138  There were 4 blocks in total. In each block, the participants competed with a class of

139  competitor that pressurised a specific habit (Fig. 2B). In the baseline block, Competitor 0 was

140  programmed to punish the participants for being impatient: the birds consistently punished a

141  participant who threw in the early interval, so the participant was incentivized to wait to avoid

142  being intercepted. In block 1, Competitor 1 punished standard choice habits, if a participant

143  selected one interval more often than the other two. In block 2, Competitor 2 predicted

144    transition habits, and punished any association between the time of the participant's current
145    throw and the time of the preceding throw. Finally, in block 3, Competitor 3 punished
146    reinforcement habits by seeking out whether the time of the current throw was associated with
147    both the time of the preceding throw and the preceding outcome. Thus, the participants played
148    against competitors who had increasingly sophisticated predictive power in each block. The
149    participants required progressive degrees of autonomy across blocks: they needed to act in a
150    way that was even more unconstrained than required by the competitors they had played
151    previously.

152    Participants did not receive any explicit instruction or explanation about what habits
153    they should avoid. The participants were never told when they should act on any given trial.
154    Instead, they could only monitor the success/failure of avoiding the birds on each trial, and
155    adapt their behaviour accordingly to try to avoid the birds on future trials. Thus, successful
156    performance under different punishment regimes would depend on implicit mechanisms of
157    adaptation rather than explicit understanding.

158    We first examined whether the predictive power of our protocol increased by checking
159    the percentage of successful bird-avoiding trials. The participants achieved near perfect
160    success rates against Competitor 0 who punished impatience (Fig. 2C; Median [Mdn] = 96.6%).
161    The participants avoided an immediate response and initiated the throw 1.5 seconds after the
162    trial starts on almost all trials (Fig. 2D upper panel). In block 1, the success rate dropped to
163    66.6%, as would be expected from purely stochastic choices (Fig. 2C; Mdn = 64.3%, $p < .001$,
164    $z = 10.69$ for blocks 0 versus 1, Wilcoxon sign rank). The success rate further decreased in
165    block 2 (Fig. 2C; Mdn = 59.0%, $p < .001$, $z = 5.74$ for blocks 1 versus 2, Wilcoxon sign rank)
166    and even further in block 3 (Fig. 2C; Mdn = 56.9%, $p = .015$, $z = 2.44$ for blocks 2 versus 3,
167    Wilcoxon sign rank). Therefore, our progressive series of punishments increasingly stressed
168    the participants' cognitive demands for avoiding the competitor.

169

**Figure 2.** Virtual competitive environments for penalising habitual actions. **A**. A trial sequence. A participant decides when to throw food within a 4.5 second time window. The food is delivered at the top-centre of the screen 1.5 sec after a key press. A virtual competitor (i.e., a flock of birds) attempts to intercept the food by adjusting the time at which it flies out of a tree. Participants win a trial if they avoid being caught by the birds. **B**. Experimental (game) design. The virtual competitor predicted which time interval participants would initiate the delivery based on their past behaviour. On each trial, the competitor punished one of three intervals, 1) early throw (0–1.5 sec), middle throw (1.5–3.0 sec) or late throw (3.0–4.5 sec). An example sequence of action intervals is shown. In the baseline block, the early interval, associated with impatience, was punished. In block 1, *standard choice habits* that favoured one interval over all others were punished (e.g., the middle interval). In block 2, *transition habits* (i.e., sequential pattern) were punished. For example, if the participant went early, middle, late, early, middle and late, the last early action would likely prime the next middle interval. In block 3, *reinforcement habits* (i.e., outcome dependence) were punished. In the example, the repetition of the same interval likely follows from a reward while a change in interval likely follows from no reward. The rewarded last late action would prime the next late interval. **C**. Success rate of avoiding the birds against each class of competitor. A dashed line denotes the chance level. For each box, the central mark represents the median, the edges of the box are the 25th and

189    75th percentiles and the whiskers are the 2.5th and 97.5th percentiles. * $p$ < significant level

190    after Bonferroni correction, Wilcoxon signed rank. N = 152. **D**. Response times before the

191    punishment of choice habits in the baseline block (upper panel) and during the punishment in

192    block 1 (lower panel). Each small dot represents a reaction time in each trial. The response

193    time data for each participant are aligned in each column.

194

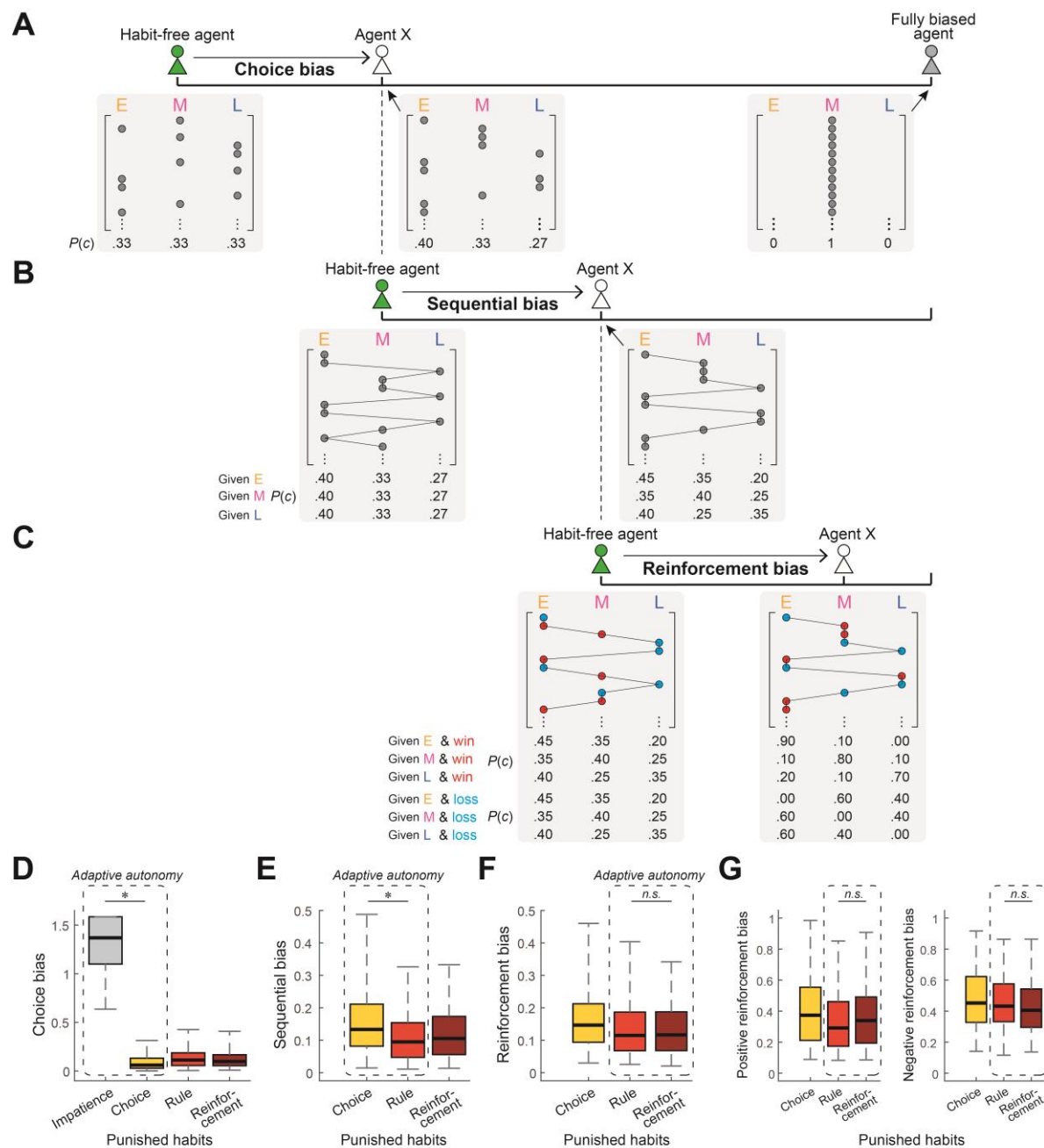195    **Do people adapt to punishments of habitual actions?**

196    We next examined the extent to which people could adapt to punishment of different

197    habits. We therefore developed a measure that reflects an individual's tendency towards a trait

198    habit in each habit family. We measured the statistical distance (Kullback-Leibler [K-L]

199    divergence) between the observed probabilities of selecting the early, middle and late intervals,

200    and the probabilities that a habit-free agent would exhibit (right arrow in Figure 3A-C). A lower

201    statistical distance means that the individual is close to the habit-free agent in terms of his

202    choice profile. We call this quantity a decision bias. A decision bias score *before* the

203    punishment of a specific habit reflects an individual's trait habit regarding when to act. We then

204    looked at the *change* in bias score when a given habitual behaviour was punished. We

205    quantified the decision bias score for standard choice habits, transition habits and

206    reinforcement habits, respectively (Figure 3A-C). A greater change in bias score would indicate

207    stronger *adaptive autonomy*, or ability to modulate the trait habit. See *Data analysis* for

208    detailsData analysis.

209    Looking at response times, the participants began to distribute action times

210    appropriately when Competitor 1 started punishing choice habits (Fig. 2D). Accordingly, their

211    choice bias—a statistical distance of the observed choice probabilities from probabilities 0.33

212    (Fig. 3A)—reduced after punishment (a dashed rectangle in Fig. 3D; Mdn = 1.37 for the

213    punishment of impatience versus Mdn = 0.06 for the punishment of choice habits, $p$ < .001, $z$ =

214    10.69, Wilcoxon sign rank). Competitor 1 did not seek transition patterns from one action to the

215    next and allowed participants to still use *transition habits*. We quantified the sequential bias by

216    considering the extent to which probabilities given the previous action are explained by one's

217    choice probabilities (Fig. 3B). We found that the sequential bias decreased after the competitor

218    pressurised transition habits (a dashed rectangle in Fig. 3E; Mdn = 0.13 for the punishment of

219    choice habits versus Mdn = 0.09 for the punishment of transition habits, $p$ < .001, $z$ = 3.82,

220    Wilcoxon sign rank). Against Competitor 3, the participants were asked to act even more freely

221    to avoid *reinforcement habits*. We evaluated the reinforcement bias by considering the extent

8

222   to which probabilities given the previous action and the previous outcome are explained by
223   probabilities given the previous action solely (Fig. 3C). The reinforcement bias did not show a
224   significant improvement (a dashed rectangle in Fig. 3F; Mdn = 0.11 for the punishment of
225   transition habits versus Mdn = 0.11 for the punishment of reinforcement habits, $p$ = .79, $z$ =
226   0.27, Wilcoxon sign rank).

227   We further tested the possibility that the participants adapted differently to the influence
228   of positive and negative reinforcements since the neural process after a positive outcome
229   stimulus is different from that after a negative outcome stimulus (Gehring & Willoughby, 2002;
230   Hajcak et al., 2006; Vickery et al., 2011), leading to a stereotypical win-stay lose-shift strategy
231   (Wang et al., 2014). We quantified the positive reinforcement bias and the negative
232   reinforcement bias separately (Suppl. Fig. 1). Nevertheless, we did not find significant
233   improvements in the positive reinforcement bias (a dashed rectangle in Fig. 3G; Mdn = 0.29 for
234   the punishment of transition habits versus Mdn = 0.34 for the punishment of reinforcement
235   habits, $p$ = .21, $z$ = -1.26, Wilcoxon sign rank) nor in the negative reinforcement bias (a dashed
236   rectangle  in Fig. 3G; Mdn = 0.43 for the punishment of transition habits versus Mdn = 0.41 for
237   the punishment of reinforcement habits, $p$ = .19, $z$ = 1.30, Wilcoxon sign rank). These results
238   suggest that people are able to become more autonomous from standard habitual choices and
239   habitual action transitions but cannot break away from outcome dependencies. That is, people
240   display habits of being guided by reinforced, such as win-stay lose-shift, even when they are
241   discouraged from doing so.
242

9

**Figure 3.** Measuring trait habits and adoptive autonomies in each habit family. In the raster plot, potential actions for the early, middle and late intervals are shown. The numerical values underneath the plot are the probabilities for choosing these three actions. A right arrow illustrates the statistical distance (i.e., pattern similarity) between a habit-free agent and a hypothetical agent X given their choice probabilities. This distance (or decision bias score) is a

250   proxy for an individual's trait habit. The lower the distance, the smaller the bias. **A.** Choice bias
251   (a proxy for *standard choice habits*). A habit-free agent would select each interval equally often
252   while a fully biased agent would select one interval on every trial. The profile of a participant's
253   choices should be somewhere in-between. **B.** Sequential bias (a proxy for *transition habits*). A
254   habit-free agent is underneath agent X in the panel (**A**) to have the identical choice
255   probabilities to him. The choice probabilities might be larger than the random probabilities 0.33.
256   However, the agent who is free from transition habits would choose their next interval
257   independently of their previous interval: whether the previous interval was early, middle or late
258   does not affect the probabilities of the next interval. Any deviations from such independent
259   choice patterns are considered residual transitions from which participants cannot break
260   (depicted by a right arrow). **C.** Reinforcement bias (a proxy of *reinforcement habits*). A habit-
261   free agent is underneath agent X in the panel (**B**) to have the identical conditional probabilities
262   to him. The agent who is free of reinforcement habits would choose their next interval
263   independently of the previous outcome: whether the previous outcome was success or failure
264   does not affect the probabilities of the next interval. Any such independent choice patterns are
265   considered residual outcome dependencies that participants cannot break (depicted by a right
266   arrow). **D-G.** A dashed rectangle highlights adaptive autonomy as the theoretical difference
267   between a pre-punishment and a post-punishment. * $p$ < significant level after Bonferroni
268   correction, Wilcoxon signed rank. On each box, the central mark represents the median, the
269   edges of the box are the 25th and 75th percentiles and the whiskers are the 2.5th and 97.5th
270   percentiles.

271

272   **Is there a common factor underlying adaptive autonomies?**

273        A change in decision bias scores between the pre-punishment and the post-
274   punishment phase provides a measure of adaptive autonomy for each habit family (dashed
275   rectangle areas in Fig. 3D-G). We considered the domain-general mechanism of cognitive
276   control, which proposes that proactive, strategic cognitive control shares its control mode
277   across tasks that recruit different cognitive elements (Braver, 2012; Braver et al., 2007; Tang et
278   al., 2022). If adaptive ability to become free of a specific habit (e.g., standard choice habits)
279   generalises to adaptive ability to become free of another habit (transition habits or
280   reinforcement habits), for example because both depend on a common, domain-general
281   mechanism, then we would find a correlation between measures of adaptive autonomy elicited
282   by different types of punishments. Looking at the correlation structure, we did not find strong or

283    even moderate correlations among them, in our sample of 152 participants (Figure 4A). This

284    suggests that the ability to voluntarily regulate one habit is not associated with the ability to

285    regulate another habit. This also suggests that our measurements are separable and evaluate

286    three distinct forms of autonomy conceptualised above, rather than a single common form of

287    autonomy. We also checked the correlation structure between the adaptive autonomy of

288    positive reinforcement bias and that of negative reinforcement bias. There was no strong

289    correlation (Figure 4B). This suggests that the ability to adapt away from a win-stay type

290    behaviour is not associated with the ability to adapt away from a lose-shift type behaviour

291    across participants. To summarise, people seem to recruit distinct cognitive capacities for

292    autonomy when unlearning standard choice habits, transition habits and reinforcement habits.

293

294



296    **Figure 4.** Correlation structure underlying measures of adaptive autonomy. **A.** A measure of

297    adaptive autonomy for each habit family (choice, transition and reinforcement) is quantified as

298    a change in decision bias scores between the pre-punishment and the post-punishment phase

299    (dashed rectangle areas in Fig. 3D-F). In a sample of 152 participants, there is no strong

300    correlation among three measure of adaptive autonomy, suggesting an adaptation is specific to

301    a particular habit family. **B.** A measure of adaptive autonomy for a positive reinforcement bias

302  is plotted against that for a negative reinforcement bias, quantified as a change in decision bias

303  scores (dashed rectangle areas in Fig. 3G).

304

**A learning process that accounts for behavioural autonomy**

306  What mechanisms could explain how participants adapted to these pressures? How

307  did they learn actions that successfully avoided their competitor? In competitive games, two

308  strategies can be taken to sustain performance. One strategy is stochastic selection by tossing

309  a coin. This mixed strategy helps with unpredictability but does so without interacting with the

310  environment or competitor. An alternative strategy attempts to predict the opponent's next

311  action based on a history of their prior actions (Hampton et al., 2008; Zhu et al., 2012). This

312  strategy is called belief learning (Camerer, 2003) – broadly speaking, the strategy adopted by

313  the virtual competitor is considered belief learning. Belief learners employ an element of

314  mentalizing because they engage in a representation of the actions and intentions of their

315  opponent (Amodio & Frith, 2006; Hampton et al., 2008). Because our task is designed to

316  produce stimulus independence, a BL strategy is a predictive way of achieving stimulus

317  independence (i.e., predict when the birds would fly and avoid them). Thus, we can reason

318  that, as a learning pathway, the participants might learn 1) a feedback-independent stochastic

319  selection strategy, or 2) a feedback-based belief leaning strategy.

320  We simulated play to test whether the BL strategy is effective at avoiding competitor's

321  predictions. We calculated the reward obtained from simulated choices that BL agents made

322  (see *Simulated play*). We also computed the simulated success rate of a simpler strategy,

323  reinforcement learning, which selects the action that was the most rewarded (Sutton & Barto,

324  2018). In principle, an RL agent knows whether their choice was rewarded or not and repeats

325  the most rewarded action until it is punished. In contrast, a BL agent knows which option a

326  competitor chose and which options they did not. For instance, if the birds intercepted the early

327  throw, a BL agent reduces the value of the early throw. At the same time, a BL agent can also

328  increase the values of the middle and late throws because the birds did not choose to

329  intercept. An RL agent can only update the value of the throw based on whether it was

330  successful or not. This difference in the internal processes allows a BL agent to update the

331  values of the options quickly, thereby prompting a frequent update of the best option. In the

332  simulated play, both RL and BL strategies achieved sufficient success rates under the

333  punishment of impatience (Fig. 5A; Mdn = 96.7% for RL versus Mdn = 98.3% for BL).

334  However, the BL strategy sustained a higher chance of winning than the RL strategy under the
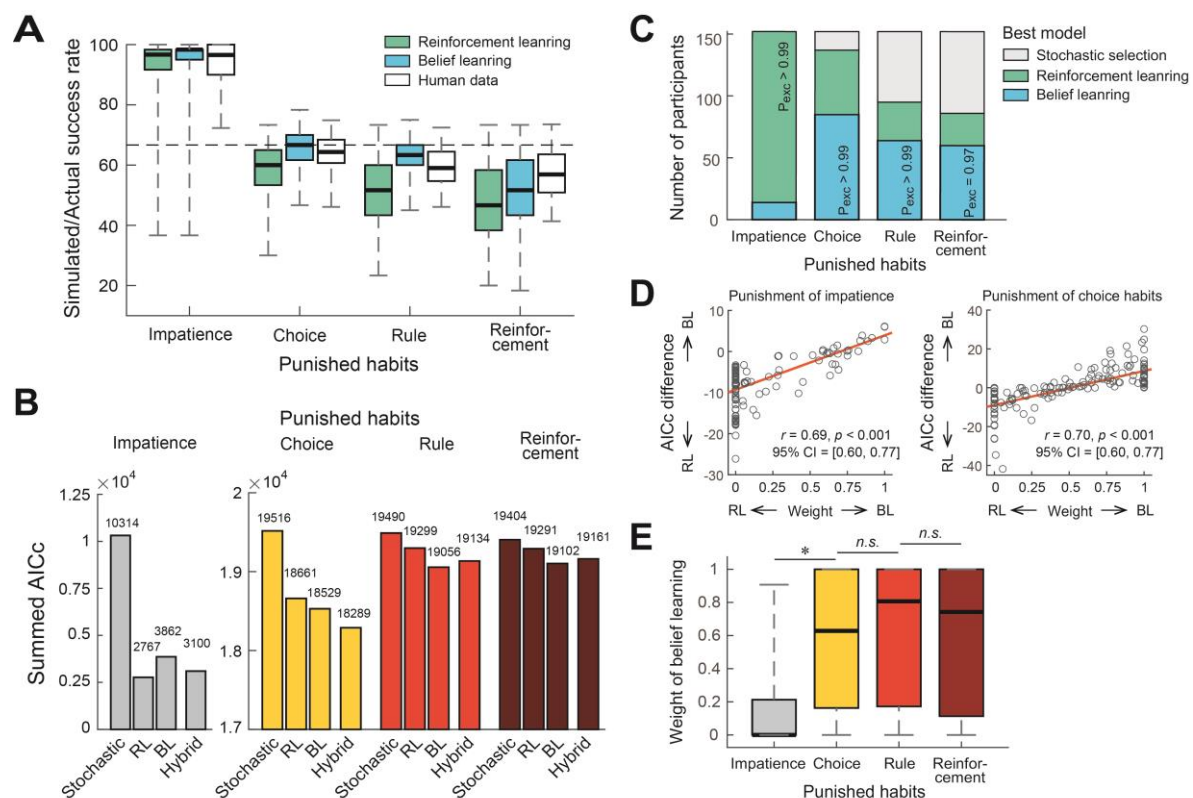
13

335   punishment of choice habits (Fig. 5A; Mdn = 60% for RL versus Mdn = 66.6% for BL) and

336   under the punishment of transition habits (Fig. 5A; Mdn = 51.6% for RL versus Mdn = 63.3%

337   for BL). The use of belief learning might account for participants' adaptions to explore new

338   actions.

339       To address which of the above strategies (stochastic selection, RL and BL) captured

340   our participants' behaviour, we fitted each model to their choice data (see *Computational*

341   *models*). We also fitted a hybrid learning rule (Camerer & Ho, 1999; Hampton et al., 2008; Zhu

342   et al., 2012) which combines reinforcement learning and belief learning. The stochastic

343   selection model had parameters that captured a participant's choice preference and

344   uncertainty, while the RL and BL models added a learning rate from feedback information. In

345   this nested structure, a better model index relative to the stochastic model would indicate the

346   presence of a feedback process. Across four blocks, we found a clear transition in the best-

347   fitting model (Fig. 5B). According to the summed AICc across participants, the RL model

348   outperformed the stochastic model and BL model under the punishment of impatience.

349   However, the BL outperformed when the competitor started punishing choice habits. The BL

350   model still outperformed when the competitor punished transition habits and reinforcement

351   habits (Fig. 5B). For each model we calculated the number of participants best fit by the model

352   (Fig. 5C) and the protected exceedance probability from the group-level Bayesian model

353   selection (Rigoux et al., 2014; Stephan et al., 2009), which is an omnibus measure of the

354   probability that the model is the best model among tested models. The protected exceedance

355   probability for the RL model to outperform the stochastic model and BL model was close to

356   100% under the punishment of impatience. The protected exceedance probability for the BL

357   model was close to 100% under the punishment of choice habits and transition habits, and this

358   was 97% under reinforcement habits (Fig. 5C).

359       As such, the hybrid learning rule fitted the data relatively well in all blocks (Fig. 5B). We

360   recovered the estimates of a relative contribution of belief learning over reinforcement learning

361   from the model fit for the hybrid rule (see *Computational models*). We first checked the

362   robustness of our estimates of the hybrid model by correlating the estimated relative weight

363   parameter to the difference in the AICcs between the RL only model and BL only model. We

364   found a strong positive correlation: the larger the weight placed on BL in the hybrid model, the

365   better the BL only model is (Fig. 5D). We then checked the estimates of the weight parameters

366   across blocks. The weight increased from punishment of impatience to punishment of choice

367   habits (Fig. 5E; Mdn = 0.00 for the punishment of impatience versus Mdn = 0.63 for the

14

368    punishment of choice habits, $p < .001$, $z = 8.32$, Wilcoxon sign rank). That is, belief learning

369    made an important contribution to participants' choices when they were punished for choice

370    habits (Fig. 5E; Mdn = 0.63), transition habits (Mdn = 0.81; versus choice habits, $p = .10$, $z =$

371    1.65) and reinforcement habits (Mdn = 0.74; versus transition habits, $p = .96$, $z = 0.04$),

372    respectively. These findings suggest a shift in learning strategies that followed the demands of

373    the competition: participants first used reward-guided behaviour when it sustained the success

374    rate. Then, once the competitor started predicting habit patterns, participants switched to

375    learning successful actions from the opponent's prior actions.

376



377

**Figure 5.** A shift in the strategy to belief learning as competitive demand increases. **A**. Real
success rate (white bars) in the actual experiment and fictive success rate (green or blue bars)
in simulated play. Agents using the reinforcement learning strategy and agents using the belief
learning strategy competed against each class of competitor. **B.** Summed AICc across
participants. Lower values of AICc are better. Stochastic: stochastic selection model. RL:
reinforcement learning model. BL: belief learning model. Hybrid: hybrid learning rule. **C.** The
number of participants best fit by the model. The $P_{exc}$ inserted in the bar denotes the protected

385    exceedance probability that supports the corresponding model. **D.** Relative weight placed on

386    belief learning over reinforcement learning captured by the hybrid model versus the AICc

387    difference between the RL model and BL model. As the weight parameter increases, the model

388    fit of belief learning improves relative to that of reinforcement learning. **E.** A transition in the

389    relative weight placed on belief learning. * $p <$ significant level after Bonferroni correction,

390    Wilcoxon signed rank. **A&E.** On each box, the central mark represents the median, the edges

391    of the box are the 25th and 75th percentiles and the whiskers are the 2.5th and 97.5th

392    percentiles.

393

394    **A shift in learning strategies enhances behavioural autonomy**

395         Because belief learning is a faster learning process than reinforcement learning, our

396    simulated play shows that belief learning indeed induces a smaller choice bias than simple

397    reinforcement learning (Suppl. Fig. 2), leading to better performance (Fig. 5A). We therefore

398    examined whether the shift in learning processes accounts for achieving a smaller choice bias

399    in the participants' data. To this end, we used a bivariate latent change score (LCS) model

400    (Carpenter et al., 2019; Kievit et al., 2018; Kievit et al., 2017; McArdle, 2009). LCS models

401    conceptualize the change in score between one time point (before punishment) and the next

402    time point (under punishment) as a latent change factor (see *Latent change score model*).

403    Under a bivariate LCS model, two factors influence the change score. The first is the extent to

404    which the reduction in bias is explained by the initial bias before punishment, which is termed

405    auto-regression in LCS models. We might expect to see a negative auto-regressive effect as a

406    consequence of a scale attenuation: individuals who started off with a larger bias score could

407    potentially have a greater reduction in the bias. While individuals those who started off with a

408    smaller bias score could have a smaller reduction in the bias because of the lower limit of the

409    scale. The second is the extent to which the reduction in bias is explained by the initial weight

410    on belief learning, which is termed cross-coupling. A bivariate LCS model would reveal a

411    negative cross-coupling effect if individuals who attempted to mind-read the competitor's

412    strategy (i.e., having a large BL weight before punishment) gained a greater reduction in the

413    bias. Moreover, by having two latent change factors, a bivariate LCS model estimates

414    correlated change: the degree to which the reduction in bias co-occurs with the change in

415    weight (i.e., a shift in learning strategies). If participants adapted their behaviour to gain a

416    smaller bias score by simply behaving randomly, the reduction in bias would not co-occur with

417    the change in weight. If, on the other hand, participants achieved a smaller bias score by

418    initiating the prediction of the competitor's strategy, these two changes would co-occur, and
419    then the bivariate LCS model would reveal a negative correlated change: gaining a greater
420    weight is associated with gaining a smaller bias.
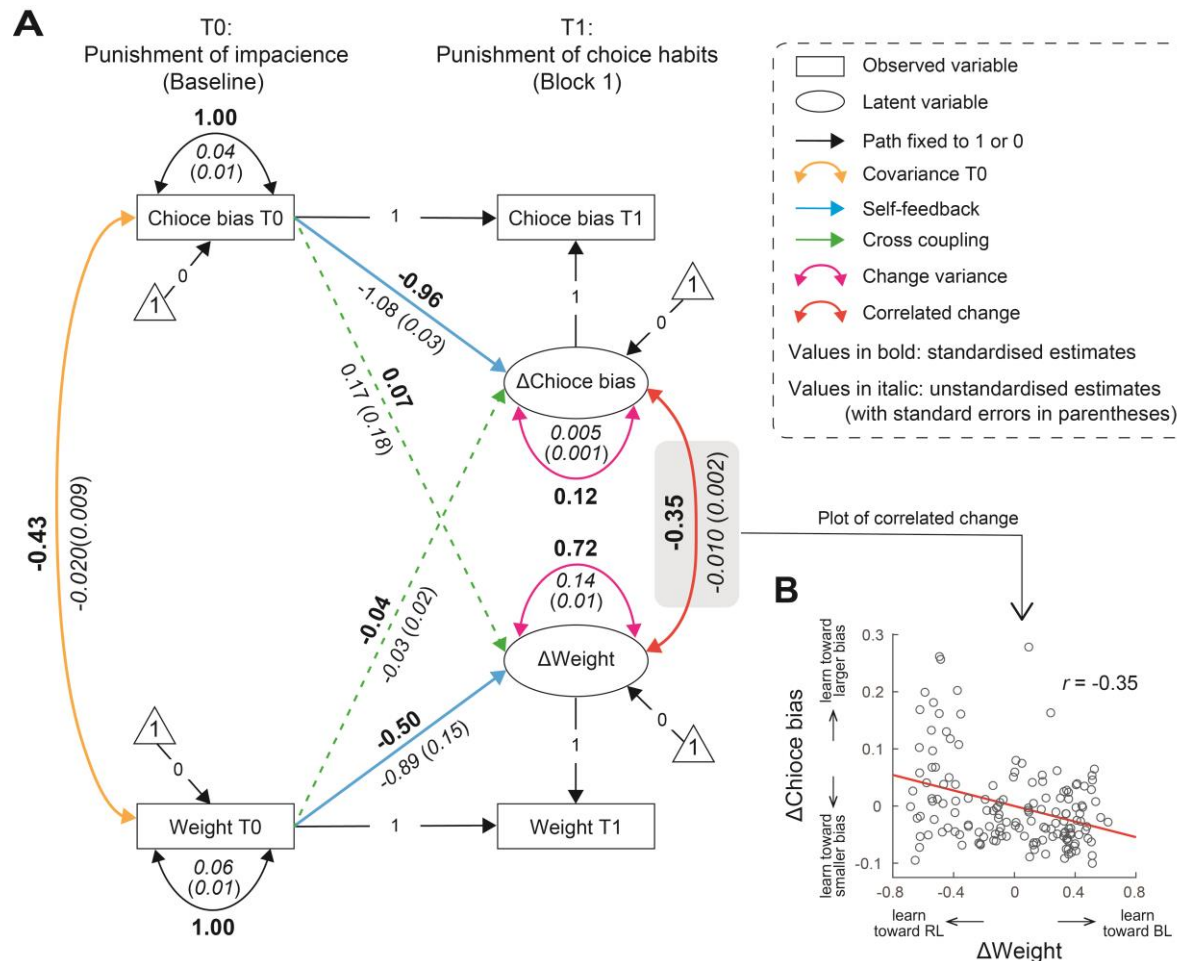
421         We investigated this inter-relationship by fitting the bivariate LCS model. Our observed
422    variables were the choice bias score and the belief learning weight estimated under the hybrid
423    learning rule. To estimate the latent change factors, we used these scores on the punishment
424    of impatience in the baseline and the punishment of standard choice habits in block 1, because
425    the participants were asked to change their choice habits between these two conditions.
426    Figure 6A illustrates fitted paths (significant paths are shown as thicker lines) from the pre-
427    punishment phase (T0: punishment of impatience) to the post-punishment phase (T1:
428    punishment of standard choice habits). This model, with fixed intercepts, shows a fit close to
429    the saturated (i.e., perfect) model ($\chi^2(4)$ = 0.00; RMSEA < 0.001, 90% confidence interval =
430    [0.000, 0.000]; CFI = 1.00; SRMR = 0.00; AIC = -325.1). We found auto-regressive effects in
431    both the choice bias score and BL weight score possibly because of a scale attenuation:
432    greater reductions in the choice bias were found in individuals who started off with a large bias,
433    and greater gains in the BL weight were found in individuals who started off with a low weight.
434    There were no significant cross-coupling effects: greater reductions in the choice bias were not
435    associated with individuals who attempted to mind-read the competitor's strategy before
436    punishment.

437         However, we found a medium negative correlated change between the change in the
438    choice bias and the change in the BL weight. This suggests that individuals who gained a
439    greater shift from reinforcement learning to belief learning exhibited a greater adaptation in
440    choice habits (Fig. 6B; standardised coefficient = -0.347, 95% confidence interval = [-0.480, -
441    0.199]). When we removed the cross-coupling paths from the model in Figure 6A, the model
442    still fitted the data well ($\chi^2(6)$ = 2.37; RMSEA < 0.001, 90% CI = [0.000, 0.053]; CFI = 1.00;
443    SRMR = 0.03; AIC = -326.8). However, removing the cross-coupling paths as well as the path
444    for the correlated change produced a bad model fit ($\chi^2(7)$ = 23.80; RMSEA = 0.126, 90% CI =
445    [0.071, 0.185]; CFI = 0.68; SRMR = 0.10; AIC = -309.0). To summarise, the computational
446    model suggests that people start forming their beliefs about when the competitor is going to act
447    and when they should act, after the competitor start responding to their own actions and
448    punishing their own choice habits. Together with structural equation modelling, the results
449    suggest that forming beliefs about the competitor's strategy helps to enhance adaptive
450    autonomy in avoiding choice habits.

451   We validated this finding in the following analysis. First, we used the choice bias score

452 and the difference in AICcs between the RL only model and BL only model as our observed

453 variables. We used these scores on the punishment of impatience in the baseline and the

454 punishment of standard choice habits in block 1. We still found that the reduction in the choice

455 bias co-occurs with a shift in strategies toward belief learning (Suppl. Fig. 3A). This result

456 validates that the correlated change in Figure 6 is unlikely to be due to the boundaries of the

457 parameter estimates under the hybrid learning rule.

458   Second, to validate whether the competitor's action indeed affected the participant's

459 action, we randomised the competitor's action in the sample we fitted the hybrid rule. This

460 permutation operation indeed disrupted the estimates of the BL weight (Suppl. Fig. 4A). In the

461 permutation sample, we did not find evidence that the reduction in the choice bias co-occurs

462 with a shift in strategies toward belief learning (Suppl. Fig. 4B&C). This result validates that the

463 competitor's action likely affected the participant's action, and that the shifting to the BL

464 strategy was likely associated with changing a pattern of choice.

465   Lastly, we checked the possibility that the reduction in the sequential bias or the

466 reinforcement bias co-occurs a shift in strategies. We did not find evidence that the reduction in

467 the sequential bias co-occurs with a shift in strategies toward belief learning, when the

468 transition habits were punished in block 2 (Suppl. Fig. 3B). Nor did we find that the reduction in

469 the reinforcement bias co-occurs with a shift in strategies toward belief learning, when the

470 reinforcement habits were punished in block 3 (Suppl. Fig. 3C). These follow-up analyses

471 confirm that our finding is not simply explained by the fact that the participant's data was used

472 to estimate both a measure of decision bias and parameters of the hybrid model.

473

474

**Figure 6.** The bivariate latent change score model of adaptive autonomy and strategic learning. **A**. Estimated parameters. The change score in the choice bias was modelled as a latent factor between the score before the punishment of choice habits and under the punishment. The change score in the belief learning weight recovered from the hybrid leaning rule was modelled similarly. Values in bold are standardised parameter estimates, and values in italic are un-standardised parameter estimates (with standard errors in parentheses). Solid lines indicate that the parameter is significant at $p < .05$. **B**. Scatter plot of correlated change. A greater adaptation in avoiding standard choice habits is correlated with a greater shift in the learning process from reinforcement learning to belief learning.

# Discussion

486

487 The human capacity for autonomous behaviour is widely asserted, and is fundamental to many
488 modern societies, but its cognitive basis is not well understood. We developed an experimental
489 paradigm that implicitly elicits autonomous behaviour, and we measured the extent to which
490 people could express autonomy by adapting their behaviour to free themselves from different
491 types of habits. We found that people can become autonomous of standard choice habits and
492 transition habits, but are limited in their ability to become free of reinforcement habits. We
493 further showed, in a large sample, that adaptive autonomy with respect to choice habits is
494 uncorrelated with adaptive autonomy with respect to transition habits and that to reinforcement
495 habits. This suggests distinct cognitive modules for these three forms of autonomy, rather than
496 a common module or a single form.  We further demonstrated the link between adaptive
497 autonomy and understanding the state of the environment: building beliefs about the
498 environmental, in our case, about a competitor's strategy, can enhance adaptive autonomy.

499

## Measures of autonomy

500

501 Traditional experimental psychology struggles to investigate autonomy because
502 traditional experiments in the studies of volition involve instructing people what they should do
503 (Baddeley, 1966; Brass & Haggard, 2007; Fleming et al., 2009; Jahanshahi et al., 1995; Libet
504 et al., 1983). The few studies that have examined human autonomous behaviour typically
505 involve competitive contexts (Forder & Dyson, 2016; Wang et al., 2014; Wong et al., 2021),
506 and have not considered subtypes of autonomy. We conceptualized three forms of autonomy,
507 as freedom from three cognitively distinct types of habit: *standard choice habits*, *transition*
508 *habits* and *reinforcement habits*. We attempted to evoke autonomous behaviour of each of
509 these three kinds using a common game-like context, and varying the competitor's strategy to
510 punish a lack of each type of autonomy. Using a statistical distance measure derived from
511 information geometry, we developed quantities that approximate a person's trait for each habit.
512 We quantified the extent to which people could break a specific habit when punished. A
513 covariance matrix underlying three adaptive autonomy measures showed no strong
514 correlations. This supported the idea that people express three distinct forms of behavioural
515 autonomy. In tasks where participants react to external stimuli quickly, it has been suggested
516 that a domain-general top-down control is used to solve different cognitive tasks (Braver, 2012;
517 Braver et al., 2007; Tang et al., 2022). However, in free, stimulus-independent action, our
518 results suggest that a domain-specific top-down control independently regulates each

519    particular form of autonomy: there are multiple ways to act freely, and it is important to consider
520    *from what* an agent is free. We studied choice biases, sequential biases and reinforcement
521    biases, but other biases to free action doubtless also exist. We showed that, for example, an
522    agent who becomes increasingly free from choice bias may be unable to free themselves from
523    the biasing effects of reinforcement.

524

525    **Relevance to classical neuropsychological tasks**

526    Our task takes a neuropsychological perspective on volitional behaviour and evokes
527    phenomena that neuropsychologists have traditionally studied using arbitrary, open-choice
528    tasks. For instance, our measure of choice bias is related to the capacity to inhibit a prepotent,
529    impulsive action (Mischel et al., 1972). People usually place costs on waiting, preferring earlier
530    rewards; a form of temporal discounting (Story et al., 2014). The sequential bias we measured
531    reflects executive control and working memory which are assessed using random number
532    generation tasks (Baddeley et al., 1998; Jahanshahi et al., 2000). In these tasks, people
533    cannot sufficiently randomise numbers and tend to seek simple rules such as repeating a digit
534    (e.g., 1,1,1), counting a digit in a natural sequence (e.g., 1,2,3) or larger with smaller inter-digit
535    gaps (e.g., 1,9,2) (Baddeley, 1966; Bar-Hillel & Wagenaar, 1991). In contrast, the capacity to
536    avoid the reinforcement bias is related to voluntary regulation of reward-seeking behaviour
537    (Bechara et al., 1994; Lejuez et al., 2002) and a balance between exploitation and exploration
538    (Cohen et al., 2007). These classical preferences are typically robust after repeating the task
539    (Neiman & Loewenstein, 2011; Ota et al., 2016) or after explicitly knowing one's own
540    behavioural trait (Ota et al., 2019). We used a competitive game to set environmental
541    constraints so that people should avoid such preferences in order to get rewards. Our results
542    demonstrate that people can balance choice frequencies and break transitions between
543    actions. Although we cannot directly compare our results with those of neuropsychological
544    studies, our results do suggest that, when pressurised, people can behave randomly and
545    autonomously more than suggested by the traditional neuropsychological literature.

546    In contrast, we found people could not avoid reinforcement habits. Neither positive
547    reinforcement bias nor negative reinforcement bias changed when penalised. A stereotypical
548    win-stay lose-shift behaviour has been shown in competitive games (Ota et al., 2020; Wang et
549    al., 2014). In particular, people are less flexible in changing lose-shift behaviour than win-stay
550    behaviour when adopting to new game rules (Forder & Dyson, 2016; Sundvall & Dyson, 2022).
551    The experience of a negative outcome generates a quicker decision and an impulsive

552 response on the next trial (Dyson et al., 2018). Indeed, event-related brain potentials show a
553 greater negative amplitude after negative outcome feedback than after positive outcome
554 feedback (Gehring & Willoughby, 2002; Hajcak et al., 2006). We note that individuals vary in
555 their ability to avoid habitual behaviour after negative reinforcement. These individual
556 differences are associated with post-error reaction times. Individuals who make quicker
557 decisions after a loss than after a win show a poorer performance than individuals who make
558 slower responses after a loss (Dyson, 2021). Therefore, overcoming impulsivity after a loss
559 may be a key aspect of volitional control for humans.

560

561 **Sustaining autonomy during interactions with the world**

562 Our model comparisons showed that participants did not achieve adaptive autonomy
563 simply by behaving randomly and stochastically. Rather, their strategy depended on the reward
564 assigned and the choices of the competitor. Reinforcement learning updates the best action by
565 a reward prediction error while the belief learning updates the best action by an action
566 prediction error, defined as a difference between the expected action competitor would take
567 and the actual action taken. Both reinforcement learning and belief learning can contribute to
568 volitional, self-generated actions, because both are determined by internal representations of
569 expected values, rather by an immediate stimulus (Frith, 2013). Furthermore, belief learning is
570 considered to recruit a mentalizing or an implicit understanding of what the other agent would
571 do (Amodio & Frith, 2006; Hampton et al., 2008). These neural substrates are often found in
572 separate neural networks: the reward prediction errors are encoded in the ventral striatum
573 (McClure et al., 2003; O'Doherty et al., 2004; Zhu et al., 2012) while the prediction errors about
574 the state of the environment are encoded in several areas including the rostral anterior
575 cingulate, the medial prefrontal cortex and the posterior superior temporal sulcus (Hampton et
576 al., 2008; Zhu et al., 2012).

577 Mentalising is a key cognitive component recruited in competitive games (Hampton et
578 al., 2008; Zhu et al., 2012). We found that people switch their strategy from reinforcement
579 learning to belief learning when the competitor started predicting their upcoming action. This
580 result suggests that people implicitly learned the likelihoods of the competitor's actions based
581 on a sampling of their past actions. Critically, belief learning was associated with enhanced
582 autonomy. In particular, successful adaptation in avoiding standard choice habits was
583 associated with a shift in learning: those individuals who shifted to learn from the likelihoods of
584 the competitor's actions rather than from reinforcement alone were able to gain greater

585   adaptive autonomy. We cannot tell whether participants discovered and explicitly represented

586   the punishment rules used by the competitor, but we speculate that explicitly understanding the

587   constraints on behaviour might be associated with increased autonomy.

588         To conclude, we have developed a new experimental paradigm and analysis pipeline

589   to study when and how human actions can become autonomous. We propose a new

590   theoretical construct of *adaptive autonomy*, meaning the capacity to free one's behavioural

591   choices from constraints of habitual responding, when a particular habit becomes dysfunctional,

592   for example due to environmental changes like the competitive pressure in our game scenarios.

593   We have shown that people can indeed express adaptive autonomy, and that they do so by

594   reducing habits of choice, habits of rule-based sequential action and habits of being guided by

595   reinforcement. These appear to be three distinct forms of adaptive autonomy, rather than a

596   single common strategy such as randomness. We show that becoming free from the effects of

597   reinforcement is particularly difficult. Finally, by showing that belief learning plays an important

598   part in boosting autonomy, we show a strong connection between autonomous action and

599   mentalising abilities.

600

601   **Limitations**

602   The three habit forms we tested were hierarchical. The standard choice habit – favouring one

603   action over all others – is more general, while the reinforcement habit is more specialised.

604   Therefore, we ordered the blocks so that the competitive game algorithms could penalise

605   habits progressively and serially. Thus, each block implicitly required participants to act more

606   freely and unpredictably than the preceding block. This fixed order may limit the generalisability

607   of our results, but the order we used is the most reasonable. The differences we observed

608   between the different forms of adaptive autonomy could be confirmed in further between-

609   participant studies.

610

611   **Empiricist view versus nativist views of human autonomy**

612   Our work is broadly compatible with an empiricist view of "free will" as opposed to a nativist

613   view. In our view, some of the key attributes historically associated with "free will", such as the

614   ability to act endogenously and purposefully, can be acquired, or at least adapted, through

615   experience.  This adaptation requires people to make novel, non-habitual, 'smart' actions in

616   certain situations. We found that people were more or less successful in adapting their trait

617   habits both at the individual level and at the level of different punishments. An individual's

23

618  degree of autonomy is unique and contingent on environmental constraints. A strongly nativist
619  view would suggest that autonomy is a state that occurs inside an individual's mind and is
620  independent of the external world. However, our results imply that being sensitive to the
621  contingencies of the external environment, and its restrictions on one's own actions is key to
622  autonomy. In this sense, autonomy can be seen a reasoned, goal-oriented response that
623  occurs within an environmental context.

624

625

# Methods
626

### Participants
627

628  One hundred and fifty-nine participants (age range = 18–45, M = 29.5 yo, SD = 7.2) were
629  recruited online via the Prolific website (https://www.prolific.co/). Participants received a basic
630  payment of £3.75 for their participation in a 30 minute experiment. They earned a bonus of up
631  to £4 based on their performance on the task. There were 95 female participants and 64 male
632  participants. Recruitment was restricted to the United Kingdom. Seven participants were
633  excluded from the analysis (Suppl. Info.) and the remaining 152 participants were analysed. All
634  procedures were approved by the Research Ethics Committee of University College London.
635  Participants gave informed consent by checking and validating the consent form.

636

### Experimental design
637

638  *Apparatus.* We used the JavaScript library jsPsych (de Leeuw, 2015) and the plugin jsPsych-
639  psychophysics (Kuroki, 2021) to program the task and hosted the experiment on the online
640  research platform Gorilla (https://gorilla.sc/) (Anwyl-Irvine et al., 2020), which participants
641  could access through their browser on their own computer. We assumed that monitor sampling
642  rates were typically around 60 Hz, with little variation across computers (Anwyl-Irvine et al.,
643  2020). The size and position of stimuli were scaled based on each participant's screen size
644  which was automatically detected. The size of stimuli reported below are for a monitor size of
645  15.6" (view point size, width x height: 1536 x 746 pixels).

646

647  *Stimuli and task.* Each trial started with a fixation cross, which appeared for 0.6–0.8 seconds.
648  The images of a tree, a flock of birds and a basket containing apples then appeared (Fig. 2A).
649  A tree (width x height: 307 x 375 pixels) was shown on the left of the screen and a flock of
650  birds (width x height: 30 x 22 pixels each) were located on the tree. A rectangular basket of

651 apples (width x height: 153 x 167 pixels) was presented in the bottom centre. After the fixation

652 cross disappeared and all images appeared, the participants were given 4.5 sec to throw the

653 food. Pressing a key initiated delivery of the food to a storage location which was located at

654 447 pixels forward from the start point. This delivery took 1.5 sec. We programmed the birds to

655 attempt to intercept and catch the food. The birds on each trial were designed to intercept the

656 food thrown within one of three intervals: 1) early throw (0–1.5 sec), 2) middle throw (1.5–3.0

657 sec) or 3) late throw (3.0–4.5 sec). After their departure, it took approximately 0.25 sec for each

658 bird to reach the storage location; the birds passed through that point. The participants

659 competed with the virtual competitor, aiming to deliver food before or after the birds reached

660 the storage location. We counted whether one of the birds overlapped with the food when the

661 delivery was completed (at the offset of moving). If this was the case, the food was caught, and

662 the participant lost a trial. If not, the food was delivered without it being caught, and the

663 participant won a trial. If no response was submitted before 4.5 sec, the food was launched

664 automatically, and a trial was terminated as a timeout. Finally, we provided a feedback

665 message: "Success!", "Fail!" or "Timeout!", which lasted for 1.0 sec. The next trial then started

666 with a fixation cross.

667 In the instructions, we emphasised the following points. First, merely reacting to the

668 absence of a stimulus – the birds resting in the tree – will not win the game because the birds

669 can travel much faster than the food. Second, merely waiting for the birds to pass is not a

670 solution because of the time constraint. Third, the birds' flying interval is not the same on every

671 trial, nor is it random. Instead, the birds can learn when the participant is likely to throw the

672 food. Therefore, it is important to predict when the birds will likely fly and to randomise your

673 throw times in order to avoid the competitor's prediction.

674

675 *Procedure.* Participants first received the instructions and viewed a set of demonstrations

676 about the task. Following some practice trials, the participants completed four blocks of the

677 game with a 1-minute break between blocks. The baseline block lasted 2.5 minutes while the

678 remaining blocks 1, 2 and 3 lasted 5 minutes each. The participants got as many throws of the

679 food as they could in the 2.5 or 5 minutes. The participants could check how much time was

680 left in each block. We used time, and not trial number, to terminate each block so that

681 participants did not respond immediately on every trial, finishing the game early. The bonus

682 payment was determined by the percentage of throws that successfully avoid birds and was

683 paid up to £1 for each block: if 40 out of 60 throws are successful, we paid £1 x 66.6% = £0.66

25

684 (average bonus, baseline: £0.94; block 1: £0.63; block 2: £0.59; block 3: £0.57). The success
685 rate and the timeout rate were included in the feedback. Nevertheless, we assumed that some
686 participants might consume time by not focusing on the game. To prevent this, we encouraged
687 participants to sustain the proportion of timeout trials under 5%.
688

689 **Competitor design**
690 The virtual competitor design was primarily inspired by primate work (Barraclough et al., 2004;
691 Lee et al., 2004) and by rat work (Tervo et al., 2014). We programmed the learning algorithm
692 (i.e., birds) to seek out behavioural patterns in the participant's choice history and to pressure
693 participants into novel behaviour. The participants could decide the time to act between 0 sec
694 (as soon as birds and food appeared) and 4.5 sec (until timeout). To make the competitor's
695 prediction simple, we discretised the time window into three intervals, 1) early interval (0–1.5
696 sec), middle interval (1.5–3.0 sec) and late interval (3.0–4.5 sec). Given past behaviour, the
697 competitor predicted which response interval a participant was likely to select. Accordingly, the
698 two other intervals were primed for winning: if the participant threw the food during the interval
699 predicted by the competitor, the participant lost. If the participant threw the food during one of
700 two other intervals, the participant won. We adjusted the birds' departure times by taking their
701 travel time (0.25 sec) and the food delivery/travel time (1.5 sec) into account: if a prediction
702 was made on the late interval, the birds departed from the tree at the period of 4.25–5.75 sec,
703 and they reached the delivery point during 4.5–6.0 sec to catch the food when it was delivered.
704 We designed four distinct competitors (Fig. 2B). First, in the baseline block, Competitor
705 0 punished participants for being impatient. In this block, the birds blocked the early throw on
706 every trial. Thus, the stimulus-absence behaviour corresponded with waiting until the middle
707 interval. Competitor 0 measured the volitional control to resist immediacy or external triggers
708 (Haggard, 2019). Second, in block 1, Competitor 1 predicted *standard choice habits* (choice
709 preferences) – which interval the participant is going to select –. On each trial, a history of the
710 participant's past ten choices was used to estimate the probabilities of selecting the early,
711 middle and late interval. The choice probabilities were then used to generate the competitor's
712 prediction on the upcoming choice. For instance, if the participant chose the early interval
713 seven times, the middle interval twice, and the late interval once, the competitor penalised the
714 early interval 70% of the time, the middle 20% of the time and the late 10% of the time. Thus,
715 Competitor 1 required participants to balance their general choice frequencies.

716     In block 2, Competitor 2 sought out *transition habits* (sequential patterns) – which

717     interval the participant is going to select after the participant made a particular response –. A

718     history of the past 60 trials was used to estimate the conditional probabilities of selecting three

719     intervals given the previous reaction time. The estimated probabilities were conditioned on the

720     last reaction time ± 0.5 sec. Suppose a participant took 2.5 sec to act in the previous trial.

721     Competitor 2 might discover that, in the past, the participant chose the early interval twice, the

722     middle interval twice, and the late interval six times after the participant had acted in 2.0-3.0

723     sec. In this case, Competitor 2 penalised the late interval 60% of the time. We assumed that

724     using the previous response time (i.e., continuous variable) is more powerful to predict the next

725     response than using the previous response interval (i.e., categorical variable). Competitor 2

726     pressured participants in avoiding habitual transition patterns. Finally, in block 3, Competitor 3

727     punished *reinforcement habits* (outcome dependence) – which interval the participant is going

728     to select after the participant made a particular response and won a trial or lost a trial –.

729     Competitor 3 used the same search algorithm as Competitor 2 with the exception that they

730     conditioned the search on the last reaction time and the last outcome. Competitor 3 required

731     participants to act independently from the previous outcome.

732

733     **Data analysis**

734     Because the birds intercepted one of the three response intervals, we mainly analysed the

735     data that was discretised into 1) the early response: responding in 0–1.5 sec, 2) the middle

736     response: responding in 1.5–3.0 sec, 3) the late response: responding in 3.0–4.5 sec (including

737     timeout).

738

739     *Quantifying trait habits.* Statistical distance is a standardised way of measuring the extent to

740     which the observed probability distribution is different from the target probability distribution.

741     We calculated the Kullback-Leibler divergence to quantify the extent to which the participant's

742     choice probability distribution is different from the choice probability distribution that a habit-free

743     agent would exhibit, a proxy of three habit families. See Figure 3.

744

745     1) Choice bias. Competitor 1 punished standard choice habits in selecting one interval more

746     often than the other two. The probabilities of choosing the early, middle and late interval for a

747     habit-free agent would be 0.33, respectively. We computed the choice probabilities $P(c)$ given

748     a history of intervals each participant chose in each block. The K-L divergence is then

27

749
$$D_{KL\ choice\ bias} = \sum_{c \in E,M,L} P(c) \log_2 \left( \frac{P(c)}{0.33} \right)$$

750

751 2) Sequential bias. Competitor 2 punished transition habits on the top of choice habits. Similar

752 to computing the choice probabilities, we computed the conditional probabilities of choosing the

753 early, middle and late interval given the interval chosen on the previous trial $P(c|c_{-1})$. We

754 measured the K-L divergence of these participant's conditional probabilities from the

755 participant's choice probabilities. The K-L divergence for each previous interval $c_{-1}$ is

756 computed as

757
$$D_{KL\ c_{-1}} = \sum_{c \in E,M,L} P(c|c_{-1}) \log_2 \left( \frac{P(c|c_{-1})}{P(c)} \right).$$

758 The total K-L divergence as a weighted sum is then

759
$$D_{KL\ sequential\ bias} = \sum_{c_{-1} \in E,M,L} P(c_{-1}) \cdot D_{KL\ c_{-1}}$$

760 Since we conditioned the K-L divergence on the previous interval chosen, we took the

761 proportion of observing that situation into account, and we weighed each divergence by this

762 prior probability. The target probabilities (i.e., habit-free agent) were set to be the participant's

763 own choice probabilities, rather than purely stochastic choices 0.33. Therefore, $P(c|c_{-1})$

764 becomes equivalent to $P(c)$ and the K-L divergence becomes zero, as long as the participant

765 selects three intervals independently from the previous choice (even if the participant favours

766 one interval). By this way, we quantified the deviation of patterns associated with the previous

767 choice from sequential patterns logically expected from the participant's own choice

768 probabilities. Competitor 2 specifically detected and punished this conditional dependence.

769

770 3) Reinforcement bias. Competitor 3 punished reinforcement habits on the top of choice habits

771 and transition habits. Similar to computing the choice probabilities, we computed the

772 conditional probabilities of choosing the early, middle and late interval given the interval chosen

773 and the outcome obtained on the previous trial $P(c|c_{-1}, o_{-1})$. We measured the K-L divergence

774 of these participant's conditional probabilities from the participant's conditional probabilities

775 given the previous interval solely. The K-L divergence for each previous interval $c_{-1}$ and each

776 previous outcome $o_{-1}$ is computed as

777
$$D_{KL\ c_{-1}, o_{-1}} = \sum_{c \in E,M,L} P(c|c_{-1}, o_{-1}) \log_2 \left( \frac{P(c|c_{-1}, o_{-1})}{P(c|c_{-1})} \right).$$

778    The total K-L divergence as a weighted sum is then

779

$$D_{KL\ reinforcement} = \sum_{c_{-1}\ \in\ E,M,L}\ \sum_{o_{-1}\ \in\ \substack{success \\ fail}} P(c_{-1}, o_{-1}) \cdot D_{KL\ c_{-1},o_{-1}}$$

780    Since we conditioned the K-L divergence on the previous interval chosen and the previous

781    outcome obtained, we took the proportion of observing that situation into account, and we

782    weighed each divergence by this joint prior probability. The target probabilities (i.e., habit-free

783    agent) were set to be the participant's own conditional probabilities given the previous interval

784    solely. Therefore, $P(c|c_{-1}, o_{-1})$ becomes equivalent to $P(c|c_{-1})$ and the K-L divergence

785    becomes zero, as long as the participant selects three intervals independently from the

786    previous outcome (even if the participant's choice depends on the previous interval). By this

787    way, we quantified the deviation of patterns associated with both the previous choice and the

788    previous outcome from patterns logically expected from the conditional dependence on the

789    previous choice solely. Competitor 3 specifically detected and punished this outcome

790    dependence. We also quantified the positive reinforcement bias and the negative

791    reinforcement bias, separately (Suppl. Fig. 1). We computed the K-L divergence of the

792    conditional probabilities given the previous interval and the previous win only or the previous

793    loss only from purely stochastic choices 0.33. Here the statistical distance can be argued as

794    the distance between the participant's post-win behaviour or post-loss behaviour and the habit-

795    free agent who is purely random. These measures were used to generate Figure 3.

796

797    *Statistical analysis.* We tested the performance difference by Wilcoxon signed rank test. The

798    alpha level of 0.05 was corrected by the number of tests we performed in each class of test

799    (Bonferroni correction).

800

801    **Computational models**

802    *Reinforcement learning.* We tested a reinforcement learning (RL) model in which an action

803    value is updated via a Rescorla-Wagner rule (Sutton & Barto, 2018). On each trial, an RL

804    agent selects an action from the early, middle or late interval $a \in E, M, L$. For an action $a$

805    selected on a trial $t$, the value of action $a$ is updated by a prediction error $\delta$:

806                                     $$\delta_t = r_t - V_t(a)$$

807    where $r_t$ is the actual reward received (1 for successfully avoiding birds and 0 for failure) and

808    $V_t(a)$ is the current expected reward for that action. The reward prediction error $\delta_t$ is then used

809    to update the value of the selected action, weighted by the learning rate $\alpha$

29

810
$$V_{t+1}(a) = V_t(a) + \alpha\delta_t.$$

811

812 *Belief learning.* In a belief learning (BL) model, a BL agent infers the opponent's state of mind –

813 what option the opponent is going to select – and decides on the action that maximises the

814 expected reward (Camerer, 2003; Hampton et al., 2008; Zhu et al., 2012). Actions $a' \in E, M, L$

815 are available for the competitor to choose. For each action $a \in E, M, L$ on trial $t$, the value of

816 that action is updated by a prediction error

817
$$\delta_t = r_t - V_t(a)$$

818 where $r_t = -1$ if $a'$ is same as $a$ (i.e., the competitor selects the same response interval as the

819 participant) while $r_t = 0$ if $a'$ is different from $a$ (i.e., the competitor selects a different response

820 interval). This prediction error is the difference between the current expected value and the

821 negative reward derived from the competitor's current choice. Therefore, the updated expected

822 value of action reflects the likelihood of the competitor's choice: the larger the value, the less

823 likely the competitor choose. The same rule with the RL model was used to update the value of

824 action $a$, with the exception that the values of all three intervals were updated on every trial.

825 Suppose that the birds repeatedly selected the early interval to intercept the food. A direct

826 observation of the birds' flight at the early interval decreases the value of the early interval. At

827 the same time, this observation implies that the birds did not fly at the middle nor the late

828 interval. This increases the values of these intervals.

829

830 *Hybrid learning.* We modelled the hybrid learning rule (aka. experience weighed attraction) as

831 a combination of reinforcement learning and belief learning (Camerer & Ho, 1999; Hampton et

832 al., 2008; Zhu et al., 2012). After updating the value of the action in each learning process, the

833 hybrid rule combines the values of the action such that

834
$$V_{t+1}(a) = (1 - w) \cdot V_{t+1}^{RL}(a) + w \cdot V_{t+1}^{BL}(a)$$

835 where one additional free parameter $w$ is used to weigh the relative contribution placed on

836 belief learning over reinforcement learning.

837 For all models, the action values were converted into the choice probabilities using the

838 soft-max function to simulate action selection,

839
$$P_t(a) = \frac{e^{\beta \cdot (V_t(a) + b(a))}}{\sum_{a \in E, M, L} e^{\beta \cdot (V_t(a) + b(a))}}$$

840 where $P_t(a)$ is the probability of choosing the interval $a$. The inverse temperature parameter $\beta$

841 scales the relative difference between the choice probabilities, which scales decision

842 uncertainty. We added the decision preference term $b$ with an exponential temporal
843 discounting (Story et al., 2014):

$$b(a) = e^{-\rho \cdot T(a)}$$

845 where $T$ is the time corresponding to the chosen interval ($T = 0$, 1.5 or 3.0 sec for the early,
846 middle or late interval, respectively). The parameter $\rho$ scales the relative preference to earlier
847 intervals, which captures an individual's temporal discounting or impatience to wait.

848 For each model, we fitted the model decision probabilities to the participant's choice
849 interval data by minimizing the negative log-likelihood of the observed choices using Bayesian
850 adaptive direct search (BADS) (Acerbi & Ma, 2017). Free parameters were optimised
851 individually for each participant and separately for each block with the following boundaries:
852 $\alpha \in [0,1]$, $\beta \in [0,20]$, $\rho \in [0,0.2]$, $w \in [0,1]$. The parameter $w$ was fixed to 0 for the
853 reinforcement learning model and fixed to 1 for the belief learning model. For the stochastic
854 selection model, we fixed the parameters $\alpha$ and $w$ to 0 so that this model could only capture
855 the participants' decision uncertainty and their choice preference. This model produced
856 constant model decision probabilities across all trials. To verify that we had found the global
857 minimum, we repeated the search process with different starting points. For model comparison,
858 we applied AICc—Akaike information criterion with a correction for finite sample size—to each
859 participant and model as the information criterion for goodness-of-fit (Burnham & Anderson,
860 1998; Hurvich & Tsai, 1989). The summed AICc across participants was reported in Figure 5B.
861

862 **Simulated play**
863 During simulated play, an RL agent and a BL agent played against the prediction algorithm
864 used by each class of competitor. The competitor's prediction was made using a history of
865 choices a simulated agent made rather than using real data. For each simulated play of 60
866 trials (which is approximately equal to the number of trials in the real game), the success rate
867 was computed. We simulated each agent's behaviour given a set of model parameters. Each
868 set of parameters was determined by an extensive grid search in the parameters' space. The
869 simulation play was repeated 3,000 times for each parameter set.
870

871 **Latent change score model**
872 Latent change score (LCS) models are the statistical framework that captures the process
873 underlying the change in the variables of interest at two measurement occasions (Carpenter et

874    al., 2019; Kievit et al., 2018; Kievit et al., 2017; McArdle, 2009). LCS models conceptualise the

875    score of variable $X$ at time point $T2$ as

876    $$X_{T2} = \beta X_{T1} + \Delta X$$

877    where the score $X_{T2}$ is a function of the score $X_{T1}$ weighted by an auto-regressive (i.e., self-

878    feedback) parameter $\beta$ and some residual $\Delta X$. By fixing the regression weight of $X_{T2}$ on $X_{T1}$ to

879    1, the change score $\Delta X$ can be simply rewritten as

880    $$\Delta X = X_{T2} - X_{T1}$$

881    In structural equation modelling, the change score can be defined as a latent factor by fixing a

882    factor loading on the score $X_{T2}$ to 1. By this mathematical manipulation, the change between

883    $T1$ and $T2$ is modelled as a latent factor. Bivariate LCS models predict the change score by an

884    auto-regressive parameter $\beta$ and a cross-coupling parameter $\alpha$:

885    $$\Delta X = \beta X_{T1} + \alpha Y_{T1}$$

886    $$\Delta Y = \beta Y_{T1} + \alpha X_{T1}$$

887    In this equation, the auto-regressive parameter $\beta$ captures the degree to which the initial score

888    $X_{T1}$ predicts (or is proportional to) the change score $\Delta X$. The cross-coupling parameter $\alpha$

889    captures the degree to which the initial score in another domain $Y_{T1}$ predicts (or is proportional

890    to) the change score $\Delta X$. Above these effects, the bivariate LCS models quantifies the

891    variance-covariance structure in the change factor, which estimates the correlated change: the

892    degree to which the change score in one domain $\Delta X$ covaries with the change score in another

893    domain $\Delta Y$ after taking auto-regressive and cross-coupling effects into account. See Kievit et

894    al., 2018 & McArdle, 2009 for reviews and Kievit et al., 2017 & Carpenter et al., 2019 for its

895    applications.

896    We examined the inter-relationships between adaptive autonomy and a shift in learning

897    strategies. There were three decision bias scores (choice bias, sequential bias and

898    reinforcement bias) and four hybrid model parameters (learning rate, decision uncertainty,

899    decision preference and relative weight of belief learning). Considering potential correlations

900    among seven variables, we controlled for the influences of the other five variables on bivariate

901    changes. In Figure 6, we used the choice bias score and the belief learning weight. We first

902    regressed all other five variables against these two scores and retaining only the residuals from

903    the regression. We then used the residual scores in the choice bias and the residual scores in

904    the belief learning weight at two measurement points to fit the bivariate LCS model. Therefore,

905    any parameter estimates in the path model were considered a mere relationship between two

906    variables included in the model. Because we fitted the bivariate LCS model to residuals,
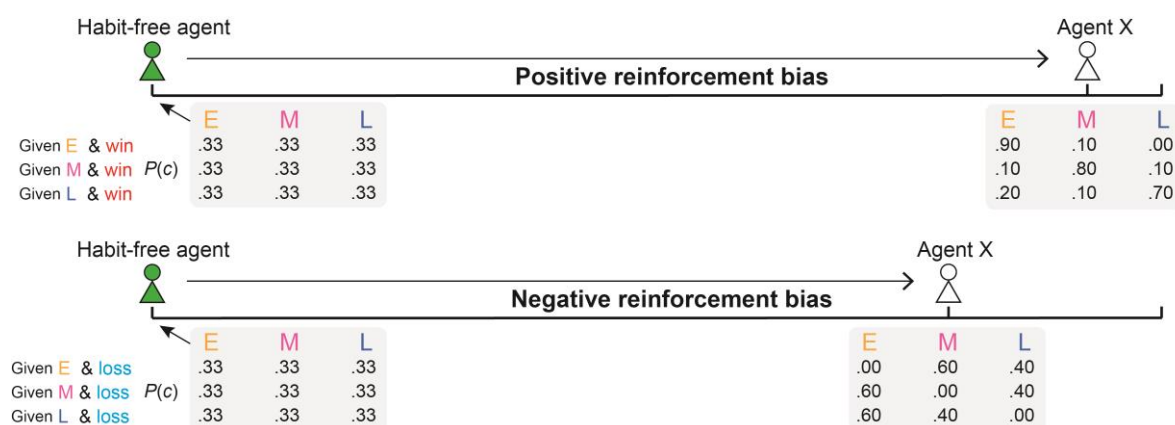
907 parameter estimates for the initial intercepts (i.e., mean initial scores) and change intercepts

908 (i.e., mean change scores) were fixed to zero.

909        Models were estimated in the lavaan package for R (version 0.6-11) (Rosseel, 2012).

910 We used maximum likelihood estimation with robust (Huber-White) standard errors and a

911 scaled test statistic. We evaluated overall model fit using the root-mean-square error of

912 approximation (RMSEA; acceptable fit: < 0.08; good fit < 0.05), the comparative fit index (CFI;

913 acceptable fit: 0.95 to 0.97; good fit > 0.97) and the standardized root-mean-square residual

914 (SRMR; acceptable fit: 0.05 to 0.10, good fit: < 0.05) (Schermelleh-Engel et al., 2003).

915

916

917 # Supplementary information
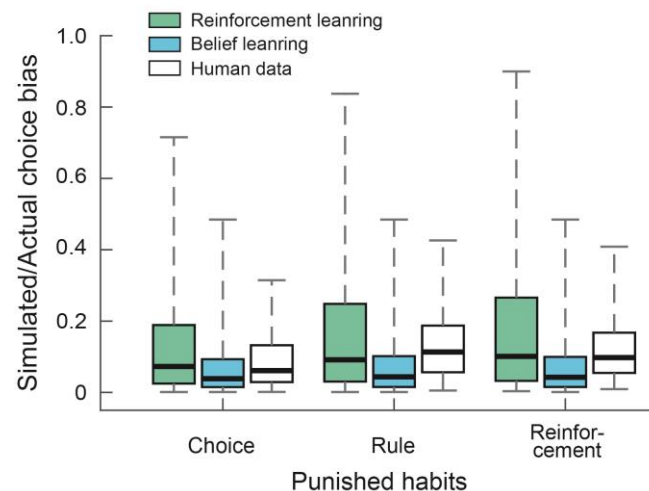
918 **Exclusion criterion**

919 We encouraged participants to sustain the percentage of timeout trials under 5%. We checked

920 the histogram of the timeout rates. Seven participants displayed a timeout above 13 %. This

921 was considerably high compared with the other participants (0–5%: 119 participants; 5–8%: 27

922 participants; 8–11%: 6 participants; 13–20%: 5 participants; >20%: 2 participants). These

923 participants might not be able to follow the instructions or might not be able to keep their

924 attention on the task, thereby removed from the analysis.

925

926

927



928
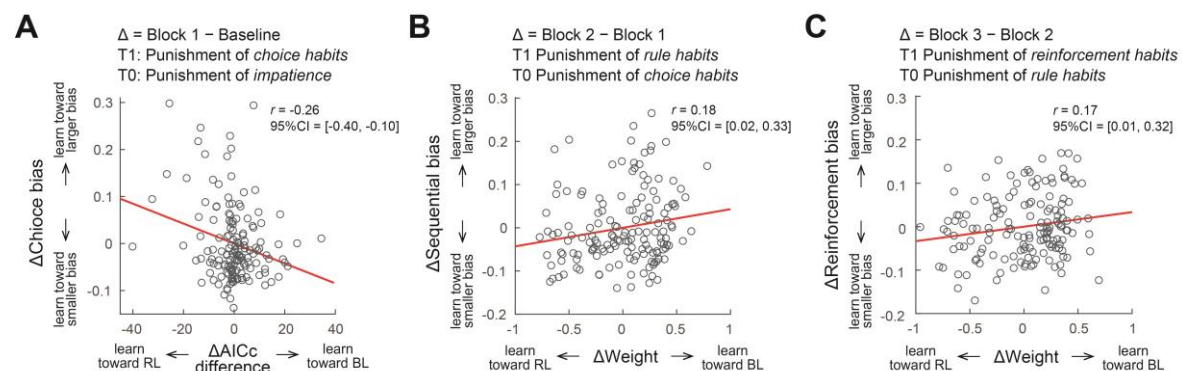
929 **Supplementary Figure 1**

930 The scores for positive reinforcement bias and negative reinforcement bias were computed

931 from the conditional probabilities of wins only and from the conditional probabilities of losses

932    only. These bias scores measure the statistical distance from random probabilities 0.33. See

933    Figure 3.

934

935



936

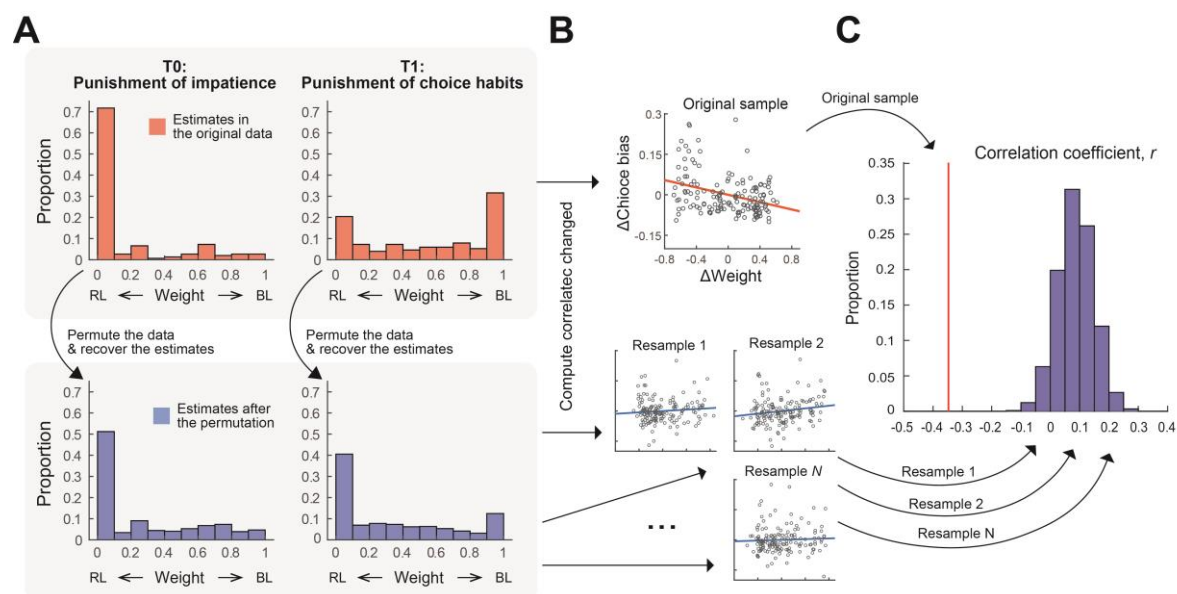**Supplementary Figure 2**

938    Real choice bias (white bars) in the actual experiment and fictive choice bias (green or blue

939    bars) in simulated play. We simulated the choice bias from agents using the reinforcement

940    learning (green) and agents using the belief learning (blue) (see *Simulated play*). Belief

941    learning produced a smaller choice bias in all punishment schemes.

942



943

**Supplementary Figure 3**

945    As supplementary results, we fitted the bivariate LCS model to the data set in panels **A-C**,

946    respectively. Estimated parameters in the path model were omitted for the sake of space.

947    Instead, we visualised the scatter plot of correlated change as we did in Figure 6B. **A**. Here we

948    used the choice bias score and the AICc difference between the RL only model and BL only

949     model, instead of the BL weight estimated under the hybrid learning rule. To estimate the latent

950     change factors, we used these observed scores on the punishment of impatience in the

951     baseline and the punishment of standard choice habits in block 1. We found a negative

952     correlated change (standardised coefficient = -0.256, 95% CI = [-0.398 -0.101]): shifting the

953     strategy toward belief learning as opposed to reinforcement learning was associated with

954     gaining greater reductions in the choice bias. **B**. Here we used the scores of sequential bias

955     and the scores of the BL weight on the punishment of choice habits in block 1 and the

956     punishment of rule habits in block 2, to estimate the latent change factors. We did not find a

957     negative correlated change (standardised coefficient = 0.180, 95% CI = [0.021 0.330]): shifting

958     the strategy toward belief learning was not associated with gaining greater reductions in the

959     sequential bias. **C**. Here we used the scores of reinforcement bias and the scores of the BL

960     weight on the punishment of rule habits in block 2 and the punishment of reinforcement habits

961     in block 3, to estimate the latent change factors. We did not find a negative correlated change

962     (standardised coefficient = 0.168, 95% CI = [0.009 0.319]): shifting the strategy toward belief

963     learning was not associated with gaining greater reductions in the reinforcement bias.

964

965



967 **Supplementary Figure 4**

968 A simulated experiment using permutation operation. **A**. In upper panels (red bars), we show

969 the proportion of the belief learning weight estimated under the hybrid learning rule. A weight

970 increases between blocks. In this model fitting, a trial sequence of rewards (success or failure)

and that of birds' flight interval (early, middle or late interval) were used to simulate the participant's choice interval (see *Computational models*). If the competitor's action is critical for the model fitting, permutating the competitor's action would disrupt the estimates of the BL weight. To validate this, we swapped the birds' flight interval. For instance, if the participant chose the late interval and the birds intercepted the middle interval, we swapped the birds' choice to the early interval on that trial. This permutation does not alter an outcome (i.e., a successful trial is still success) but does randomise the competitor's action. We did not swap the birds' choice for unsuccessful trials. In each iteration, we permutated the birds' flight interval 75% of the time and estimated the BL weight under the hybrid rule. We repeated this procedure 100 times for each participant and for each block. In lower panels (blue bars), the parameter estimates recovered from the permutation sample are shown. The BL weights are right-skewed and indeed differ from the original estimates. **B**. In the upper panel, we show a scatter plot of the correlated change: how reductions in the choice bias co-occurs with a shift in learning strategies, replotted from Figure 6B. The correlation coefficient in this original sample ($r$ = -0.347) is plotted as a vertical red line in the panel **C**. In lower panels, we plot the same correlated change but derived from the permutation sample. We only show three representative plots. In each sample, we computed the standardised coefficient of the correlated change. The proportion of these correlation coefficients is shown as a blue histogram in the panel **C**. The 95% confidence interval ranged from -0.03 to 0.20. The original correlation coefficient is significantly different from the permutated coefficients. This analysis validates the robustness of estimating the belief learning weight and the robustness of estimating the coefficient in the correlated change.

# Competing interests

The authors declare no competing interests.

# Author's contributions

Conceptualization, KO and PH; data collection, KO; investigation, KO, LC, and PH; formal analysis, KO; writing – original draft, KO; writing – review & editing, KO, LC and PH; supervision, LC and PH; funding acquisition, PH.

# Acknowledgments

1003 This research was supported by the John Templeton Foundation and the Fetzer

1004 Foundation awarded to PH.

1005

# **Data and code availability**

1007 All data and analysis codes are available here:

1008

# **References**

1010 Acerbi, L., & Ma, W. J. (2017). Practical Bayesian optimization for model fitting with Bayesian

1011 adaptive direct search. *Advances in neural information processing systems, 30*, 1834-1844.

1012

1013 Amodio, D. M., & Frith, C. D. (2006, Apr). Meeting of minds: the medial frontal cortex and social

1014 cognition. *Nat Rev Neurosci, 7*(4), 268-277. https://doi.org/10.1038/nrn1884

1015

1016 Anwyl-Irvine, A. L., Massonnie, J., Flitton, A., Kirkham, N., & Evershed, J. K. (2020, Feb). Gorilla in

1017 our midst: An online behavioral experiment builder. *Behav Res Methods, 52*(1), 388-407.

1018 https://doi.org/10.3758/s13428-019-01237-x

1019

1020 Baddeley, A., Emslie, H., Kolodny, J., & Duncan, J. (1998, Nov). Random generation and the

1021 executive control of working memory. *Q J Exp Psychol A, 51*(4), 819-852.

1022 https://doi.org/10.1080/713755788

1023

1024 Baddeley, A. D. (1966, May). The capacity for generating information by randomization. *Q J Exp*

1025 *Psychol, 18*(2), 119-129. https://doi.org/10.1080/14640746608400019

1026

1027 Bar-Hillel, M., & Wagenaar, W. A. (1991). The perception of randomness. *Advances in Applied*

1028 *Mathematics, 12*(4), 428-454. https://doi.org/10.1016/0196-8858(91)90029-i

1029

1030 Barandiaran, X. E., & Di Paolo, E. A. (2014). A genealogical map of the concept of habit. *Front Hum*

1031 *Neurosci, 8*, 522. https://doi.org/10.3389/fnhum.2014.00522

1032

1033 Barraclough, D. J., Conroy, M. L., & Lee, D. (2004, Apr). Prefrontal cortex and decision making in a

1034 mixed-strategy game. *Nat Neurosci, 7*(4), 404-410. https://doi.org/10.1038/nn1209

1035

Bechara, A., Damasio, A. R., Damasio, H., & Anderson, S. W. (1994). Insensitivity to future consequences following damage to human prefrontal cortex. *Cognition, 50*(1-3), 7-15. https://doi.org/10.1016/0010-0277(94)90018-3

Brass, M., & Haggard, P. (2007, Aug 22). To do or not to do: the neural signature of self-control. *J Neurosci, 27*(34), 9141-9145. https://doi.org/10.1523/JNEUROSCI.0924-07.2007

Braver, T. S. (2012, Feb). The variable nature of cognitive control: a dual mechanisms framework. *Trends Cogn Sci, 16*(2), 106-113. https://doi.org/10.1016/j.tics.2011.12.010

Braver, T. S., Gray, J. R., & Burgess, G. C. (2007). *Explaining the many varieties of working memory variation: Dual mechanisms of cognitive control.* . Oxford University Press. https://doi.org/https://doi.org/10.1093/acprof:oso/9780195168648.003

Burnham, K. P., & Anderson, D. R. (1998). Model Selection and Multimodel Inference. *Springer, New York, NY.*

Camerer, C., & Ho, T. (1999). Experience-weighted Attraction Learning in Normal Form Games. *Econometrica, 67*(4), 827-874. https://doi.org/10.1111/1468-0262.00054

Camerer, C. F. (2003). *Behavioral Game Theory: Experiments in Strategic Interaction*. Princeton Univ Press.

Carpenter, J., Sherman, M. T., Kievit, R. A., Seth, A. K., Lau, H., & Fleming, S. M. (2019, Jan). Domain-general enhancements of metacognitive ability through adaptive training. *J Exp Psychol Gen, 148*(1), 51-64. https://doi.org/10.1037/xge0000505

Cohen, J. D., McClure, S. M., & Yu, A. J. (2007, May 29). Should I stay or should I go? How the human brain manages the trade-off between exploitation and exploration. *Philos Trans R Soc Lond B Biol Sci, 362*(1481), 933-942. https://doi.org/10.1098/rstb.2007.2098

de Leeuw, J. R. (2015, Mar). jsPsych: a JavaScript library for creating behavioral experiments in a Web browser. *Behav Res Methods, 47*(1), 1-12. https://doi.org/10.3758/s13428-014-0458-y

1070 Dolan, R. J., & Dayan, P. (2013, Oct 16). Goals and habits in the brain. *Neuron, 80*(2), 312-325.
1071    https://doi.org/10.1016/j.neuron.2013.09.007

1072

1073 Du, Y., Krakauer, J. W., & Haith, A. M. (2022, May). The relationship between habits and motor
1074    skills in humans. *Trends Cogn Sci, 26*(5), 371-387. https://doi.org/10.1016/j.tics.2022.02.002

1075

1076 Dyson, B. J. (2021, Feb 3). Variability in competitive decision-making speed and quality against
1077    exploiting and exploitative opponents. *Sci Rep, 11*(1), 2859. https://doi.org/10.1038/s41598-
1078    021-82269-2

1079

1080 Dyson, B. J., Sundvall, J., Forder, L., & Douglas, S. (2018, Oct). Failure generates impulsivity only
1081    when outcomes cannot be controlled. *J Exp Psychol Hum Percept Perform, 44*(10), 1483-1487.
1082    https://doi.org/10.1037/xhp0000557

1083

1084 Fleming, S. M., Mars, R. B., Gladwin, T. E., & Haggard, P. (2009, Oct). When the brain changes its
1085    mind: flexibility of action selection in instructed and free choices. *Cereb Cortex, 19*(10), 2352-
1086    2360. https://doi.org/10.1093/cercor/bhn252

1087

1088 Forder, L., & Dyson, B. J. (2016, Sep 23). Behavioural and neural modulation of win-stay but not
1089    lose-shift strategies as a function of outcome value in Rock, Paper, Scissors. *Sci Rep, 6*, 33809.
1090    https://doi.org/10.1038/srep33809

1091

1092 Frith, C. (2013, Sep). The psychology of volition. *Exp Brain Res, 229*(3), 289-299.
1093    https://doi.org/10.1007/s00221-013-3407-6

1094

1095 Gehring, W. J., & Willoughby, A. R. (2002, Mar 22). The medial frontal cortex and the rapid
1096    processing of monetary gains and losses. *Science, 295*(5563), 2279-2282.
1097    https://doi.org/10.1126/science.1066893

1098

1099 Haggard, P. (2008, Dec). Human volition: towards a neuroscience of will. *Nat Rev Neurosci, 9*(12),
1100    934-946. https://doi.org/10.1038/nrn2497

1101

1102 Haggard, P. (2019, Jan 4). The Neurocognitive Bases of Human Volition. *Annu Rev Psychol, 70*, 9-
1103    28. https://doi.org/10.1146/annurev-psych-010418-103348

1104

1105    Hajcak, G., Moser, J. S., Holroyd, C. B., & Simons, R. F. (2006, Feb). The feedback-related
1106        negativity reflects the binary evaluation of good versus bad outcomes. *Biol Psychol, 71*(2), 148-
1107        154. https://doi.org/10.1016/j.biopsycho.2005.04.001

1109    Hampton, A. N., Bossaerts, P., & O'Doherty, J. P. (2008, May 6). Neural correlates of mentalizing-
1110        related computations during strategic interactions in humans. *Proc Natl Acad Sci U S A,*
1111        *105*(18), 6741-6746. https://doi.org/10.1073/pnas.0711099105

1113    Hurvich, C. M., & Tsai, C.-L. (1989). Regression and time series model selection in small samples.
1114        *Biometrika, 76*(2), 297-307. https://doi.org/10.1093/biomet/76.2.297

1116    Jahanshahi, M., Dirnberger, G., Fuller, R., & Frith, C. D. (2000, Dec). The role of the dorsolateral
1117        prefrontal cortex in random number generation: a study with positron emission tomography.
1118        *Neuroimage, 12*(6), 713-725. https://doi.org/10.1006/nimg.2000.0647

1120    Jahanshahi, M., Jenkins, I. H., Brown, R. G., Marsden, C. D., Passingham, R. E., & Brooks, D. J.
1121        (1995, Aug). Self-initiated versus externally triggered movements. I. An investigation using
1122        measurement of regional cerebral blood flow with PET and movement-related potentials in
1123        normal    and    Parkinson's    disease    subjects. *Brain,    118    ( Pt    4)*,    913-933.
1124        https://doi.org/10.1093/brain/118.4.913

1126    Jenkins, I. H., Jahanshahi, M., Jueptner, M., Passingham, R. E., & Brooks, D. J. (2000, Jun). Self-
1127        initiated versus externally triggered movements. II. The effect of movement predictability on
1128        regional    cerebral    blood    flow. *Brain,    123    ( Pt    6)*,    1216-1228.
1129        https://doi.org/10.1093/brain/123.6.1216

1131    Kievit, R. A., Brandmaier, A. M., Ziegler, G., van Harmelen, A. L., de Mooij, S. M. M., Moutoussis,
1132        M., Goodyer, I. M., Bullmore, E., Jones, P. B., Fonagy, P., Consortium, N., Lindenberger, U., &
1133        Dolan, R. J. (2018, Oct). Developmental cognitive neuroscience using latent change score
1134        models:    A    tutorial    and    applications. *Dev    Cogn    Neurosci,    33*,    99-117.
1135        https://doi.org/10.1016/j.dcn.2017.11.007

1137    Kievit, R. A., Lindenberger, U., Goodyer, I. M., Jones, P. B., Fonagy, P., Bullmore, E. T., & Dolan, R.
1138        J. (2017, Jun). Mutualistic coupling between vocabulary and reasoning supports cognitive

1139    development during late adolescence and early adulthood *Psychological Science 28*(10), 1419-
1140    1431. https://doi.org/https://doi.org/10.1177/0956797617710785
1141
1142    Kuroki, D. (2021, Feb). A new jsPsych plugin for psychophysics, providing accurate display duration
1143    and stimulus onset asynchrony. *Behav Res Methods, 53*(1), 301-310.
1144    https://doi.org/10.3758/s13428-020-01445-w
1145
1146    Lashley, K. S. (1951). The problem of serial order in behavior. In L. A. Jeffress (Ed.), *Cerebral*
1147    *mechanisms in behavio*r (pp. 112–131). *Wiley, New York*.
1148
1149    Lee, D., Conroy, M. L., McGreevy, B. P., & Barraclough, D. J. (2004, Dec). Reinforcement learning
1150    and decision making in monkeys during a competitive game. *Brain Res Cogn Brain Res, 22*(1),
1151    45-58. https://doi.org/10.1016/j.cogbrainres.2004.07.007
1152
1153    Lee, D., McGreevy, B. P., & Barraclough, D. J. (2005, Oct). Learning and decision making in
1154    monkeys during a rock-paper-scissors game. *Brain Res Cogn Brain Res, 25*(2), 416-430.
1155    https://doi.org/10.1016/j.cogbrainres.2005.07.003
1156
1157    Lejuez, C. W., Read, J. P., Kahler, C. W., Richards, J. B., Ramsey, S. E., Stuart, G. L., Strong, D.
1158    R., & Brown, R. A. (2002, Jun). Evaluation of a behavioral measure of risk taking: the Balloon
1159    Analogue Risk Task (BART). *J Exp Psychol Appl, 8*(2), 75-84. https://doi.org/10.1037//1076-
1160    898x.8.2.75
1161
1162    Libet, B., Gleason, C. A., Wright, E. W., & Pearl, D. K. (1983, Sep). Time of conscious intention to
1163    act in relation to onset of cerebral activity (readiness-potential). The unconscious initiation of a
1164    freely voluntary act. *Brain, 106 (Pt 3)*, 623-642. https://doi.org/10.1093/brain/106.3.623
1165
1166    McArdle, J. J. (2009). Latent variable modeling of differences and changes with longitudinal data.
1167    *Annu Rev Psychol, 60*, 577-605. https://doi.org/10.1146/annurev.psych.60.110707.163612
1168
1169    McClure, S. M., Berns, G. S., & Montague, P. R. (2003). Temporal Prediction Errors in a Passive
1170    Learning Task Activate Human Striatum. *Neuron, 38*(2), 339-346. https://doi.org/DOI:
1171    10.1016/s0896-6273(03)00154-5
1172

1173  Mischel, W., Ebbesen, E. B., & Zeiss, A. R. (1972, Feb). Cognitive and attentional mechanisms in
1174      delay of gratification. *J Pers Soc Psychol, 21*(2), 204-218. https://doi.org/10.1037/h0032198
1175

1176  Neiman, T., & Loewenstein, Y. (2011, Dec 6). Reinforcement learning in professional basketball
1177      players. *Nat Commun, 2*, 569. https://doi.org/10.1038/ncomms1580
1178

1179  O'Doherty, J., Dayan, P., Schultz, J., Deichmann, R., Friston, K., & Dolan, R. J. (2004, Apr 16).
1180      Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science, 304*(5669),
1181      452-454. https://doi.org/10.1126/science.1094285
1182

1183  Ota, K., Shinya, M., & Kudo, K. (2016, Nov 21). Sub-optimality in motor planning is retained
1184      throughout 9 days practice of 2250 trials. *Sci Rep, 6*, 37181. https://doi.org/10.1038/srep37181
1185

1186  Ota, K., Shinya, M., Maloney, L. T., & Kudo, K. (2019, Oct 16). Sub-optimality in motor planning is
1187      not improved by explicit observation of motor uncertainty. *Sci Rep, 9*(1), 14850.
1188      https://doi.org/10.1038/s41598-019-50901-x
1189

1190  Ota, K., Tanae, M., Ishii, K., & Takiyama, K. (2020, Jan 22). Optimizing motor decision-making
1191      through competition with opponents. *Sci Rep, 10*(1), 950. https://doi.org/10.1038/s41598-019-
1192      56659-6
1193

1194  Rigoux, L., Stephan, K. E., Friston, K. J., & Daunizeau, J. (2014, Jan 1). Bayesian model selection
1195      for    group    studies    -    revisited.    *Neuroimage,    84*,    971-985.
1196      https://doi.org/10.1016/j.neuroimage.2013.08.065
1197

1198  Robbins, T. W., & Costa, R. M. (2017, Nov 20). Habits. *Curr Biol, 27*(22), R1200-R1206.
1199      https://doi.org/10.1016/j.cub.2017.09.060
1200

1201  Rosenbaum, D. A., Cohen, R. G., Jax, S. A., Weiss, D. J., & van der Wel, R. (2007, Aug). The
1202      problem of serial order in behavior: Lashley's legacy. *Hum Mov Sci, 26*(4), 525-554.
1203      https://doi.org/10.1016/j.humov.2007.04.001
1204

1205  Rosseel, Y. (2012). lavaan: AnRPackage for Structural Equation Modeling. *Journal of Statistical
1206      Software, 48*(2). https://doi.org/10.18637/jss.v048.i02
1207

1208     Schermelleh-Engel, K., Moosbrugger, H., & Müller, H. (2003). Evaluating the fit of structural
1209        equation models: Tests of significance and descriptive goodness-of-fit measures. *Methods of*
1210        *Psychological Research Online, 8*, 23–74.

1211

1212     Stephan, K. E., Penny, W. D., Daunizeau, J., Moran, R. J., & Friston, K. J. (2009, Jul 15). Bayesian
1213        model selection for group studies. *Neuroimage, 46*(4), 1004-1017.
1214        https://doi.org/10.1016/j.neuroimage.2009.03.025

1215

1216     Story, G. W., Vlaev, I., Seymour, B., Darzi, A., & Dolan, R. J. (2014). Does temporal discounting
1217        explain unhealthy behavior? A systematic review and reinforcement learning perspective. *Front*
1218        *Behav Neurosci, 8*, 76. https://doi.org/10.3389/fnbeh.2014.00076

1219

1220     Sundvall, J., & Dyson, B. J. (2022). Breaking the bonds of reinforcement: Effects of trial outcome,
1221        rule consistency and rule complexity against exploitable and unexploitable opponents. *PLoS*
1222        *One, 17*(2), e0262249. https://doi.org/10.1371/journal.pone.0262249

1223

1224     Sutton, R. S., & Barto, A. G. (2018). Reinforcement learning: An Introduction, 2nd edition. *MIT*
1225        *Press, Cambridge*.

1226

1227     Tang, R., Bugg, J. M., Snijder, J. P., Conway, A. R., & Braver, T. S. (2022, Aug 10). The Dual
1228        Mechanisms of Cognitive Control (DMCC) project: Validation of an online behavioural task
1229        battery. *Q J Exp Psychol (Hove)*, 17470218221114769.
1230        https://doi.org/10.1177/17470218221114769

1231

1232     Tervo, D. G. R., Proskurin, M., Manakov, M., Kabra, M., Vollmer, A., Branson, K., & Karpova, A. Y.
1233        (2014, Sep 25). Behavioral variability through stochastic choice and its gating by anterior
1234        cingulate cortex. *Cell, 159*(1), 21-32. https://doi.org/10.1016/j.cell.2014.08.037

1235

1236     Vickery, T. J., Chun, M. M., & Lee, D. (2011, Oct 6). Ubiquity and specificity of reinforcement
1237        signals throughout the human brain. *Neuron, 72*(1), 166-177.
1238        https://doi.org/10.1016/j.neuron.2011.08.011

1239

1240     Wang, Z., Xu, B., & Zhou, H. J. (2014, Jul 25). Social cycling and conditional responses in the
1241        Rock-Paper-Scissors game. *Sci Rep, 4*, 5830. https://doi.org/10.1038/srep05830

1242

1243  Wong, A., Merholz, G., & Maoz, U. (2021, Oct 19). Characterizing human random-sequence
1244      generation in competitive and non-competitive environments using Lempel-Ziv complexity. *Sci
1245      Rep, 11*(1), 20662. https://doi.org/10.1038/s41598-021-99967-6
1246
1247  Wood, W., & Runger, D. (2016). Psychology of Habit. *Annu Rev Psychol, 67*, 289-314.
1248      https://doi.org/10.1146/annurev-psych-122414-033417
1249
1250  Worthy, D. A., Hawthorne, M. J., & Otto, A. R. (2013, Apr). Heterogeneity of strategy use in the
1251      Iowa gambling task: a comparison of win-stay/lose-shift and reinforcement learning models.
1252      *Psychon Bull Rev, 20*(2), 364-371. https://doi.org/10.3758/s13423-012-0324-9
1253
1254  Zhu, L., Mathewson, K. E., & Hsu, M. (2012, Jan 31). Dissociable neural representations of
1255      reinforcement and belief prediction errors underlie strategic learning. *Proc Natl Acad Sci U S A,
1256      109*(5), 1419-1424. https://doi.org/10.1073/pnas.1116783109
1257
1258
1259
1260