

# 1 **Single-Trial Inhibition of Anterior Cingulate Disrupts Model-based** 2 **Reinforcement Learning in a Two-step Decision Task.**

3  
4 Thomas Akam<sup>1,2</sup>, Inês Rodrigues-Vaz<sup>1,3</sup>, Xiangyu Zhang<sup>4</sup>, Michael Pereira<sup>1</sup>, Rodrigo Oliveira<sup>1</sup>, Peter  
5 Dayan<sup>4</sup>, Rui M. Costa<sup>1,3</sup>.

6 Affiliations:

7 1. Champalimaud Neuroscience Program, Champalimaud Centre for the Unknown, Lisbon, Portugal

8 2. Department of Experimental Psychology, Oxford University, Oxford, UK

9 3. Department of Neuroscience, Zuckerman Mind Brain Behavior Institute, Columbia University, New  
10 York, NY, USA.

11 4. RIKEN-MIT Center for Neural Circuit Genetics at the Picower Institute for Learning and Memory,  
12 Department of Biology and Department of Brain and Cognitive Sciences. Massachusetts Institute of  
13 Technology, Cambridge, Massachusetts, USA.

14 5. Gatsby Computational Neuroscience Unit, UCL, London, UK.

15  
16 **Abstract:**

17  
18 The anterior cingulate cortex (ACC) is implicated in learning the value of actions, and thus in allowing  
19 past outcomes to influence the current choice. However, it is not clear whether or how it  
20 contributes to the two major ways such learning is thought to happen: model-based mechanisms  
21 that learn action-state predictions and use these to infer action values; or model-free mechanisms  
22 which learn action values directly through reward prediction errors. Having confirmed, using a  
23 classical probabilistic reversal learning task, that optogenetic inhibition of ACC neurons on single  
24 trials indeed affected reinforcement learning, we examined the consequence of this manipulation in  
25 a novel two-step decision task designed to dissociate model-free and model-based learning  
26 mechanisms in mice. On the two-step task, silencing spared the influence of the trial outcome but  
27 reduced the influence of the experienced state transition. Analysis using reinforcement learning  
28 models indicated that ACC inhibition disrupted model-based RL mechanisms.

## 30 Introduction:

31

32 The anterior cingulate cortex (ACC) has long been implicated in reward guided decision making  
33 (Rushworth et al., 2004; Rushworth and Behrens, 2008). ACC neurons encode diverse decision  
34 variables (Cai and Padoa-Schioppa, 2012; Ito et al., 2003; Matsumoto et al., 2003; Sul et al., 2010),  
35 but ACC has been particularly associated with action reinforcement (Hadland et al., 2003; Kennerley  
36 et al., 2006; Rudebeck et al., 2008). However, instrumental learning is not a unitary phenomenon  
37 but rather is thought to be mediated by parallel control systems which use different computational  
38 principles to evaluate choices (Balleine and Dickinson, 1998; Daw et al., 2005; Dolan and Dayan,  
39 2013). It has recently become a pressing problem to understand the neural underpinnings of these  
40 controllers and their interactions.

41 In familiar environments when executing well practiced actions, behaviour is apparently controlled  
42 by a habitual system thought to employ model-free reinforcement learning (RL) (Sutton and Barto,  
43 1998). Model-free RL uses reward prediction errors to acquire or cache preferences between  
44 actions. However, when the environment or motivational state changes, model-free preferences can  
45 become out of date, and actions are instead determined by a goal-directed system believed to  
46 follow the precepts of model-based RL (Sutton and Barto, 1998). Model-based RL learns a predictive  
47 model of the consequences of actions, i.e. the states and rewards to which they typically  
48 immediately lead, and evaluates options by using the model to simulate or otherwise estimate their  
49 resulting long-run outcomes. Such a dual controller approach is beneficial because model-free and  
50 model-based RL possess complementary strengths, the former allowing quick and computationally  
51 cheap decision making at the cost of slower adaptation to changes in the environment, the latter  
52 flexible and efficient use of new information at the cost of computational effort and decision speed.

53 On specific anatomical and physiological grounds, we hypothesised that ACC is a component of the  
54 model-based control system. Firstly, the ACC provides a massive input to posterior dorsomedial  
55 striatum (Oh et al., 2014; Hintiryan et al., 2016), a region critical for model-based control as assessed  
56 through outcome-devaluation (Yin et al., 2005a, 2005b; Hilario et al., 2012). Secondly, decision  
57 related signals in ACC suggest that it plays a role in representing task contingencies beyond model-  
58 free cached values (Daw et al., 2011; Cai and Padoa-Schioppa, 2012; Karlsson et al., 2012; O'Reilly et  
59 al., 2013; Doll et al., 2015). We therefore sought to test the role of ACC in a reward guided decision  
60 task able to dissociate model-based and model-free mechanisms.

61 The classical approach to dissociating the systems in the laboratory involves outcome devaluation  
62 (Adams and Dickinson, 1981; Colwill and Rescorla, 1985). A subject is first trained to perform an

63 action to receive a reward. The reward is then devalued, e.g. through pairing with illness, and the  
64 subject's subsequent tendency to perform the action is tested in extinction, i.e. without further  
65 rewards being delivered. If the action is mediated by a model-based prediction of the specific  
66 outcome to which it leads, devaluing that outcome will reduce the tendency to perform the action.  
67 If, on the other hand, the action is mediated by a cached model-free action value, devaluation will  
68 have no effect (Balleine and Dickinson, 1998; Daw et al., 2005). Learned actions often transition  
69 from being devaluation sensitive or goal-directed early in learning to being devaluation insensitive or  
70 habitual after extensive training (Dickinson et al., 1983; Dickinson, 1985). Lesion and inactivation  
71 studies using outcome devaluation indicate that goal-directed and habitual behaviours rely on  
72 partially separate cortical-basal ganglia circuits (Balleine et al., 2003; Killcross and Coutureau, 2003;  
73 Ostlund and Balleine, 2005; Yin et al., 2004, 2005b; Hilario et al., 2012; Gremel and Costa, 2013a,  
74 2013b).

75 Unfortunately, outcome devaluation has limitations as a paradigm for decision neuroscience. Firstly,  
76 the critical devaluation test during which behavioural strategies are dissociated must be short  
77 because it is performed in extinction, limiting the number of choices or actions performed.  
78 Secondly, devaluation is a unidirectional single-shot manipulation of value. Neurophysiology thrives  
79 on behavioural paradigms that generate large decision datasets with parametric variation of decision  
80 variables. However, in workhorse tasks such as perceptual decision making or probabilistic reversal  
81 learning, the only uncertainty about the outcome of each decision is whether reward will be directly  
82 delivered. Thus, model-based prediction of future state and model-free prediction of future reward  
83 are ineluctably confounded.

84 Instead, at least for human subjects, novel tasks have recently been developed which aim to  
85 distinguish model-free and model-based reasoning in a stable manner over many trials. These tasks  
86 generally require subjects to take multiple steps through a decision tree to reach rewards, thus  
87 licensing the simulation-based search that is characteristic of the model-based controller (Daw et al.,  
88 2011; Simon and Daw, 2011; Huys et al., 2012). The most widely used is the so called two-step task  
89 (Daw et al., 2011), in which a choice between two actions leads probabilistically to one of two  
90 different states, in which further actions lead probabilistically to reward. Daw's two-step task has  
91 been used to assess the influence on behavioural strategy of behavioural (Otto et al., 2013, 2014)  
92 and neuronal manipulations (Wunderlich et al., 2012; Smittenaar et al., 2013), genetic factors (Doll  
93 et al., 2016), psychiatric illness (Sebold et al., 2014; Voon et al., 2015), and variants have also been  
94 used to examine more mechanistic aspects of interaction between the systems (Lee et al., 2014;  
95 Keramati et al., 2016; Doll et al., 2015). There is substantial interest from a number of groups in  
96 developing versions of the task for animal subjects to permit the use of more powerful neuroscience

97 tools (Miller et al. Soc. Neurosci. Abstracts 2013, 855.13, Groman et al. Soc. Neurosci. Abstracts  
98 2014, 558.19, Miranda et al. Soc. Neurosci. Abstracts 2014 756.09).

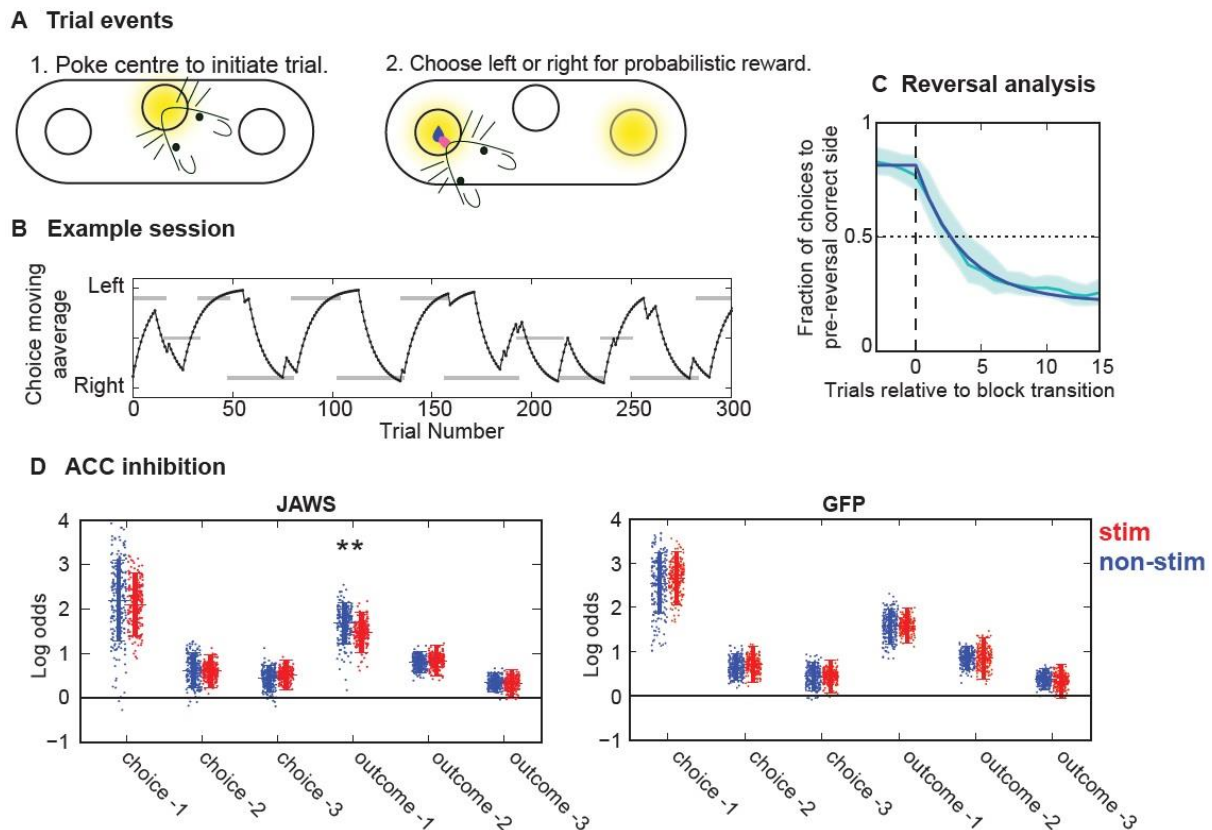
99 Here, we report our adaptation of the two-step task to study model-based and model-free learning  
100 in mice, and the use of our novel variant to probe the involvement of the anterior cingulate cortex  
101 (ACC), a region expected to be centrally involved. Based on an in depth computational analysis  
102 (Akam et al., 2015), we substantially modified the implementation and structure of the task,  
103 developing a new version in which both the reward probabilities in the leaf states of the decision  
104 tree and the action-state transition probabilities change over time. Here, detailed characterisation  
105 of subjects' behaviour indicated that, as in the human version, choices were guided by a mixture of  
106 model-based and model-free RL. However, we also observed a number of previously unexplored  
107 characteristics, including forgetting about actions that were not chosen, perseverative influences  
108 that spanned multiple trials, and representation of actions both in terms of the choices they  
109 represent and the motor output they require.

110 We found that optogenetic silencing of ACC neurons on individual trials reduced the influence of the  
111 experienced state transition on subsequent choice without affecting the influence of the trial  
112 outcome (rewarded or not). Analysis using RL models suggested this effect was due to reduced  
113 influence of model-based RL on ACC inhibition trials. For comparison purposes we performed the  
114 same ACC manipulation in a standard probabilistic reversal learning task, where it reduced the  
115 influence of the previous trial outcome on subsequent choice. These data are consistent with  
116 subjects using a combination of model-based and model-free RL in both tasks, but with the two-step  
117 task uniquely allowing a dissociation of their respective contributions to choice behaviour.

## 118 Results:

### 119 *Single-trial inhibition of ACC impairs probabilistic reversal learning.*

120 To confirm that ACC is involved in reward-guided decision making in mice, we first assessed whether  
121 optogenetic silencing of ACC neurons affected decision making in a standard probabilistic reversal  
122 learning task (Figure 1). Mice were trained to initiate each trial in a central nose-poke port which  
123 was flanked by left and right poke ports (Figure 1A). Trial initiation caused the left and right pokes  
124 to light up and subjects then chose between them for the chance of obtaining a water reward.  
125 Reward probabilities changed in blocks, with three block types; *left good* (left=0.75/right=0.25),  
126 *neutral* (0.5/0.5) and *right good* (0.25/0.75). Subject's choices tracked which option had higher  
127 reward probability (Figure 1B, C), choosing the correct option at the end of non-neutral blocks with  
128 probability  $0.80 \pm 0.04$  (mean  $\pm$  SD), and adapting to reversals in the reward probability with a time  
129 constant of 3.57 trials (exponential fit tau).



130  
 131 **Figure 1. Optogenetic silencing of ACC in probabilistic reversal learning task.** **A)** Diagram of apparatus and  
 132 trial events. **B)** Example session, black line shows exponential moving average ( $\tau = 8$  trials) of choices, grey  
 133 bars indicate reward probability blocks with y position of bar indicating whether left or right side has high  
 134 reward probability or a neutral block. **C)** Choice probability trajectories around reversal in reward probabilities:  
 135 Pale blue line – average trajectory, dark blue line – exponential fit, shaded area – cross-subject standard  
 136 deviation. **D)** Logistic regression analysis showing predictor loadings for stimulated (red) and non-stimulated  
 137 (blue) trials, for the ACC JAWS (left panel) and GFP controls (right panel). Bars indicate  $\pm 1$  standard deviation  
 138 of the population level distributions, dots indicate maximum a posteriori session fits. \*\* indicates significant  
 139 difference ( $P < 0.01$ ) between stimulated and non-stimulated trials.

140 The following figure supplements are available for figure 1.

141 [Figure supplement 1. JAWS inhibition of ACC neurons.](#)

142 [Figure supplement 2. Average JAWS expression.](#)

143

144 We silenced the activity of ACC neurons on individual trials using the red-shifted halorhodopsin  
 145 JAWS (Chuong et al., 2014). An AAV viral vector expressing JAWS-GFP under the CaMKII promoter  
 146 was injected bilaterally into ACC of experimental animals ( $n=10$  JAWS), while control animals ( $n=10$ )  
 147 were injected with an AAV expressing GFP under the CaMKII promoter. Illumination was provided  
 148 by a high power red LED chronically implanted above the cortical surface (Figure 1 - figure  
 149 supplement 1). Electrophysiological recordings in animals implanted with micro-wire bundles ( $n=2$ )  
 150 confirmed that red light (50mW, 630nm) from the implanted LEDs robustly inhibited ACC neurons  
 151 (Figure 1- figure supplement 1). ACC neurons were inhibited using JAWS on a randomly selected  
 152 1/6 trials, with a minimum of two non-stimulated trials between each stimulated trial. Stimulation  
 153 was delivered from when subjects poked in the side poke and received the trial outcome until the

154 subsequent choice. The dataset comprised 12855 stimulated and 65186 non-stimulated trials for  
155 the JAWS animals and 11096 stimulated and 55913 non-stimulated trials for the controls.

156 We assessed the effect of ACC silencing using a logistic regression analysis with previous choices and  
157 outcomes as regressors. We separately analysed choices made during stimulation and on non-  
158 stimulated trials and used permutation tests to identify significant differences between the predictor  
159 loadings in the two conditions (Figure 1D). Previous choices predicted current choice with  
160 decreasing loading at increasing lag relative to the current trial. Obtaining reward further predicted  
161 repeating the rewarded choice, again with decreasing loading at increasing lag. ACC inhibition  
162 significantly reduced the influence of the most recent outcome (i.e., whether reward was received)  
163 on subsequent choice (permutation test  $P = 0.004$  uncorrected,  $P = 0.024$  Bonferroni corrected for 6  
164 predictors), but did not affect the influence of either previous choices or earlier outcomes ( $P > 0.18$   
165 uncorrected). Light stimulation did not affect the influence of previous outcomes or choices on  
166 subsequent choice in the GFP controls ( $P > 0.38$  uncorrected) and the stimulation-by-group  
167 interaction was significant for the influence of the most recent outcome on choice ( $P = 0.014$ ,  
168 permutation test).

169 These data indicate that transient ACC silencing disrupted reward-guided decision making in the  
170 probabilistic reversal learning task, however this task does not discriminate whether this was due to  
171 an effect on model-free mechanisms which learn action values directly, or model-based mechanisms  
172 which learn action-state transition probabilities and use these to guide choice. We therefore  
173 performed the same optogenetic manipulation in a multi-step decision task designed to dissociate  
174 the contribution of model-based and model-free reinforcement learning.

#### 175 *Development of a novel two-step task for mice*

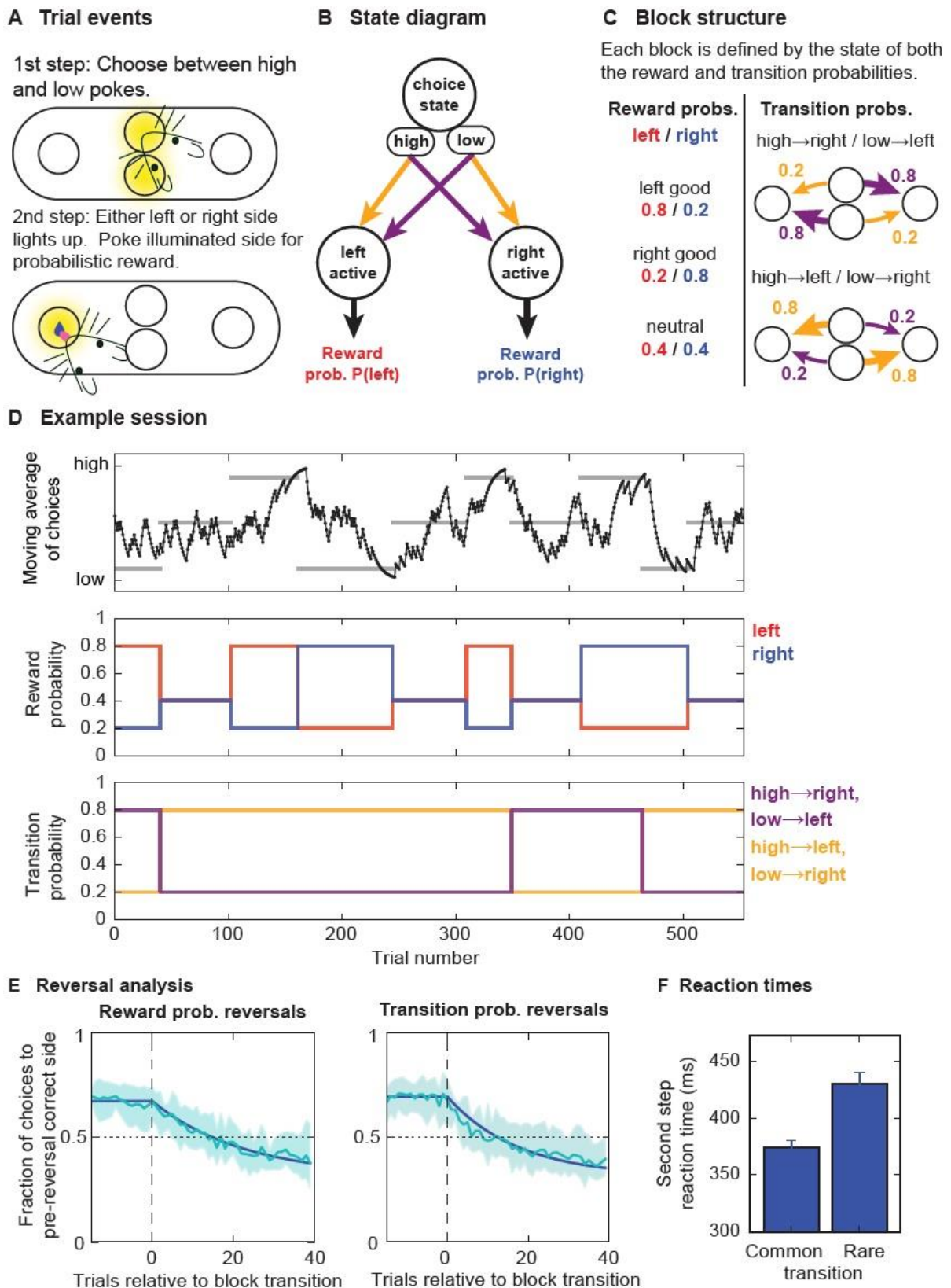
176 The task was based on that developed for humans by Daw et al. (2011) but both the physical format  
177 in which it was presented to subjects and the task structure were heavily adapted for use with mice.  
178 We first summarise changes to the task structure and their rationale before detailing the task  
179 implementation. As in the Daw two-step task, our version consisted of a choice between two 'first-  
180 step' actions which lead probabilistically to one of two 'second-step' states where reward could be  
181 obtained. Unlike the Daw task, in each second-step state there was a single action rather than a  
182 choice between two actions available, reducing the number of reward probabilities the subject must  
183 track from four to two (Figure 2 – figure supplement 1). In the original task, the stochasticity of the  
184 state transitions and reward probabilities caused both model-based and model-free control to  
185 obtain rewards at a rate negligibly different from random choice at the first-step (Akam et al., 2015;  
186 Kool et al., 2016). To promote task engagement, we increased the contrast between good and bad

187 options by using a block-based reward probability distribution rather than the random walks used in  
188 the original, and by increasing the probability of common state transitions (see below) from 0.7 to  
189 0.8. The final, and most significant, structural change was the introduction of reversals in the  
190 transition probabilities mapping the first-step actions to the second-step states. This step was taken  
191 to preclude subjects developing habitual strategies consisting of mappings from second-step states  
192 in which rewards had recently been obtained to specific actions at the first step (e.g. rewards in  
193 state  $X \rightarrow$  chose action  $x$ , where action  $x$  is that which commonly leads to state  $X$ ). Such strategies  
194 can, in principle, generate behaviour that looks very similar to model-based control despite not using  
195 a forward model which predicts the future state given chosen action (see Akam et al. (2015) for a  
196 detailed discussion).

197 We implemented the task using a set of four nose-poke ports: a low and a high poke in the centre,  
198 flanked by a left and a right poke (Figure 2A). Each trial started with the central pokes lighting up,  
199 mandating a choice. The resulting action led probabilistically to one of two states termed 'left-  
200 active' and 'right-active', in which respectively the left or right poke was illuminated. The subject  
201 then had to poke the illuminated side to gain a probabilistic water reward (Figure 2A,B). A 1 second  
202 inter-trial interval started from when the subject exited the side port at the end of the trial. The next  
203 trial then started with the illumination of the central pokes.

204 Both the transition probabilities linking the first-step actions to the second-step states, and the  
205 reward probabilities in each second-step state, changed in blocks (Figure 2C, D), such that each block  
206 was defined by the state of both the transition and reward probabilities. There were three possible  
207 states of the reward probabilities: *left good* (left=0.8/right=0.2), *neutral* (0.4/0.4) and *right good*  
208 (0.2/0.8). There were two possible states of the transition probabilities: *high  $\rightarrow$  right / low  $\rightarrow$  left*, in  
209 which the high poke commonly (80% of trials) gave access to the right-active state and the low poke  
210 commonly gave access to the left-active state, and *high  $\rightarrow$  left / low  $\rightarrow$  right* in which the high poke  
211 commonly gave access to the left-active, and the low poke commonly gave access to the right-active  
212 state. In either case, on 20% of trials, a rare transition occurred such that each first-step action gave  
213 access to the state commonly reached from the other first-step action. At block transitions, either  
214 the reward probabilities or the transition probabilities changed, except on transitions to neutral  
215 blocks, 50% of which were accompanied by a change in the transition probabilities (See Fig S3 for full  
216 block transition structure). Reversals in which first-step action (high or low) had higher reward  
217 probability, could therefore occur either due to the reward probabilities of the second-step states  
218 reversing, or due to the transition probabilities linking the first-step actions to the second-step states  
219 reversing. Block transitions were triggered based on a behavioural criterion (see methods) which  
220 resulted in block lengths of  $63.6 \pm 31.7$  (mean  $\pm$  SD) trials.





221

222 **Figure 2. Two-step task.** **A)** Diagram of apparatus and trial events. **B)** State diagram of task. **C)** Block  
 223 structure, left side shows the three possible states of the reward probabilities, right side shows the two  
 224 possible states of the transition probabilities. **D)** Example session: Top panel - Exponential moving average ( $\tau$   
 225 = 8 trials) of choices. Horizontal grey bars show blocks, with correct choice (high, low or neutral) indicated by y  
 226 position of bars. Middle panel – reward probabilities in left active (red) and right active (blue) states. Bottom



227 panel – Transition probabilities linking first-step actions (high, low pokes) to second step states (left/right  
228 active). **E)** Reversal analysis: Pale blue line – average trajectory, dark blue line – exponential fit, shaded area –  
229 cross-subject standard deviation. Left panel - reversals in reward probability, right panel – reversals in  
230 transition probabilities. **F)** Second step reaction times following common and rare transitions - i.e. the time  
231 between the first step choice and side poke entry. Error bars show cross-subject SEM.

232 The following figure supplements are available for figure 2.

233 [Figure supplement 1](#). Comparison of original and new two-step task structures.

234 [Figure supplement 2](#). Block transition probabilities.

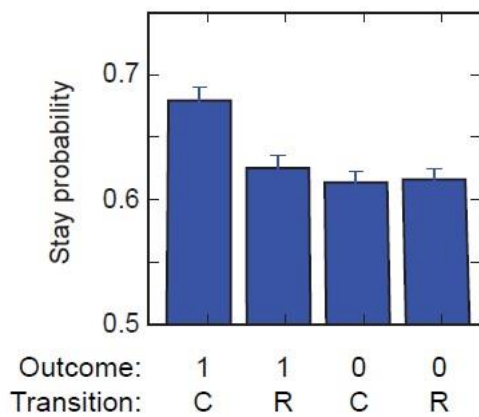
235 [Figure supplement 3](#). Body weight trajectory across training.

236

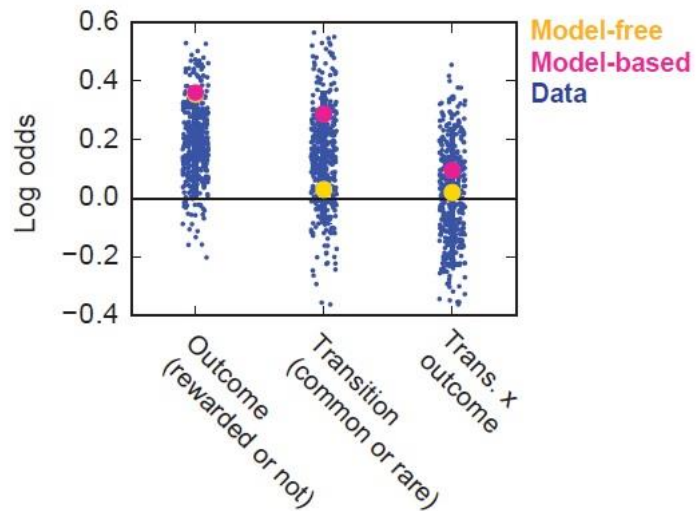
237 Subjects learned the task in 3 weeks with minimal shaping (see methods) and performed an average  
238 of  $576 \pm 174$  (mean  $\pm$  SD) trials per day thereafter. The baseline behavioural dataset consisted of  
239 sessions from day 22 of training onwards from 17 subjects, for a total of 400 sessions and 230237  
240 trials. Subject's choices tracked which first-step action had higher reward probability (Figure 2D,E),  
241 choosing the correct option at the end of non-neutral blocks with probability  $0.68 \pm 0.03$  (mean  $\pm$   
242 SD). Choice probabilities adapted faster ( $P = 0.009$ , bootstrap test) following block transitions in  
243 which the action-state transition probabilities reversed (exponential fit  $\tau = 17.6$  trials), compared  
244 with block transitions in which the reward probabilities in the two second-step states reversed ( $\tau =$   
245  $22.7$  trials, Figure 2E). Reaction times at the second step, i.e. the latency from when the left or right  
246 side illuminated till the subject poked in the corresponding port, were faster following common than  
247 rare transitions ( $P = 2.8 \times 10^{-8}$ , paired t-test) (Figure 2F).

248 The choice probability trajectories around reversals show that subjects tracked which choice is best,  
249 but do not discriminate whether they used model-based or model-free RL. Both strategies are  
250 capable of tracking the best option, but do so in different ways: a model-based strategy learns  
251 estimates of the transition-probabilities linking the first-step actions to second-step states, and the  
252 reward probabilities in these states, and calculates the expected value of choosing each first-step  
253 action by combining these. By contrast, a model-free strategy directly learns action values for the  
254 first-step actions through the reward prediction errors that occur when the second-step is reached,  
255 and, via what is known as an eligibility trace, when the outcome (rewarded or not) is obtained after  
256 the second-step. As these different strategies learn different representations of the world, which  
257 are updated in different ways based on experienced events, it may be possible to dissociate them  
258 based on the fine structure of how events on each trial affect subsequent choices. We employ both  
259 of the two analysis approaches that are traditionally employed to do this: logistic regression showing  
260 how events on each trial affect subsequent choices, and direct fitting to the behavioural data of  
261 combined model-based and model-free reinforcement learning models. We detail these approaches  
262 below, and use them to unpick the effects of silencing the ACC.

### A Stay probability analysis



### B Logistic regression: Data and simulation.



263 **Figure 3. Stay probability and logistic regression analyses.** **A)** Stay probability analysis. Fraction of trials the  
264 subject repeated the same choice following each combination of outcome (rewarded (1) or not (0)) and  
265 transition (common (C) or rare (R)). Error bars show cross-subject SEM. **B)** Logistic regression loadings for  
266 predictors; *outcome* (tendency to repeat choices following reward), *transition* (tendency to repeat choices  
267 following common transitions) and *transition-outcome interaction* (tendency to repeat choices following  
268 rewarded common transition trials and non-rewarded rare transition trials), comparing subject's data (blue)  
269 with simulated data from a model-free (yellow) and model-based (pink) agent fit to the subjects behaviour.  
270 For subjects data; blue bars indicate  $\pm 1$  standard deviation of the population level distributions, blue dots  
271 indicate maximum a posteriori (MAP) session fits. The full set of predictor loadings is shown in figure  
272 supplement 1.  
273

274 The following figure supplements are available for figure 3.

275 [Figure supplement 1. Full logistic regression model fit.](#)

276

277 *Logistic regression analysis to disambiguate model-based versus model-free strategies*

278 The simplest picture of behaviour is the raw so-called stay probabilities of repeating the first-step  
279 choice for the four possible combinations of transition and outcome (Figure 3A). Subjects were most  
280 likely to repeat choices following rewarded common transition trials, with a lower stay probability on  
281 rewarded rare-transition trials and non-rewarded trials. Logistic regression analyses of the  
282 relationship between choice and trial events test the nature of the interaction between transition  
283 and outcome, as this has historically been taken indicative of model-based reasoning. However,  
284 drawing such conclusions requires including various additional predictors in the model to capture  
285 strong, potentially confounding, effects. Some of these are conventional – for instance,  
286 accommodating perseveration or alternation between first-step choices and other direct biases of  
287 choice. However, we recently showed (Akam et al., 2015) the necessity of including an additional  
288 predictor which promotes repeating correct choices, as this avoids the effect of untoward  
289 correlations.

290 We therefore performed a logistic regression analysis which predicted stay probability as a function  
291 of trial events (outcome, transition and their interaction), with four additional regressors: the  
292 regressor discussed above which promoted repeating correct choices, a regressor which promoted  
293 repeating the previous choice, and two regressors capturing choice biases discussed below (Figure  
294 3B, Figure 3 – figures supplement 1). Positive loading on the outcome predictor indicated that  
295 receiving reward was reinforcing (i.e. predicted staying) ( $P < 0.001$ , bootstrap confidence interval).  
296 Positive loading on the transition predictor indicated that experiencing common transitions was also  
297 reinforcing ( $P < 0.001$ ). Loading on the transition-outcome interaction predictor was not significantly  
298 different from zero ( $P = 0.79$ ). The absence of transition-outcome interaction has been used in the  
299 context of the traditional Daw two-step task (Daw 2011) to suggest that behaviour is model-free.  
300 However, we have previously shown (Akam et al. 2015) that this depends on the subjects not  
301 learning the transition probabilities from the transitions they experience. Such fixedness is  
302 reasonable for the traditional task, for which the probabilities are fixed and are known to be so by  
303 the human subjects. It is not for our task. Our analysis (Akam et al. 2015) suggests that when model-  
304 learning is included, loading in the logistic regression analysis is shifted off the interaction predictor  
305 and onto the outcome and transition predictors.

306 To understand more precisely the implications of this analysis, we simulated the behaviour of a  
307 model-based and a model-free RL agent, with the parameters of both fit to the behavioural data,  
308 and performed the logistic regression analysis on the data simulated from both models (Figure 2B).  
309 Data simulated from the model-free agent showed a large loading on the outcome regressor (i.e.  
310 rewards were reinforcing), but minimal loading on the transition and transition-outcome interaction  
311 regressors. By contrast, data simulated from the model-based agent showed a large loading on both  
312 outcome and transition predictors (i.e. both rewards and common transitions were reinforcing), and  
313 a small loading on the interaction predictor. The robust loading on the transition predictor observed  
314 in the experimental data is therefore consistent with subjects using model-based control as a  
315 component of their behavioural strategy.

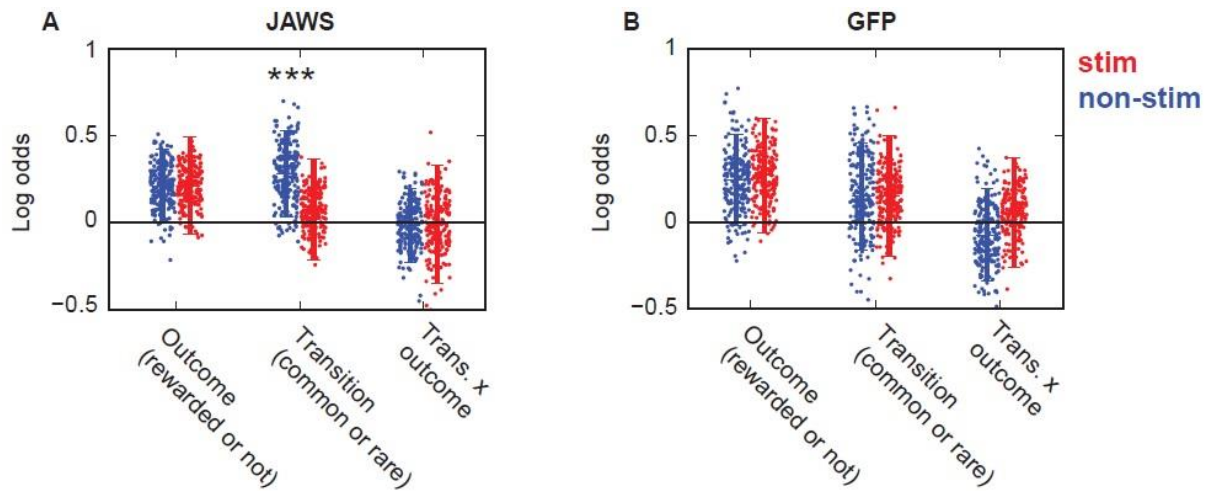
316 In addition to the three predictors reflecting the influence of the previous trial's events, positive  
317 loading on the 'stay' predictor (Figure 3 – figure supplement 1,  $P < 0.001$ ), indicated an overall  
318 tendency to repeat choices, consistent with the raw stay probabilities (Fig 3a). The 'correct'  
319 predictor also showed positive loading ( $P < 0.001$ ) indicating that subjects were more likely to repeat  
320 choices to the correct, i.e. higher reward probability option irrespective of the experienced trial  
321 outcome. Subjects showed a small bias towards the high poke ( $P < 0.001$ ) suggesting that the  
322 physical layout of the pokes made this action somewhat easier to execute. We included a second  
323 bias predictor which captured asymmetry in subject's bias dependent on the side they finish the

324 previous trial on, i.e. a positive loading on this predictor promoted a bias towards the high poke if  
325 the previous trial ended on the left side, and towards the low poke if the previous trial ended on the  
326 right side. We term this a ‘rotational’ bias as positive loading promotes clockwise movement around  
327 the set of pokes (e.g. left→high, right→low), while negative loading promotes counter-clockwise  
328 movement. Though loading on this predictor was not on average different from zero ( $P = 0.092$ ), it  
329 exhibited a substantial spread across the population of sessions such that a subset of sessions  
330 showed a strong rotational bias in either direction. Including this predictor substantially improved  
331 integrated Bayes Information Criterion (iBIC) scores for the regression model ( $\Delta$  iBIC = 2639)  
332 indicating it captured a real feature of the data. Subjects may have developed this form of bias  
333 because it is the simplest fixed response pattern that was not penalised by the block transition rule:  
334 As block transitions were triggered based on a moving average of correct choices, developing an  
335 overall bias for the high or low poke resulted in the favoured poke spending most of the time as the  
336 bad option. Rotational bias may therefore be a default action which could be quickly executed when  
337 there was little evidence to suggest one option was better than the other.

### 338 *Single-Trial Anterior Cingulate silencing in the two-step task impairs model based strategies*

339 Parameters for optogenetic silencing in the two-step task were as closely as possible matched to  
340 those used in the probabilistic reversal learning task, with the same viral vector, injection sites and  
341 light stimulation. Again, optogenetic inhibition was delivered on a randomly selected 1/6 of trials,  
342 with a minimum of two non-stimulated trials between each stimulation trial. Inhibition was  
343 delivered from when the subject entered the side poke and received the trial outcome until the  
344 subsequent choice. The JAWS dataset comprised 11 animals with 12827 stimulated and 64523 non-  
345 stimulated trials, the GFP control dataset 12 animals, 11663 stimulated and 59408 non-stimulated  
346 trials.

347 We evaluated the effect of ACC inhibition on behaviour by performing the logistic regression analysis  
348 separately for choices which occurred during stimulation and on non-stimulated trials. As in the  
349 baseline dataset, both experimental and control animals showed positive loading on both the  
350 outcome and transition predictors on non-stimulated trials, indicating that both receiving reward  
351 and experiencing common transitions was reinforcing (Figure4 A,B). Optogenetic inhibition of ACC  
352 neurons reduced the influence of the previous state transition (common or rare) on subjects  
353 subsequent choice ( $P < 0.0002$  uncorrected permutation test,  $P < 0.0006$  Bonferroni corrected for  
354 multiple comparison of 3 predictors, stimulation by group interaction  $P = 0.029$ ), but did not affect  
355 the influence of the previous reward ( $P = 0.94$  uncorrected), or the transition-outcome interaction ( $P$   
356 = 0.90 uncorrected).



357  
358 **Figure 4. Optogenetic silencing of ACC in two-step task. A)** Logistic regression analysis of ACC inhibition  
359 dataset showing loadings for the outcome, transition and transition-outcome interaction predictors for choices  
360 made on stimulated (red) and non-stimulated (blue) trials. **B)** As (a) but for GFP control animals. \*\*\* indicates  
361 significant difference ( $P < 0.001$ ) between stimulated and non-stimulated trials.

362 The following figure supplements are available for figure 4.

363 [Figure supplement 1. ACC inhibition stay probabilities.](#)

364 [Figure supplement 2. ACC inhibition full logistic regression model fits.](#)

365 [Figure supplement 3. ACC inhibition reaction times.](#)

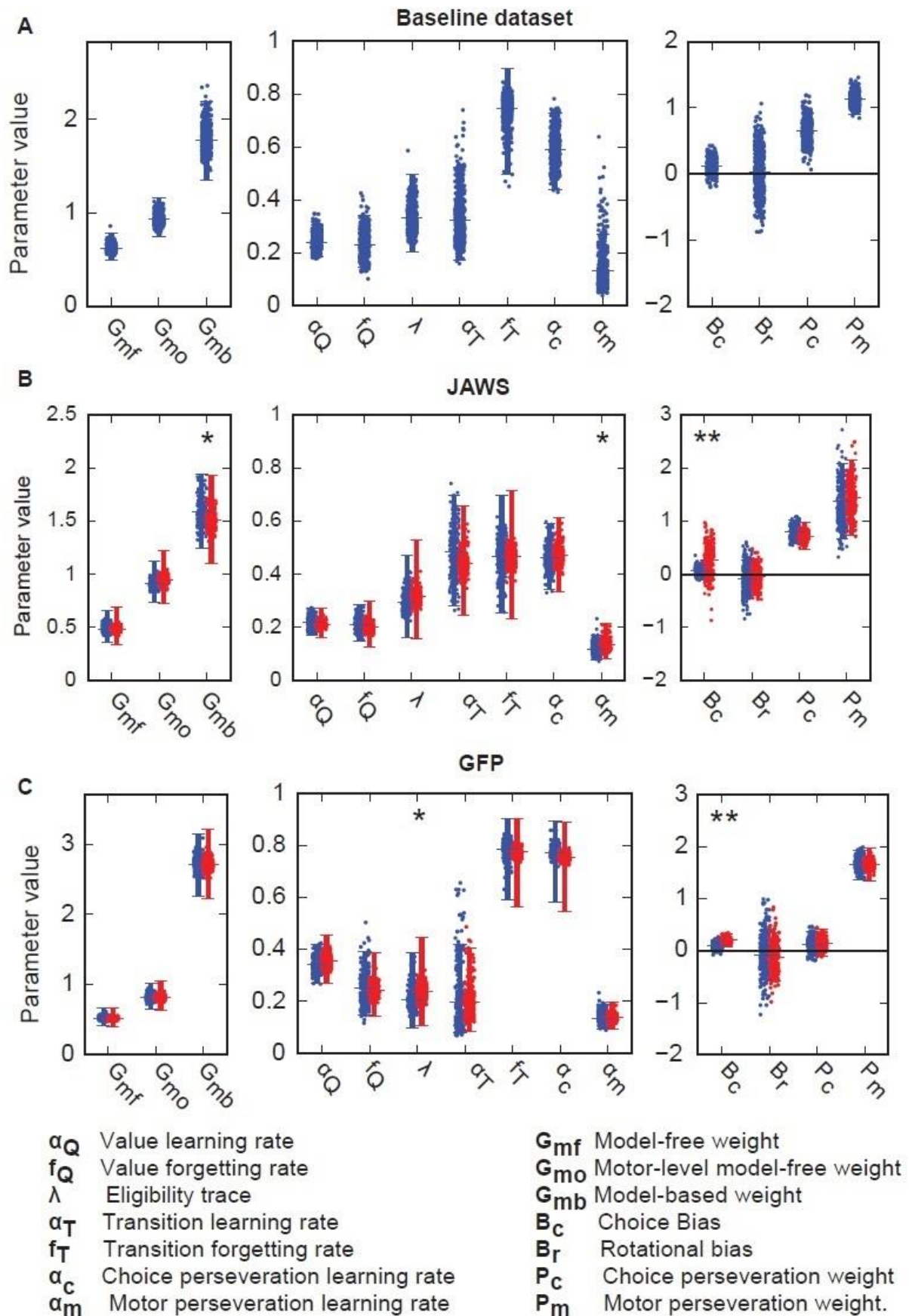
366

367 This selective reduction in influence of the previous state transition while sparing the influence of  
368 the previous trial outcome is consistent with a shift from model-based towards model-free control  
369 as it is the transition predictor which most strongly differentiates behaviour generated by these two  
370 strategies (Figure 3B). Neither outcome, transition nor transition-outcome interaction predictors  
371 were affected by light stimulation in the GFP controls (Bonferroni corrected  $P > 0.2$ ). In both  
372 experimental and control groups, light stimulation produced a small but significant bias towards the  
373 high poke, potentially reflecting an orienting response to the light (Bonferroni corrected  $P < 0.0015$ )  
374 (Figure 4 – figure supplement 1). Reaction times were not affected by light stimulation in either  
375 group (Paired t-test  $P > 0.36$ ) (Figure 4 – figure supplement 2).

#### 376 *Reinforcement learning model analysis*

377 To gain a sharper picture of the baseline behaviour and the effects of ACC silencing, we fitted and  
378 compared RL models to the respective datasets. Using our large baseline dataset, we performed an  
379 in-depth comparison of different RL models, as detailed in the supplementary material. Here, we  
380 summarise the principal findings. Our starting point was the RL agent used in the original Daw two-  
381 step task (Daw et al., 2011) in which behaviour is generated by a mixture of model-based and model-  
382 free strategies. Since the state transition probabilities change over time in our task, we modified the  
383 model to include ongoing learning about the transition probabilities.





384  
 385  
 386  
 387

**Figure 5. Reinforcement learning model fitting: A)** Parameter values for best fitting RL model on baseline dataset. Bars indicate  $\pm 1$  standard deviation of the population level distributions, dots indicate maximum a posteriori session fits. **B)** Reinforcement learning model fit to ACC inhibition dataset whose parameters take



388 separate values on stimulated (red) and non-stimulated (blue) trials. **C)** As (b) but for GFP control animals. \*  
389 indicates significant difference ( $P < 0.05$ ) between stimulated and non-stimulated trials, \*\* indicates  $P < 0.01$ .

390 The following figure supplements are available for figure 5.

391 [Figure supplement 1. Baseline dataset BIC score model comparison.](#)

392 [Figure supplement 2. Alternative RL model fits.](#)

393 [Figure supplement 3. Simulating effects of stimulation.](#)

394

395 As with human behaviour on the Daw two-step task, the model (Figure 5A, Figure 5 - figure  
396 supplement 1) that best fit our baseline dataset used a mixture of model-based and model-free  
397 control. However, model comparison indicated the existence of a number of further structural  
398 features that have not previously been reported in models used for the Daw two-step task:  
399 forgetting about the values and state transitions for not-chosen actions, action perseveration effects  
400 spanning multiple trials, and representation of actions both at the level of the choice they represent  
401 (e.g. high poke) and the motor action they require (e.g. left  $\rightarrow$  high movement). These are discussed  
402 in detail in the supplementary material. Taken together, the additional features produced a very  
403 substantial improvement in fit quality ( $\Delta$  iBIC = 11018) over the model which lacked them (Figure 5 –  
404 figure supplements 1,2).

405 In seeking to use the model that fit the baseline dataset most parsimoniously to identify what aspect  
406 of learning or control was disrupted by ACC stimulation, we therefore had to understand their  
407 potential disrupting effects on telling apart model-based and model-free behaviour from data. As  
408 we also discuss in the supplementary material, this is a significant concern because either  
409 perseveration or model-free RL occurring at the level of motor actions rather than choices can  
410 generate loading on the transition predictor in the logistic regression (Figure 5 – figure supplement  
411 3), breaking the simple pattern observed in figure 3B whereby only model-based RL gives substantial  
412 loading on the transition predictor.

413 We therefore sought to understand what aspect of learning or control was affected by the ACC  
414 inhibition by fitting a version of the RL model to the stimulation dataset in which parameters were  
415 free to take different values on stimulated and non-stimulated trials. In the JAWS animals (Figure  
416 5B), the weighting parameter for the model-based system, which controls how strongly model-based  
417 action values influence choice, was significantly reduced on stimulation trials ( $P = 0.021$ ,  
418 permutation test). This was not observed in control GFP animals ( $P = 0.348$ ). We also found that  
419 the learning rate for motor-level perseveration was increased in stimulation trials ( $P = 0.01$ ). The  
420 absolute size of the effects were not large, though this is likely influenced by the fitting procedure  
421 we used whereby we fit a version of the model in which parameters were constrained to take the  
422 same value on stimulated and unstimulated trials and then used this fit as the starting conditions for

423 fitting the full model (see methods). Consistent with the logistic regression analyses, bias towards  
424 the high poke was significantly higher in both JAWS and GFP control animals on stimulation trials ( $P$   
425  $< 0.001$ ), which likely reflects a bias caused by the light. The control animals also showed a  
426 significantly higher value for the eligibility trace parameter on stimulated trials ( $P = 0.027$ ).

427 Taken in isolation this model fitting analysis would not be taken as robust support for an effect of  
428 ACC inhibition on model-based control because the effects would not survive multiple comparison  
429 correction for the large number of model parameters. However, we are not using this analysis to  
430 demonstrate the existence of an effect, but rather to test a hypothesis and probe the nature of the  
431 effect found in the regression analysis. Therefore, the lack of multiple comparison correction is  
432 appropriate here. We know that ACC inhibition affected some aspect of learning or control which  
433 causes experiencing a common transition to promote repeating the preceding choice (Figure 4A).  
434 Standard model-free RL does not predict any effect of transition type on choice while model-based  
435 RL does (Figure 3B), however we found that such an influence could also be generated by other  
436 factors, specifically perseveration or model-free RL occurring at the level of motor actions (Figure 5 –  
437 figure supplement 3). The RL analysis of the stimulation data supports the hypothesis that it is  
438 reduced influence of model-based RL on choice that explains the effect observed in the regression  
439 analysis as the weighting parameter for the model-based component was reduced on stimulation  
440 trials. The increased learning rate for motor-level perseveration should if anything increase loading  
441 on the transition predictor and hence could not explain the regression analysis effect. The  
442 probabilistic reversal learning task further argues against the effect of ACC inhibition being on  
443 outcome independent perseveration at the motor-level as in this task ACC inhibition reduced the  
444 influence of the most recent outcome.

## 445 Discussion:

446

447 We developed a novel two-step decision task for rodents that was designed to dissociate model-  
448 based and model-free RL. We used this task to probe the effect on reward guided behaviour of  
449 silencing ACC neurons, finding that optogenetic inhibition on individual trials reduced the influence  
450 of the experienced state transition, but not the trial outcome, on subsequent choice. Analysis using  
451 RL models suggested these effects were due to a disruption of model-based control.

452 The task was adapted from the two-step decision making task developed for human subjects by Daw  
453 and colleagues (Daw et al., 2011). The Daw two-step has been widely adopted because it offers the  
454 possibility of dissociating control strategies during ongoing learning and decision making, and  
455 generates large decision datasets well suited to behavioural modelling, manipulations and

456 neurophysiology. However, in Akam et al. (2015) and here, we identified and addressed a significant  
457 challenge for the presently popular programme of developing versions of this task for animal  
458 subjects – that subjects may develop habitual mappings from where rewards are received to first  
459 step actions (referred to as extended state representations) which can generate behaviour that  
460 closely resembles model-based strategies. This is a particular concern in animal studies due to the  
461 different way subjects learn the task. Human subjects participating in the Daw two-step task are  
462 given detailed information about the structure of the task beforehand such that they start with a  
463 largely correct model, and then perform a limited number (~200) of trials. By contrast, animal  
464 subjects are typically extensively trained to reach the required performance level before recordings  
465 or manipulations are performed, giving ample opportunity to learn alternative strategies. In  
466 humans, extensive training renders apparently model-based behaviour resistant to a cognitive load  
467 manipulation (Economides et al., 2015) which normally disrupts model-based control (Otto et al.,  
468 2013), suggesting that it is possible to develop automatized strategies which closely resemble  
469 planning.

470 Motivated by this concern, we modified the task structure, introducing reversals into the transition  
471 probabilities mapping the first-step actions to the second-step states. This breaks the long term  
472 predictive relationship between where rewards are obtained and which first-step action has higher  
473 value, precluding a habit-like strategy that exploits this simple relationship, but not confounding a  
474 model-based strategy beyond requiring ongoing learning about the current state of the transition  
475 probabilities. The resulting task is quite complex compared with typical rodent decision tasks, and it  
476 is notable that mice are capable not just of learning it, but of doing so in a few weeks with minimal  
477 shaping. A further advantage of introducing reversals in the transition probabilities is that over the  
478 course of a session, the action-state transition probabilities, first-step action-values, and second-step  
479 state values are mutually decorrelated from each other. This should provide rich opportunity for  
480 future work identifying these decision variables in neural activity.

481 Our approach to developing a rodent two-step task contrasts with that taken by Miller et al. (Miller  
482 et al., 2016b) who retained the fixed transition probabilities of the original Daw task. Model-free use  
483 of extended state representations can produce a similar pattern of regression loadings to those  
484 observed by Miller et al., but interpreted by them in model-based terms. Indeed, the rats in the  
485 Miller et al. study showed little or no evidence of classical model-free behaviour leading to their  
486 conclusion that the behaviour is dominated by model-based planning. This might be surprising as  
487 even humans who have been explicitly told the correct structure of the Daw two-step task show an  
488 approximately even mix of model-based and model-free strategies.

489 Using our large baseline dataset, we performed a detailed characterisation of subject's behaviour on  
490 the new task, including an extensive process of RL model comparison. This indicated that subjects  
491 used a mixture of model-based and model-free RL, consistent with human subjects on the Daw two-  
492 step task. The model comparison also revealed a number of unexpected features of the behaviour;  
493 forgetting about value and state transition probabilities for not chosen actions, perseveration effects  
494 spanning multiple trials, and representation of actions both in terms of the choice they represent  
495 and the motor action they require. We are not aware of studies which have yet compared models  
496 including these elements on human two-step task data.

497 In retrospect, given the finding that representations at the motor-level influenced choice behaviour,  
498 the physical implementation of the task we used had a significant shortcoming: The action required  
499 to execute a given first step choice was different depending on the state reached at the second step  
500 on the previous trial. This caused unnecessary ambiguity in interpreting regression loadings in terms  
501 of control strategy and should be remedied in future work with this class of tasks by modifying the  
502 physical layout of the apparatus.

503 As a target for silencing, we chose the cingulate cortex between AP +1 and AP -0.5 (Figure 2 – figure  
504 supplement 2), which a recent cytoarchitectural study classifies as straddling the boundary between  
505 anterior-cingulate regions 24a and 24b and mid-cingulate regions 24a' and 24b' (Vogt and Paxinos,  
506 2014). Although it has not hitherto been studied in the context of distinguishing actions and habits,  
507 there are anatomical, physiological and lesion-based reasons in rodents, monkeys and humans for  
508 considering this particular role for the structure. First, neurons in rat (Sul et al., 2010) and monkey  
509 (Ito et al., 2003; Matsumoto et al., 2003; Kennerley et al., 2011; Cai and Padoa-Schioppa, 2012) ACC  
510 carry information about chosen actions, reward, action values and prediction errors during decision  
511 making tasks. Where reward type (juice flavour) and size were varied independently (Cai and Padoa-  
512 Schioppa, 2012), a subset of ACC neurons encoded the chosen reward type rather than the reward  
513 value, consistent with a role in learning action-state relationships. In a probabilistic decision making  
514 task in which reward probabilities changed in blocks, neuronal representations in rat ACC underwent  
515 abrupt changes when subjects detected a possible block transition (Karlsson et al., 2012). This  
516 suggests that the ACC may represent the block structure of the task, a form of world model used to  
517 guide action selection, albeit one based on learning about latent states of the world (Gershman and  
518 Niv, 2010; Akam et al., 2015), rather than the forward action-state transition model of classical  
519 model-based RL.

520 Second, neuroimaging in the Daw two-step task has identified representation of model-based value  
521 in the BOLD signal in anterior- and mid-cingulate regions (Daw et al., 2011; Doll et al., 2015).

522 Likewise, neuroimaging in a saccade task in which subjects constructed and updated a model of the  
523 location of target appearance observed ACC activation when subjects updated an internal model of  
524 where saccade targets were likely to appear, (O'Reilly et al., 2013).

525 Third, ACC lesions in macaques produce deficits in tasks which require learning of action-outcome  
526 relationships (Hadland et al., 2003; Kennerley et al., 2006; Rudebeck et al., 2008), though the designs  
527 do not identify whether it is representation of the value or other dimensions of the outcome which  
528 were disrupted. Lesions of rodent ACC produce selective deficits in cost benefit decision making  
529 where subjects must weigh up effort against reward size (Walton et al., 2003; Rudebeck et al.,  
530 2006); however, again, the associative structures concerned are not clear.

531 Finally, the ACC provides a massive innervation to the posterior dorsomedial striatum (Oh et al.,  
532 2014; Hintiryan et al., 2016), a region necessary for learning and expression of goal directed action  
533 as assessed by outcome devaluation (Yin et al., 2005a, 2005b; Hilario et al., 2012).

534 We duly found that silencing ACC neurons on individual trials produced a selective change in how  
535 the previous trials events affected choice, reducing the influence of the previous state transition,  
536 while sparing the influence of reward. This appeared to be due reduced influence of model-based  
537 control on stimulated trials.

538 Recent discussion has focussed on whether ACC plays a direct role in decision making by calculating  
539 decision variables such as the expected value of possible courses of action, or a higher level function  
540 of deciding how much computational effort to expend on a decision (Kolling et al., 2016; Shenhav et  
541 al., 2016). Our results do not discriminate between these theories, because a shift in the balance  
542 between model-based and model-free control could occur either due to directly disrupting the  
543 model-based controller, or disrupting a higher-level system which arbitrated between their usage.

544 In sum, we suggest that our study offers a pioneering example of both the prospects and perils for  
545 the development of a new class of behavioural neuroscience investigations. We showed that it is  
546 possible to fashion sophisticated behavioural tasks that even mice can acquire quickly and  
547 effectively, thus affording all the benefits of modern genetic tools. However, in doing so, we showed  
548 the necessity for examining the behaviour in painstaking detail, lest one be misled by surface  
549 characteristics. We then provided suitably qualified support for the involvement of a key region of  
550 the brain in a cognitive trade-off of great contemporary interest. Our methods should offer rich  
551 opportunities for addressing this and other questions concerning the implementation and  
552 interaction of different neural control systems.

553

## 554 **Methods:**

555

556 *Animals.* All procedures were reviewed and performed in accordance with the Champalimaud Centre  
557 for the Unknown Ethics Committee guidelines. 59 male C57BL mice aged between 2 – 3 months at  
558 the start of experiments were used in the study. Mice were housed socially, except for 1 week in  
559 individual housing post-surgery where applicable. Animals were housed under a 12 hours light/dark  
560 cycle with experiments performed during the light cycle. 17 subjects were used in the two-step task  
561 baseline behaviour dataset. 14 subjects (8 JAWS, 6 GFP controls) were used for the two-step task  
562 ACC manipulation only. 14 subjects (8 JAWS, 6 GFP controls) were used for the probabilistic reversal  
563 learning task ACC manipulation only. 14 subjects (8 JAWS, 6 GFP controls) were first trained and  
564 tested on the two-step ACC manipulation, then retrained for a week on the probabilistic reversal  
565 learning task and tested on the ACC manipulation in this task. 7 JAWS-GFP animals were excluded  
566 from the study due to poor or mislocated JAWS expression. In the group that was tested on both  
567 tasks, 1 Jaws and 2 control animals were lost from the study before optogenetic manipulation on the  
568 probabilistic reversal learning task due to failure of the LED implants. The resulting group sizes for  
569 the optogenetic manipulation experiments were as reported in the results section.

## 570 *Behaviour*

571 Mice were placed on water restriction 48 hours before the first behavioural training session, and  
572 given 1 hour *ad libitum* access to water in their home cage 24 hours before the first training session.  
573 Mice received 1 training session per day of duration 1.5 – 2 hours, and were trained 6 days per week  
574 with 1 hour *ad libitum* water access in their home cage on their day off. During behavioural training  
575 mice had access to dry chow in the testing apparatus as we found this increased the number of trials  
576 performed and amount of water consumed. On days when mice were trained they typically  
577 received all their water in the task (typically 0.5-1.25ml), but additional water was provided as  
578 required to maintain a body weight >85% of their pre-restriction weight. Under this protocol,  
579 bodyweight typically dropped to ~90% of pre-restriction level in the first week of training, then  
580 gradually increased over weeks to reach a steady state of ~95-105% pre-restriction body weight  
581 (Figure 2 – figure supplement 3).

582 Behavioural experiments were performed in 14 custom made 12x12cm operant chambers using  
583 pyControl (<http://pycontrol.readthedocs.io/en/latest/>), a behavioural experiment control system  
584 built around the Micropython microcontroller. The pyControl task definition files are included in  
585 supplementary material. The apparatus, trial structure and block structure of the two-step task are



586 described in the results section. Block transitions were triggered based on subject's behaviour,  
587 occurring 20 trials after an exponential moving average ( $\tau = 8$  trials) of subject's choices crossed a  
588 75% correct threshold. The 20 trial delay between the threshold crossing and block transition  
589 allowed subjects performance at the end of blocks to be assessed without selection bias due to the  
590 block transition rule. In neutral blocks where there was no correct choice, block transitions occurred  
591 with 0.1 probability on each trial after the 40<sup>th</sup>, giving a mean neutral block length of 50 trials.  
592 Subjects started each session with the reward and transition probabilities in the same state that the  
593 previous session finished on.

594 Subjects encountered the full trial structure from the first day of training. The only task parameters  
595 that were changed over the course of training were the reward and state transition probabilities and  
596 the reward sizes. These were changed to gradually increase task difficulty over days of training, with  
597 the typical trajectory of parameter changes as follows:

Session number	Reward size (ul)	Transition probabilities (common / rare)	Reward probabilities (good / bad side)
1	10	0.9 / 0.1	First 40 trials all rewarded, subsequently 0.9 / 0.1
2 - 4	10	0.9 / 0.1	0.9 / 0.1
5 - 6	6.5	0.9 / 0.1	0.9 / 0.1
7 - 8	4	0.9 / 0.1	0.9 / 0.1
9 - 12	4	0.8 / 0.2	0.9 / 0.1
13+	4	0.8 / 0.2	0.8 / 0.2

598

599 The trials structure and block structure of the probabilistic reversal learning task are described in the  
600 results section. Block transitions from non-neutral blocks were triggered 10 trials after an  
601 exponential moving average ( $\tau = 8$  trials) crossed a 75% correct threshold. Block transitions from  
602 neutral blocks occurred with probability 0.1 on each trial after the 15<sup>th</sup> of the block to give an  
603 average neutral block length of 25 trials.

#### 604 *Optogenetic Inhibition*

605 Experimental animals were injected bilaterally with *AAV5-CamKII-Jaws-KGC-GFP-ER2* (UNC vector  
606 core, titre:  $5.9 \times 10^{12}$ ) using 16 injections each of 50nL (total 800nL) spread across 4 injection tracks  
607 (2 per hemisphere) at coordinates: AP: 0, 0.5, ML:  $\pm 0.4$ , DV: -1, -1.2, -1.4, -1.6mm relative to dura.  
608 Control animals were injected with *AAV5-CaMKII-GFP* (UNC vector core, titre:  $2.9 \times 10^{12}$ ) at the same

609 coordinates. Injections were performed at a rate of 4.6nL/5 seconds, using a Nanojet II (Drummond  
610 Scientific) with bevelled glass micropipettes of tip diameter 50-100um. A circular craniotomy of  
611 diameter 1.8mm was centred on AP: 0.25, ML: 0, and a high power red led (Cree XLamp XP-E2) was  
612 positioned above the craniotomy touching the dura. The LED was mounted on a custom designed  
613 insulated metal substrate PCB (Figure 1 – figure supplement 1A). The LEDs were powered using a  
614 custom designed constant current LED driver built around the AL8805 integrated circuit. Light  
615 stimulation (50mW, 630nm) was delivered on stimulation trials from when the subject entered the  
616 side poke until the subsequent choice, up to a maximum of 6 seconds. Stimulation was delivered on  
617 a randomly selected 17% of trials, with a minimum of 2 non-stimulated trials between each  
618 stimulation trial followed by a 0.25 probability of stimulation on each subsequent trial. At the end of  
619 behavioural experiments, animals were sacrificed and perfused with paraformaldehyde (4%). The  
620 brains were sectioned in 50um coronal slices and the location of viral expression was characterised  
621 with fluorescence microscopy (Figure 1 – figure supplement 2).

622 Two animals were injected unilaterally with the JAWS-GFP virus using the coordinates described  
623 above and implanted with the LED implant and a movable bundle of 16 tungsten micro-wires of  
624 23µm diameter (Innovative-Neurophysiology) to record unit activity. After 4 weeks of recovery,  
625 recording sessions were performed at 24 hour intervals and the electrode bundle was advanced by  
626 50 um after each session, covering a depth range of 300 – 1300um from dura over the course of  
627 recordings. During recording sessions mice were free to move inside a sound attenuating chamber.  
628 Light pulses (50mW power, 5 second duration) were delivered at random intervals with a mean  
629 inter-stimulus interval of 30 seconds. Neural activity was acquired using a Plexon recording system  
630 running Omniplex v. 1.11.3. The signals were digitally recorded at 40000 Hz and subsequently band-  
631 pass filtered between 200 Hz and 3000 Hz. Following filtering, spikes were detected using an  
632 amplitude threshold set at twice the standard deviation of the bandpass filtered signal. Initial  
633 sorting was performed automatically using Kilosort (Pachitariu et al., 2016). The results were refined  
634 via manual sorting based on waveform characteristics, PCA and inter-spike interval histogram.  
635 Clusters were classified as single units if well separated from noise and other units and the spike rate  
636 in the 2ms following each spike was less than 1% of the average spike rate.

637 *Behavioural analysis:* All analysis of behaviour was performed in Python, full analysis code and  
638 behavioural data is included in supplementary material.

#### 639 *Logistic regression model*

640 The logistic regression model for the two-step task predicted the probability of choosing the high  
641 poke as a function events on the previous trial using the following set of predictors:

<b>Variables used to define two-step task regression predictors</b>	
<i>C</i>	+1 if previous choice to high poke, -1 if previous choice to low poke
<i>O</i>	+1 if previous trial rewarded, -1 if previous trial not rewarded
<i>T</i>	+1 if previous trial had common transition, -1 if previous trial had rare transition
<i>R</i>	+1 if previous choice to correct (higher reward probability) option, -1 if previous choice to incorrect (lower reward probability) option, 0 if neutral block
<b>Predictors used in two-step task logistic regression</b>	
<i>Bias: high/low</i>	1 for all trials. (Promotes choosing high poke)
<i>Bias: clockwise /counter-clockwise</i>	0.5 if previous trial ended on left side, -0.5 if right side. (Promotes choosing high following trials ending on left, low following trials ending on the right)
<i>Stay</i>	0.5 <i>C</i> (Promotes repeating previous Choice)
<i>Correct</i>	0.5 <i>C R</i> (Promotes repeating correct choices)
<i>Outcome</i>	0.5 <i>C O</i> (Promotes repeating rewarded choices)
<i>Transition</i>	0.5 <i>C T</i> (Promotes repeating choices following common transitions)
<i>Transition outcome interaction</i>	0.5 <i>C T O</i> (Promotes repeating choices following rewarded common transitions and non-rewarded rare transitions).

642

643 Note, regression predictors were scaled to take values of  $\pm 0.5$  such that the loading are in units of  
 644 log-odds. The two-step task logistic regression excluded the first 20 trials after each reversal in the  
 645 transition probabilities as it is ambiguous which transitions are common and rare at this point. This  
 646 resulted in ~9% of trials being excluded from the logistic regression analysis.

647 The logistic regression analysis for the probabilistic reversal learning task predicted the probability of  
 648 choosing the left poke as a function of events on the previous 3 trials, using the following set of  
 649 predictors:

<b>Variables used to define probabilistic reversal learning task regression predictors</b>	
$C_{-t}$	1 if left poke chosen on trial $-t$ , -1 if right poke chosen.
$O_{-t}$	1 trial $-t$ rewarded, -1 if trial $-t$ not rewarded.
<b>Predictors used in probabilistic reversal learning task logistic regression</b>	
$Bias$	1 for all trials (Promotes choosing left poke)
$Choice_{-t}$	$0.5 C_{-t}$ for $t \in \{1,2,3\}$ (Promotes repeating choices)
$Outcome_{-t}$	$0.5 C_{-t}O_{-t}$ for $t \in \{1,2,3\}$ (Promotes repeating rewarded choices)

650

651 *Reinforcement learning modelling:*

652 The following variables and parameters were used in the RL models:

<b>RL model variables</b>	
$R$	Reward obtained on trial (0 or 1)
$a_1$	Action taken at first step (high or low poke)
$a_2$	Action taken at second step (left or right poke)
$a'_1$	Action not taken at first step (high or low poke)
$a'_2$	Action not taken at second step (left or right poke)
$m_1$	Motor-level action taken at first step (e.g. left $\rightarrow$ high)
$m'_1$	Motor-level action not taken at first step
$s_1$	First step state (choice state)
$s_2$	Second step state (left-active or right-active)
$s'_2$	State not reached at second step (left-active or right-active)

$Q_{mf}(s, a)$	Model-free action value for action $a$ in state $s$
$Q_{mo}(s_1, m)$	Motor-level model-free action value for motor action $m$ following in state $s_1$
$P(s a)$	Estimated transition probability of reaching state $s$ after taking action $a$
$C(s_1, a)$	Choice perseveration variable
$M(s_1, m)$	Motor perseveration variable
$B(s_1, a_i)$	Choice bias variable
$R(s_1, m_i)$	Rotational bias variable.
<b>RL model parameters</b>	
$\alpha_Q$	Value learning rate
$f_Q$	Value forgetting rate
$\lambda$	Eligibility trace parameter
$\alpha_T$	Transition learning rate
$f_T$	Transition forgetting rate
$\alpha_c$	Learning rate for choice perseveration
$\alpha_m$	Learning rate for motor-level perseveration
$G_{mf}$	Model-free action value weight
$G_{mo}$	Motor-level model free action value weight
$G_{mb}$	Model-based action value weight
$B_c$	Choice bias (high/low)
$B_r$	Rotational bias (clockwise/counter-clockwise)
$P_c$	Choice perseveration strength

$P_m$	Motor-level perseveration strength
-------	------------------------------------

653

654 *RL Model equations:*

655 Model-free RL: The action value update used by the model-free RL component was:

$$656 \quad Q_{mf}(s_1, a_1) \leftarrow (1 - \alpha_Q)Q_{mf}(s_1, a_1) + \alpha_Q \left( Q_{mf}(s_2, a_2) + \lambda (R - Q_{mf}(s_2, a_2)) \right)$$

$$657 \quad Q_{mf}(s_2, a_2) \leftarrow (1 - \alpha_Q)Q_{mf}(s_2, a_2) + \alpha_Q R$$

658 In models that included value forgetting this value of not chosen actions was updated as:

$$659 \quad Q_{mf}(s_1, a'_1) \leftarrow (1 - f_Q)Q_{mf}(s_1, a'_1)$$

$$660 \quad Q_{mf}(s'_2, a'_2) \leftarrow (1 - f_Q) Q_{mf}(s'_2, a'_2)$$

661 Model-based RL: The model-based component updated its estimate of the state transition  
662 probabilities mapping first-step action to second-step state as:

$$663 \quad P(s_2|a_1) \leftarrow (1 - \alpha_T)P(s_2|a_1) + \alpha_T$$

$$664 \quad P(s'_2|a_1) \leftarrow (1 - \alpha_T)P(s'_2|a_1)$$

665 In models that included transition probability forgetting, the state transition probabilities for the not  
666 chosen action decayed towards a uniform distribution as:

$$667 \quad P(s_2|a'_1) \leftarrow (1 - f_T)P(s_2|a'_1) + 0.5f_T$$

$$668 \quad P(s'_2|a'_1) \leftarrow (1 - f_T)P(s'_2|a'_1) + 0.5f_T$$

669 At the start of each trial, model-based first step action values were calculated as:

$$670 \quad Q_{mb}(s_1, a_i) = Q(s_1, a_i) = \sum_j P(s_j|a_i)Q_{mf}(s_j, a_2)$$

671 Motor-level model-free RL: Agents which included motor-level model-free RL learned values for the  
672 first step actions represented as motor movements (e.g. left  $\rightarrow$  high). The motor movement  $m_i$  for a  
673 given choice  $a_i$  (high or low) at the first step is dependent on the second-step state (left or right) the  
674 previous trial ended on. Motor-level model-free action values were updated as:

$$675 \quad Q_{mo}(s_1, m_1) \leftarrow (1 - \alpha_Q)Q_{mo}(s_1, m_1) + \alpha_Q \left( Q_{mf}(s_2, a_2) + \lambda (R - Q_{mf}(s_2, a_2)) \right)$$

676 In models with motor-level model-free RL and value forgetting, all motor-level model-free values  
677 except that of the action taken decayed as:



678  $Q_{mo}(s_1, m'_1) \leftarrow (1 - f_Q)Q_{mo}(s_1, m'_1)$

679 Perseveration: Choice perseveration was modelled using variables  $C(s_1, a)$  which reflected the  
680 previous choice history. In models using a single trial choice kernel these were updated as:

681  $C(s_1, a_1) \leftarrow 0.5$

682  $C(s_1, a'_1) \leftarrow 0$

683 In models which used an exponential choice kernel,  $C(s_1, a)$  were updated as:

684  $C(s_1, a_1) \leftarrow (1 - \alpha_c)C(s_1, a_1) + 0.5 \alpha_c$

685  $C(s_1, a'_1) \leftarrow (1 - \alpha_c)C(s_1, a'_1)$

686 In models which used motor-level perseveration this was modelled using variables  $M(s_1, m)$  which  
687 reflected the previous history of motor actions at the first step. The motor-preservation variable for  
688 the motor action executed was updated as:

689  $M(s_1, m_1) \leftarrow (1 - \alpha_m)M(s_1, a_1) + 0.5 \alpha_m$

690 The motor perseveration variables for all other motor actions was updated as:

691  $M(s_1, a'_1) \leftarrow (1 - \alpha_m)M(s_1, a'_1)$

692 Biases: A bias for the high/low poke was modelled with a bias variable  $B$  which took values:

693  $B(s_1, a_i) = 0.5$  if  $a_i$  is high poke,  $-0.5$  if  $a_i$  is low poke.

694 The rotational bias (see results section) was modelled with a variable  $R(m_i)$  which took values:

695  $R(s_1, m_i) = 0.5$  if  $m_i$  is a clockwise movement (left  $\rightarrow$  high or right  $\rightarrow$  low)

696  $R(s_1, m_i) = -0.5$  if  $m_i$  is a counter-clockwise movement (left  $\rightarrow$  low or right  $\rightarrow$  high)

697 Combined action values: Model-free, motor-level model-free and model-based action values were  
698 combined with perseveration and bias terms to give the net action values that drove choice  
699 behaviour.

700  $Q_{net}(s_1, a_i) = G_{mf}Q_{mf}(s_1, a_i) + G_{mo}Q_{mo}(s_1, m_i) + G_{mb}Q_{mb}(s_1, a_i) + P_c C(s_1, a_i) + P_m M(s_1, m_i)$   
701  $+ B_c B(s_1, a_i) + B_r R(s_1, m_i)$

702 Where  $G_{mf}$ ,  $G_{mo}$  and  $G_{mb}$  are weights controlling the influence of respectively the model-free,  
703 motor-level model-free and model-based action values,  $P_c$  &  $P_m$  control the strength of choice- and  
704 motor-level perseveration, and  $B_c$  &  $B_r$  control the strength of choice and rotational biases,  $m_i$  is

705 that motor action which equates to choice  $a_i$  given the second step state reached on the previous  
 706 trial.

707 Given the net action values for the two first step actions, choice probability was given by the softmax  
 708 decision rule:

709 Probability of choosing action  $a_i = \frac{e^{Q_{net}(s_1, a_i)}}{\sum_j e^{Q_{net}(s_1, a_j)}}$

710 *Hierarchical modelling:*

711 Both the logistic regression analyses and reinforcement learning model fitting used a Bayesian  
 712 hierarchical modelling framework (Huys et al., 2011), in which parameter vectors  $\mathbf{h}_i$  for individual  
 713 sessions were assumed to be drawn from Gaussian distributions at the population level with means  
 714 and variance  $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$ . The population level prior distributions were set to their maximum  
 715 likelihood estimate:

716  $\boldsymbol{\theta}^{ML} = \operatorname{argmax}_{\boldsymbol{\theta}} \{p(D|\boldsymbol{\theta})\}$

717 
$$= \operatorname{argmax}_{\boldsymbol{\theta}} \left\{ \prod_i^N \int d\mathbf{h}_i p(D_i|\mathbf{h}_i) p(\mathbf{h}_i|\boldsymbol{\theta}) \right\}$$

718 Optimisation was performed using the Expectation-Maximisation algorithm with a Laplace  
 719 approximation for the E-step at the k-th iteration given by:

720  $p(\mathbf{h}_i^k | D_i) = N(\mathbf{m}_i^k, \mathbf{V}_i^k)$

721  $\mathbf{m}_i^k = \operatorname{argmax}_{\mathbf{h}} \{p(D_i|\mathbf{h})p(\mathbf{h}|\boldsymbol{\theta}^{k-1})\}$

722 Where  $N(\mathbf{m}_i^k, \mathbf{V}_i^k)$  is a normal distribution with mean  $\mathbf{m}_i^k$  given by the maximum a posteriori value  
 723 of the session parameter vector  $\mathbf{h}_i$  given the population level means and variance  $\boldsymbol{\theta}^{k-1}$ , and the  
 724 covariance  $\mathbf{V}_i^k$  given by the inverse Hessian of the likelihood around  $\mathbf{m}_i^k$ . For simplicity we assumed  
 725 that the population level covariance  $\boldsymbol{\Sigma}$  had zero off-diagonal terms. For the k-th M-step of the EM  
 726 algorithm the population level prior distribution parameters  $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$  are updated as:

727 
$$\boldsymbol{\mu}^k = \frac{1}{N} \sum_{i=1}^N \mathbf{m}_i^k$$

728 
$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{i=1}^N [(\mathbf{m}_i^k)^2 + \mathbf{V}_i^k] - (\boldsymbol{\mu}^k)^2$$

729 Parameters were transformed before inference to enforce constraints ( $0 < \{G_{mf}, G_{mo}, G_{mb}\}$ ,  $0 <$   
730  $\{\alpha_Q, f_Q, \lambda, \alpha_T, f_T, \alpha_c, \alpha_m\} < 1$ ).

731 To avoid local minima reinforcement learning models fits were repeated 16 times with the means of  
732 the population level prior distributions initialised to random values, the repeat with the best  
733 likelihood was then used.

734 *Model comparison:*

735 To compare the goodness of fit for models with different numbers of parameters we used the  
736 integrated Bayes Information Criterion (iBIC) score. The iBIC score is related to the model log  
737 likelihood  $p(D|M)$  as:

$$\begin{aligned} 738 \log p(D|M) &= \int d\boldsymbol{\theta} p(D|\boldsymbol{\theta})p(\boldsymbol{\theta}|M) \\ 739 &\approx -\frac{1}{2}iBIC = \log p(D|\boldsymbol{\theta}^{ML}) - \frac{1}{2}|M|\log |D| \end{aligned}$$

740 Where  $|M|$  is the number of fitted parameters of the prior,  $|D|$  is the number of data points (total  
741 choices made by all subjects) and iBIC is the integrated BIC score. The log data likelihood given  
742 maximum likelihood parameters for the prior  $\log p(D|\boldsymbol{\theta}^{ML})$  is calculated by integrating out the  
743 individual session parameters:

$$\begin{aligned} 744 \log p(D|\boldsymbol{\theta}^{ML}) &= \sum_i^N \log \int d\mathbf{h} p(D_i|\mathbf{h})p(\mathbf{h}|\boldsymbol{\theta}^{ML}) \\ 745 &\approx \sum_i^N \log \frac{1}{K} \sum_{j=1}^K p(D_i|\mathbf{h}^j) \end{aligned}$$

746 Where the integral is approximated as the average over  $K$  samples drawn from the prior  $p(\mathbf{h}|\boldsymbol{\theta}^{ML})$ .  
747 Bootstrap 95% confidence intervals were estimated for the iBIC scores by resampling from the  
748 population of samples drawn from the prior.

749 *Permutation testing:*

750 Permutation testing was used to assess the significance of differences in model fits between  
751 stimulated and non-stimulated trials. For the logistic regression analyses, the regression model was  
752 fit separately to stimulated and non-stimulated trials to give two sets of population level parameters  
753  $\boldsymbol{\theta}_s = \{\boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s\}$  and  $\boldsymbol{\theta}_n = \{\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n\}$ , where  $\boldsymbol{\theta}_s$  are the parameters for the stimulated trials and  $\boldsymbol{\theta}_n$  are  
754 the parameters for the non-stimulated trials. The distance between the population level means for  
755 the stimulated and non-stimulated conditions were calculated as:

756 
$$\Delta_{true} = |\mu_s - \mu_n|$$

757 An ensemble of L permuted datasets was then created by shuffling the labels on trials such that  
758 trials were randomly assigned to the ‘stimulated’ and ‘non-stimulated’ conditions. The model was fit  
759 separately to the stimulated and non-stimulated trials for each permuted dataset and the distance  
760 between population level means in the stimulated and non-stimulated conditions was calculated for  
761 each permuted dataset i as:

762 
$$\Delta_{perm}^i = |\mu_s^i - \mu_n^i|$$

763 The distribution of distances  $\Delta_{perm}$  over the population of permuted datasets approximates the  
764 distribution of distances under the null hypothesis that stimulation does not affect the model  
765 parameters. The P-values for the observed distances  $\Delta_{true}$  are then given by:

766 
$$P = \frac{1}{L} \sum_{i=1}^L x^i$$

767 where  $x^i = 1$  for  $\Delta_{perm}^i \geq \Delta_{true}$ ,  $x^i = 0$  for  $\Delta_{perm}^i < \Delta_{true}$

768 In addition to testing for a significant main effect of the stimulation we tested for significant  
769 stimulation by group interaction. We first evaluated the true distance between the effect sizes for  
770 the two groups as:

771 
$$\Delta_{true} = |(\mu_s^{JAWS} - \mu_n^{JAWS}) - (\mu_s^{GFP} - \mu_n^{GFP})|$$

772 The approximate distribution of this distance under the null hypothesis that there was no difference  
773 between the groups was evaluated by creating an ensemble of permuted datasets in which we  
774 randomly assigned subjects to the JAWS and GFP groups and the interaction P value was calculated  
775 as above.

776 For reinforcement learning models, the model cannot be fitted separately to stimulated and non-  
777 stimulated trials because of the serial dependence of decision variables from trial to trial. We  
778 therefore created RL models where all or a subset of the model parameters took separate values on  
779 stimulated and non-stimulated trials, such that if the base model had n parameters the resulting  
780 model had 2n parameters. To test for significant differences between parameters on stimulated and  
781 non-stimulated trials, the model was fit to give a set of population level parameters  $\theta = \{\mu, \Sigma\}$ , of  
782 which a subset  $\mu_s, \Sigma_s$  were active on stimulation trials and their counterparts  $\mu_n, \Sigma_n$  were active on  
783 non-stimulation trials. As before the distances between the stimulated and non-stimulated  
784 parameter values were calculated as  $\Delta_{true} = |\mu_s - \mu_n|$  and permutation testing otherwise proceeded  
785 as described above for the regression models.

786 The following procedure was used to minimise problems with local minima when these high  
787 parameter count RL models were fitted to stimulation data. We first fitted a version of the model in  
788 which the parameters were the same for stimulated and non-stimulated trials. This fit was repeated  
789 16 times with randomised initial values for the population level prior means. The fit with the best  
790 likelihood across repeats was used to initialise the population level prior distribution for the full  
791 model in which parameters were free to take different values on stimulated and non-stimulated  
792 trials, such the stim and non-stim parameters started the fitting procedure with the same values.  
793 For permutation testing the same initial fit was used for the true and permuted datasets. To ensure  
794 that permutation test results were not dependent on the specific initial fit found, the whole  
795 procedure was repeated 20 times and the mean P value across the 20 repeats was taken.  
796 Permutation tests were run on the Oxford Advanced Research Computing (ARC) facility.

797 *Bootstrap test for reversal analysis:*

798 The speed of behavioural adaptation to reversals in the transition and reward probabilities was  
799 evaluated by fitting exponentials to the average choice probability trajectories following each type of  
800 reversal (Figure 1E). To test whether adaptation following reversals in transition probabilities was  
801 significantly faster than that following reversals in reward probabilities, we constructed a bootstrap  
802 confidence interval for the difference  $\Delta_{\tau} = \tau_R - \tau_T$ , where  $\tau_R$  and  $\tau_T$  are respectively the  
803 exponential time constants following reversals in the reward and transition probabilities. The  
804 bootstrap confidence interval was evaluated by creating an ensemble of L resampled datasets by  
805 drawing subjects with replacement from the set of subjects that comprised the baseline dataset.  
806 The bootstrap P-value was then evaluated as:

807 
$$P = \frac{1}{L} \sum_{i=1}^L x^i$$

808 where  $x^i = 1$  for  $\Delta_{\tau} < 0$ ,  $x^i = 0$  for  $\Delta_{\tau} \geq 0$ .

809 Logistic regressions of simulated data:

810 To evaluate the logistic regression loadings expected for a model-based and model-free agent on the  
811 task (Figure 2B), we first fitted each agent type to our baseline behavioural dataset using the  
812 hierarchical framework outlined above. The agents used were a model-free agent with eligibility  
813 traces and value forgetting, and a model-based agent with value and transition probability  
814 forgetting. We then simulated data (4000 sessions each of 500 trials) from each agent, drawing  
815 parameters for each session from the fitted population level distributions for that agent. We

816 performed the logistic regression on the simulated data, again using the hierarchical framework as  
817 for the logistic regression analysis of experimental data.

### 818 *Simulating effects of single trial inhibition*

819 In Figure 5 – figure supplement 3 we simulated the effects of lesioning on ‘stimulation’ trials  
820 individual components of that RL model found to give the best fit to the baseline dataset. This was  
821 done by setting the weighting parameter for the relevant component to zero on stimulation trials,  
822 removing its influence on choice on that trial. The components lesioned and their respective  
823 weighting parameters were; choice-level model-free RL ( $G_{mf}$ ), motor-level model-free RL ( $G_{mo}$ ),  
824 model-based RL ( $G_{mb}$ ), motor-level perseveration ( $P_m$ ). For each lesion simulation, a simulated  
825 dataset (4000 sessions each of 500 trials) was generated using parameters for each session drawn  
826 from the population level distribution of the model fit to the baseline dataset. The logistic  
827 regression analysis of the simulated data was performed as on the experimental data by fitting the  
828 regression model separately to choices made on stimulated and non-stimulated trials.

829

### 830 Acknowledgements:

831

832 The authors thank Zach Mainen, Joe Patton, Mark Walton, Tim Behrens, Nathaniel Daw, Kevin Miller  
833 and Bruno Miranda for discussions about the work. The authors acknowledge the use of the  
834 University of Oxford Advanced Research Computing (ARC) facility  
835 (<http://dx.doi.org/10.5281/zenodo.22558>).

836

### 837 Competing interests:

838

839 The authors have no competing interests to report.

840



841 **References:**

842

843 Adams, C.D., and Dickinson, A. (1981). Instrumental responding following reinforcer devaluation. *Q.*  
844 *J. Exp. Psychol. Sect. B* *33*, 109–121.

845 Akaishi, R., Umeda, K., Nagase, A., and Sakai, K. (2014). Autonomous Mechanism of Internal Choice  
846 Estimate Underlies Decision Inertia. *Neuron* *81*, 195–206.

847 Akam, T., Costa, R., and Dayan, P. (2015). Simple Plans or Sophisticated Habits? State, Transition and  
848 Learning Interactions in the Two-Step Task. *PLoS Comput Biol* *11*, e1004648.

849 Balleine, B.W., and Dickinson, A. (1998). Goal-directed instrumental action: contingency and  
850 incentive learning and their cortical substrates. *Neuropharmacology* *37*, 407–419.

851 Balleine, B.W., Killcross, A.S., and Dickinson, A. (2003). The Effect of Lesions of the Basolateral  
852 Amygdala on Instrumental Conditioning. *J. Neurosci.* *23*, 666–675.

853 Cai, X., and Padoa-Schioppa, C. (2012). Neuronal encoding of subjective value in dorsal and ventral  
854 anterior cingulate cortex. *J. Neurosci.* *32*, 3791–3808.

855 Chuong, A.S., Miri, M.L., Busskamp, V., Matthews, G.A.C., Acker, L.C., Sørensen, A.T., Young, A.,  
856 Klapoetke, N.C., Henninger, M.A., Kodandaramaiah, S.B., et al. (2014). Noninvasive optical inhibition  
857 with a red-shifted microbial rhodopsin. *Nat. Neurosci.* *17*, 1123–1129.

858 Colwill, R.M., and Rescorla, R.A. (1985). Postconditioning devaluation of a reinforcer affects  
859 instrumental responding. *J. Exp. Psychol. Anim. Behav. Process.* *11*, 120–132.

860 Daw, N.D., Niv, Y., and Dayan, P. (2005). Uncertainty-based competition between prefrontal and  
861 dorsolateral striatal systems for behavioral control. *Nat. Neurosci.* *8*, 1704–1711.

862 Daw, N.D., Gershman, S.J., Seymour, B., Dayan, P., and Dolan, R.J. (2011). Model-based influences on  
863 humans' choices and striatal prediction errors. *Neuron* *69*, 1204–1215.

864 Dickinson, A. (1985). Actions and habits: the development of behavioural autonomy. *Philos. Trans. R.*  
865 *Soc. B Biol. Sci.* *308*, 67–78.

866 Dickinson, A., Nicholas, D.J., and Adams, C.D. (1983). The effect of the instrumental training  
867 contingency on susceptibility to reinforcer devaluation. *Q. J. Exp. Psychol.* *35*, 35–51.

868 Dolan, R.J., and Dayan, P. (2013). Goals and Habits in the Brain. *Neuron* *80*, 312–325.

869 Doll, B.B., Duncan, K.D., Simon, D.A., Shohamy, D., and Daw, N.D. (2015). Model-based choices  
870 involve prospective neural activity. *Nat. Neurosci.* *18*, 767–772.

871 Doll, B.B., Bath, K.G., Daw, N.D., and Frank, M.J. (2016). Variability in Dopamine Genes Dissociates  
872 Model-Based and Model-Free Reinforcement Learning. *J. Neurosci.* *36*, 1211–1222.

873 Economides, M., Kurth-Nelson, Z., Lübbert, A., Guitart-Masip, M., and Dolan, R.J. (2015). Model-  
874 Based Reasoning in Humans Becomes Automatic with Training. *PLoS Comput Biol* *11*, e1004463.

875 Gershman, S.J., and Niv, Y. (2010). Learning latent structure: carving nature at its joints. *Curr. Opin.*  
876 *Neurobiol.* *20*, 251–256.

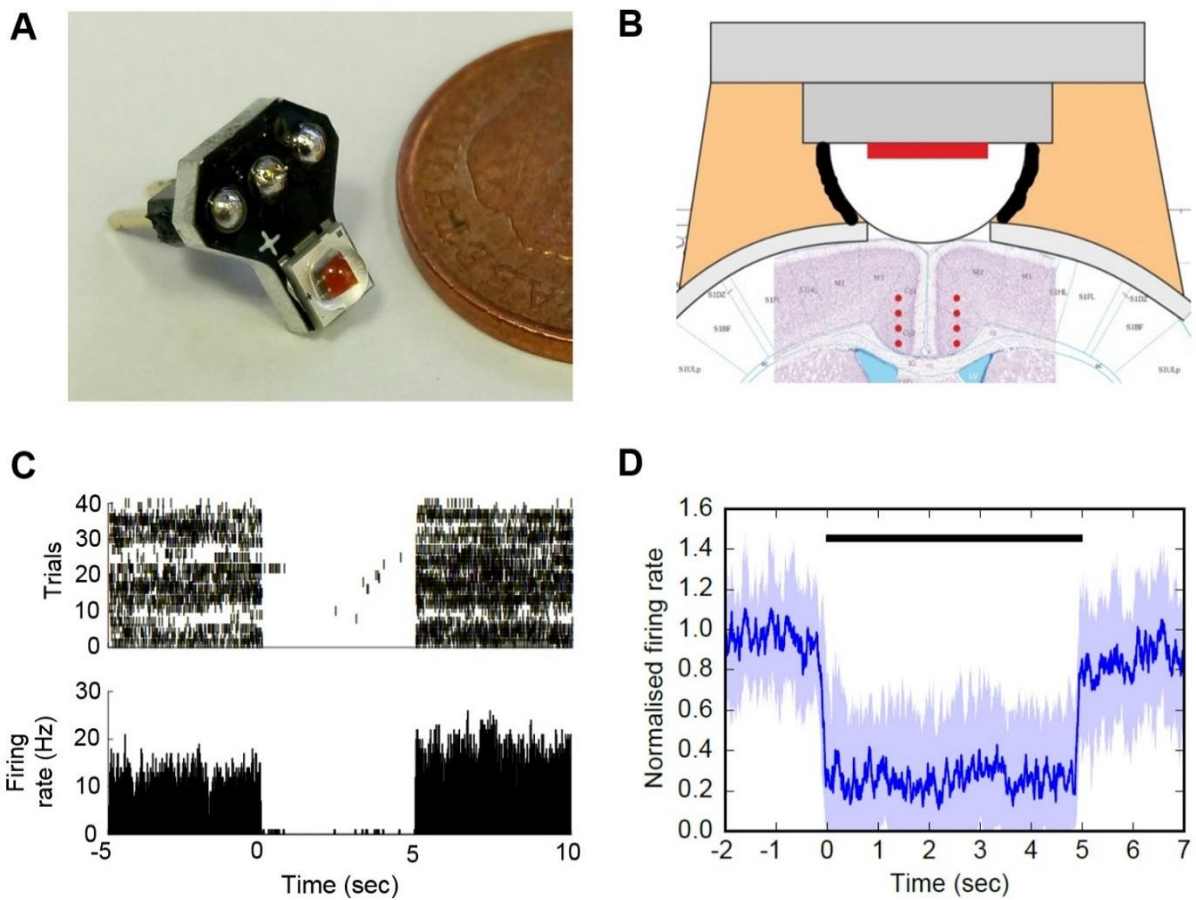
- 877 Gold, J.I., Law, C.-T., Connolly, P., and Bennur, S. (2008). The Relative Influences of Priors and  
878 Sensory Evidence on an Oculomotor Decision Variable During Perceptual Learning. *J. Neurophysiol.*  
879 *100*, 2653–2668.
- 880 Gremel, C.M., and Costa, R.M. (2013a). Premotor cortex is critical for goal-directed actions. *Front.*  
881 *Comput. Neurosci.* *7*.
- 882 Gremel, C.M., and Costa, R.M. (2013b). Orbitofrontal and striatal circuits dynamically encode the  
883 shift between goal-directed and habitual actions. *Nat. Commun.* *4*.
- 884 Hadland, K.A., Rushworth, M.F.S., Gaffan, D., and Passingham, R.E. (2003). The Anterior Cingulate  
885 and Reward-Guided Selection of Actions. *J. Neurophysiol.* *89*, 1161–1164.
- 886 Hilario, M., Holloway, T., Jin, X., and Costa, R.M. (2012). Different dorsal striatum circuits mediate  
887 action discrimination and action generalization. *Eur. J. Neurosci.* *35*, 1105–1114.
- 888 Hintiryan, H., Foster, N.N., Bowman, I., Bay, M., Song, M.Y., Gou, L., Yamashita, S., Bienkowski, M.S.,  
889 Zingg, B., Zhu, M., et al. (2016). The mouse cortico-striatal projectome. *Nat. Neurosci.*
- 890 Huys, Q.J.M., Cools, R., Gölzer, M., Friedel, E., Heinz, A., Dolan, R.J., and Dayan, P. (2011).  
891 Disentangling the Roles of Approach, Activation and Valence in Instrumental and Pavlovian  
892 Responding. *PLoS Comput Biol* *7*, e1002028.
- 893 Huys, Q.J.M., Eshel, N., O’Nions, E., Sheridan, L., Dayan, P., and Roiser, J.P. (2012). Bonsai trees in  
894 your head: how the Pavlovian system sculpts goal-directed choices by pruning decision trees. *PLoS*  
895 *Comput. Biol.* *8*, e1002410.
- 896 Ito, M., and Doya, K. (2009). Validation of decision-making models and analysis of decision variables  
897 in the rat basal ganglia. *J. Neurosci.* *29*, 9861–9874.
- 898 Ito, M., and Doya, K. (2015). Distinct Neural Representation in the Dorsolateral, Dorsomedial, and  
899 Ventral Parts of the Striatum during Fixed- and Free-Choice Tasks. *J. Neurosci.* *35*, 3499–3514.
- 900 Ito, S., Stuphorn, V., Brown, J.W., and Schall, J.D. (2003). Performance Monitoring by the Anterior  
901 Cingulate Cortex During Saccade Countermanding. *Science* *302*, 120–122.
- 902 Karlsson, M.P., Tervo, D.G., and Karpova, A.Y. (2012). Network resets in medial prefrontal cortex  
903 mark the onset of behavioral uncertainty. *Science* *338*, 135–139.
- 904 Kennerley, S.W., Walton, M.E., Behrens, T.E.J., Buckley, M.J., and Rushworth, M.F.S. (2006). Optimal  
905 decision making and the anterior cingulate cortex. *Nat Neurosci* *9*, 940–947.
- 906 Kennerley, S.W., Behrens, T.E., and Wallis, J.D. (2011). Double dissociation of value computations in  
907 orbitofrontal and anterior cingulate neurons. *Nat. Neurosci.* *14*, 1581–1589.
- 908 Killcross, S., and Coutureau, E. (2003). Coordination of Actions and Habits in the Medial Prefrontal  
909 Cortex of Rats. *Cereb. Cortex* *13*, 400–408.
- 910 Kolling, N., Wittmann, M.K., Behrens, T.E.J., Boorman, E.D., Mars, R.B., and Rushworth, M.F.S.  
911 (2016). Value, search, persistence and model updating in anterior cingulate cortex. *Nat. Neurosci.*  
912 *19*, 1280–1285.

- 913 Kool, W., Cushman, F.A., and Gershman, S.J. (2016). When Does Model-Based Control Pay Off? *PLOS*  
914 *Comput Biol* 12, e1005090.
- 915 Matsumoto, K., Suzuki, W., and Tanaka, K. (2003). Neuronal correlates of goal-based motor selection  
916 in the prefrontal cortex. *Science* 301, 229–232.
- 917 Miller, K., Shenhav, A., and Ludvig, E. (2016a). Habits without Values. bioRxiv 067603.
- 918 Miller, K.J., Botvinick, M.M., and Brody, C.D. (2016b). Dorsal hippocampus plays a causal role in  
919 model-based planning. bioRxiv 096594.
- 920 Oh, S.W., Harris, J.A., Ng, L., Winslow, B., Cain, N., Mihalas, S., Wang, Q., Lau, C., Kuan, L., Henry,  
921 A.M., et al. (2014). A mesoscale connectome of the mouse brain. *Nature* 508, 207–214.
- 922 O’Reilly, J.X., Schüffelgen, U., Cuell, S.F., Behrens, T.E., Mars, R.B., and Rushworth, M.F. (2013).  
923 Dissociable effects of surprise and model update in parietal and anterior cingulate cortex. *Proc. Natl.*  
924 *Acad. Sci.* 110, E3660–E3669.
- 925 Ostlund, S.B., and Balleine, B.W. (2005). Lesions of medial prefrontal cortex disrupt the acquisition  
926 but not the expression of goal-directed learning. *J. Neurosci.* 25, 7763.
- 927 Otto, A.R., Gershman, S.J., Markman, A.B., and Daw, N.D. (2013). The Curse of Planning Dissecting  
928 Multiple Reinforcement-Learning Systems by Taxing the Central Executive. *Psychol. Sci.* 24, 751–761.
- 929 Otto, A.R., Skatova, A., Madlon-Kay, S., and Daw, N.D. (2014). Cognitive control predicts use of  
930 model-based reinforcement learning. *J. Cogn. Neurosci.*
- 931 Pachitariu, M., Steinmetz, N., Kadir, S., Carandini, M., and Harris, K.D. (2016). Kilosort: realtime  
932 spike-sorting for extracellular electrophysiology with hundreds of channels. bioRxiv 061481.
- 933 Paxinos, G., and Franklin, K.B. (2007). *The mouse brain in stereotaxic coordinates* -3rd Edition  
934 (Academic Press).
- 935 Rudebeck, P.H., Walton, M.E., Smyth, A.N., Bannerman, D.M., and Rushworth, M.F.. (2006). Separate  
936 neural pathways process different decision costs. *Nat. Neurosci.* 9, 1161–1168.
- 937 Rudebeck, P.H., Behrens, T.E., Kennerley, S.W., Baxter, M.G., Buckley, M.J., Walton, M.E., and  
938 Rushworth, M.F.S. (2008). Frontal Cortex Subregions Play Distinct Roles in Choices between Actions  
939 and Stimuli. *J. Neurosci.* 28, 13775–13785.
- 940 Rushworth, M.F.S., and Behrens, T.E.J. (2008). Choice, uncertainty and value in prefrontal and  
941 cingulate cortex. *Nat. Neurosci.* 11, 389–397.
- 942 Rushworth, M.F.S., Walton, M.E., Kennerley, S.W., and Bannerman, D.M. (2004). Action sets and  
943 decisions in the medial frontal cortex. *Trends Cogn. Sci.* 8, 410–417.
- 944 Sebold, M., Deserno, L., Nebe, S., Schad, D.J., Garbusow, M., Hägele, C., Keller, J., Jünger, E.,  
945 Kathmann, N., Smolka, M., et al. (2014). Model-Based and Model-Free Decisions in Alcohol  
946 Dependence. *Neuropsychobiology* 70, 122–131.
- 947 Shenhav, A., Cohen, J.D., and Botvinick, M.M. (2016). Dorsal anterior cingulate cortex and the value  
948 of control. *Nat. Neurosci.* 19, 1286–1291.

- 949 Simon, D.A., and Daw, N.D. (2011). Neural Correlates of Forward Planning in a Spatial Decision Task  
950 in Humans. *J. Neurosci.* *31*, 5526–5539.
- 951 Smittenaar, P., FitzGerald, T.H.B., Romei, V., Wright, N.D., and Dolan, R.J. (2013). Disruption of  
952 Dorsolateral Prefrontal Cortex Decreases Model-Based in Favor of Model-free Control in Humans.  
953 *Neuron*.
- 954 Sul, J.H., Kim, H., Huh, N., Lee, D., and Jung, M.W. (2010). Distinct roles of rodent orbitofrontal and  
955 medial prefrontal cortex in decision making. *Neuron* *66*, 449–460.
- 956 Sutton, R.S., and Barto, A.G. (1998). Reinforcement learning: An introduction (The MIT press).
- 957 Thorndike, E.L. (1911). *Animal intelligence: Experimental studies*.
- 958 Vogt, B.A., and Paxinos, G. (2014). Cytoarchitecture of mouse and rat cingulate cortex with human  
959 homologues. *Brain Struct. Funct.* *219*, 185–192.
- 960 Voon, V., Derbyshire, K., Rück, C., Irvine, M.A., Worbe, Y., Enander, J., Schreiber, L.R.N., Gillan, C.,  
961 Fineberg, N.A., Sahakian, B.J., et al. (2015). Disorders of compulsivity: a common bias towards  
962 learning habits. *Mol. Psychiatry* *20*, 345–352.
- 963 Walton, M.E., Bannerman, D.M., Alterescu, K., and Rushworth, M.F.. (2003). Functional  
964 specialization within medial frontal cortex of the anterior cingulate for evaluating effort-related  
965 decisions. *J. Neurosci.* *23*, 6475.
- 966 Wunderlich, K., Smittenaar, P., and Dolan, R.J. (2012). Dopamine Enhances Model-Based over  
967 Model-Free Choice Behavior. *Neuron* *75*, 418–424.
- 968 Yin, H.H., Knowlton, B.J., and Balleine, B.W. (2004). Lesions of dorsolateral striatum preserve  
969 outcome expectancy but disrupt habit formation in instrumental learning. *Eur. J. Neurosci.* *19*, 181–  
970 189.
- 971 Yin, H.H., Knowlton, B.J., and Balleine, B.W. (2005a). Blockade of NMDA receptors in the  
972 dorsomedial striatum prevents action–outcome learning in instrumental conditioning. *Eur. J.*  
973 *Neurosci.* *22*, 505–512.
- 974 Yin, H.H., Ostlund, S.B., Knowlton, B.J., and Balleine, B.W. (2005b). The role of the dorsomedial  
975 striatum in instrumental conditioning. *Eur. J. Neurosci.* *22*, 513–523.

976

977 Figure supplements:



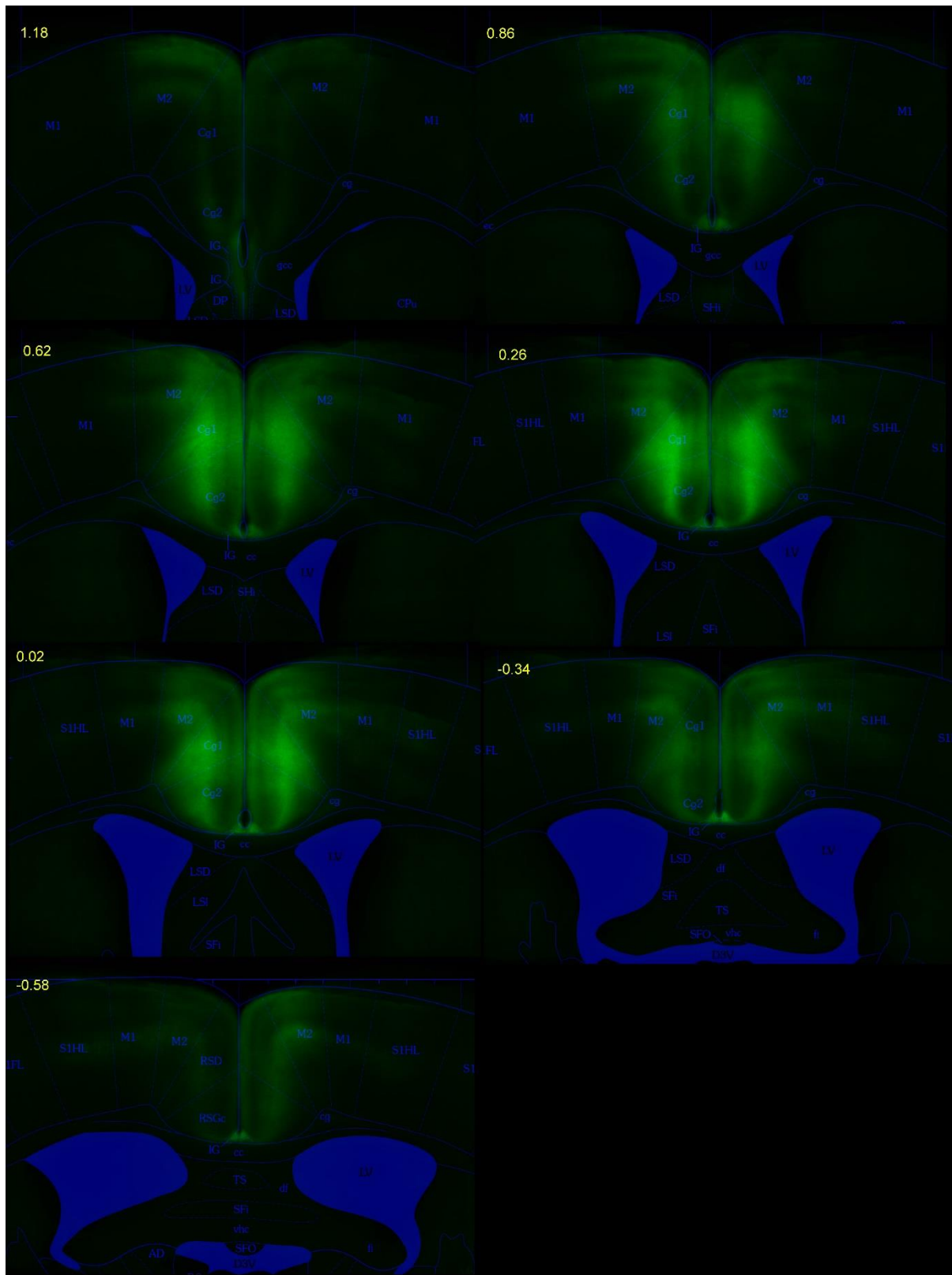
978

979 **Figure 1 - figure supplement 1. JAWS inhibition of ACC neurons.** **A)** LED implant. **B)** Implantation diagram, red dots indicate location of virus injections. **C)** Inhibition of example cell, top panel – spike raster, bottom panel average firing rate. **D)** Normalised firing rate for significantly inhibited cells (Kruskal-Wallis  $P < 0.05$ , 67/249 cells), dark blue line – median, shaded area 25 – 75 percentile.

980

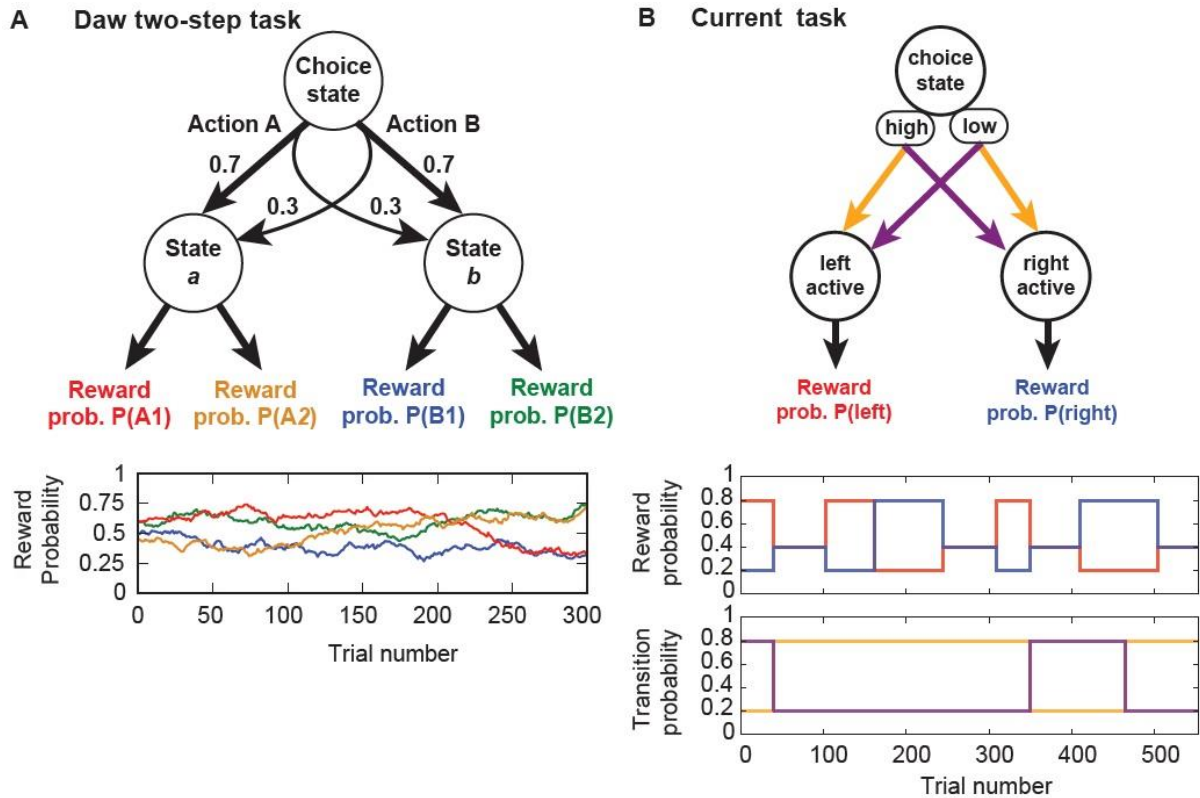
981

982



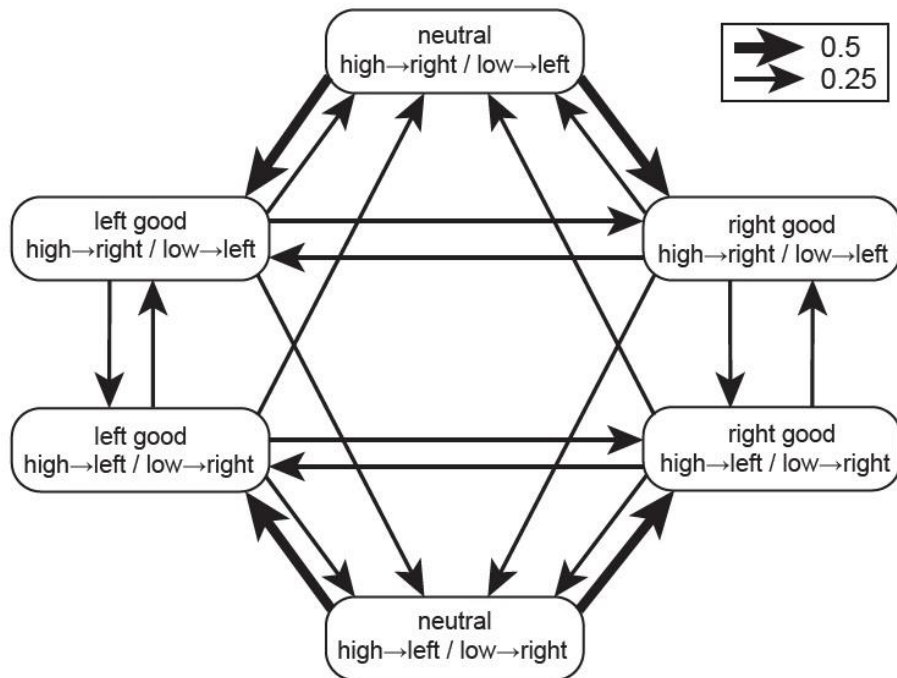
983  
984  
985  
986





987  
988  
989  
990

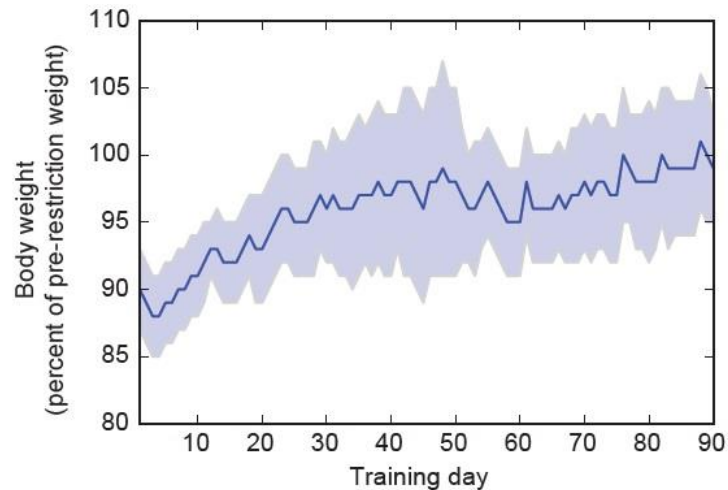
**Figure 2 - figure supplement 1. Comparison of original and new two-step task structures.** A) State diagram of the original Daw two step task with example reward probability trajectories. B) State diagram of the two-step task used in the current study with example reward probability and transition probability trajectories.



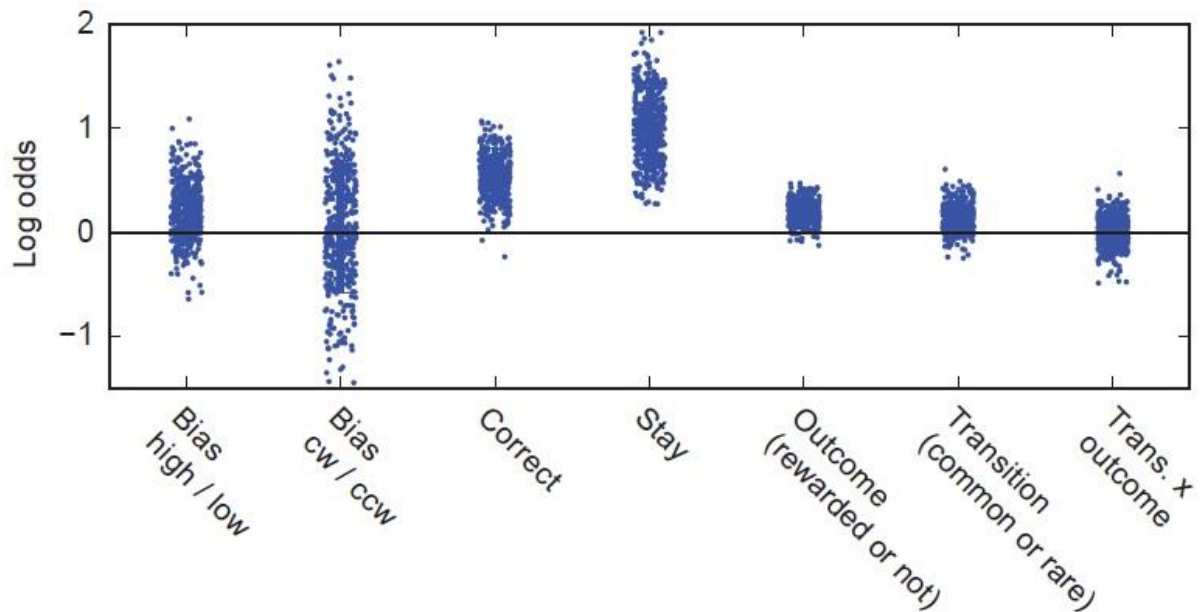
991  
992

**Figure 2 - figure supplement 2. Block transition probabilities.** Diagram of block transition probabilities for the two-step task used in the current study.



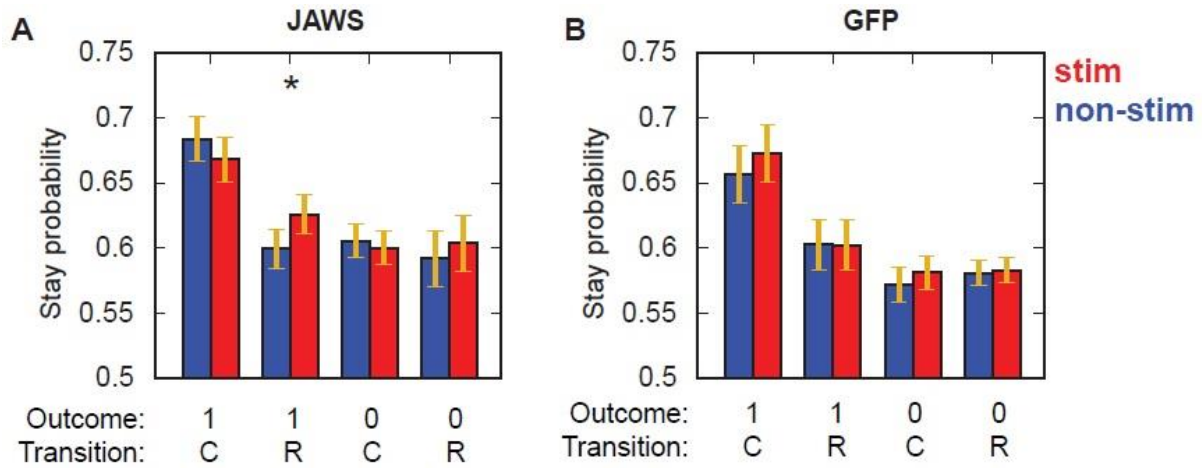


993 **Figure 2 - figure supplement 3. Body weight trajectory across training:** Mean (blue line) and standard-  
994 deviation (shaded area) of subject's body weight trajectory across days of training.



995 **Figure 3 - figure supplement 1. Full logistic regression model fit.** Fit of the logistic regression model to the  
996 baseline dataset showing loadings for all 7 parameters. Bars indicate  $\pm 1$  standard deviation of the population  
997 level distributions, dots indicate maximum a posteriori session fits. Predictors: *Bias high/low* – tendency to  
998 choose the high poke, *bias clockwise /counter-clockwise* – tendency to choose high following left and low  
1000 following right, *Correct* – tendency to choose the correct option, i.e. that option which commonly leads to  
1001 state with higher reward probability, *Stay* – tendency to repeat choices irrespective of subsequent trial  
1002 events, *Outcome* – tendency to repeat choices following reward, *Transition* – tendency to repeat choices  
1003 following common transitions, *Transition-outcome interaction* – tendency to repeat choices following  
1004 rewarded common transition trials and non-rewarded rare transition trials.

1005



1006

1007

1008

1009

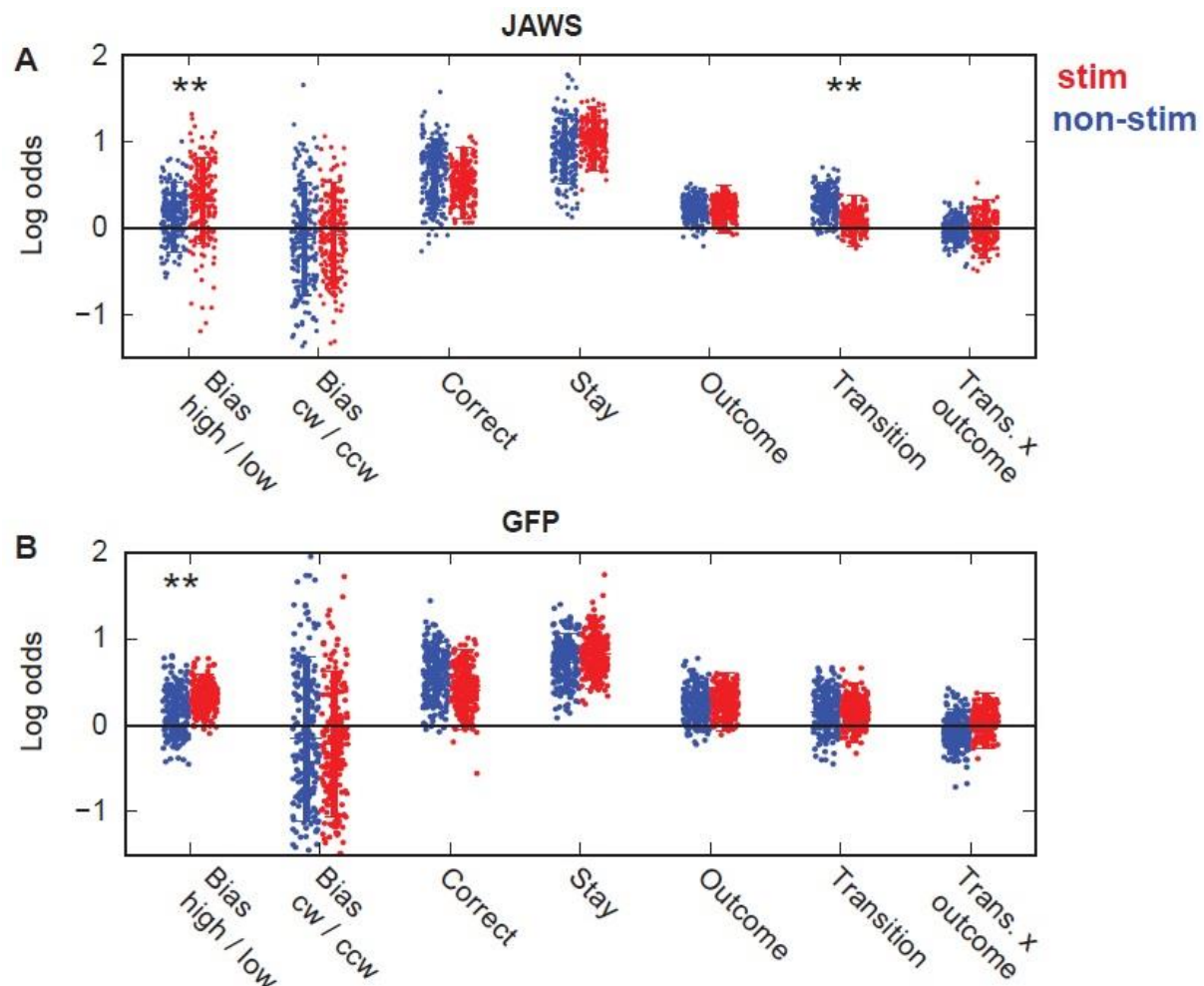
1010

1011

1012

1013

**Figure 4 - figure supplement 1. ACC inhibition stay probabilities** Stay probability analysis for JAWS (A) and GFP control (B) animals showing fraction of trials the subject repeated the same choice following each combination of outcome (rewarded (1) or not (0)) and transition (common (C) or rare (R)). Stay probabilities were evaluated separately for trials with (red) and without (blue) light stimulation delivered from the trial outcome to the subsequent choice. Error bars show cross-subject SEM. \* indicates paired t-test P value < 0.05.



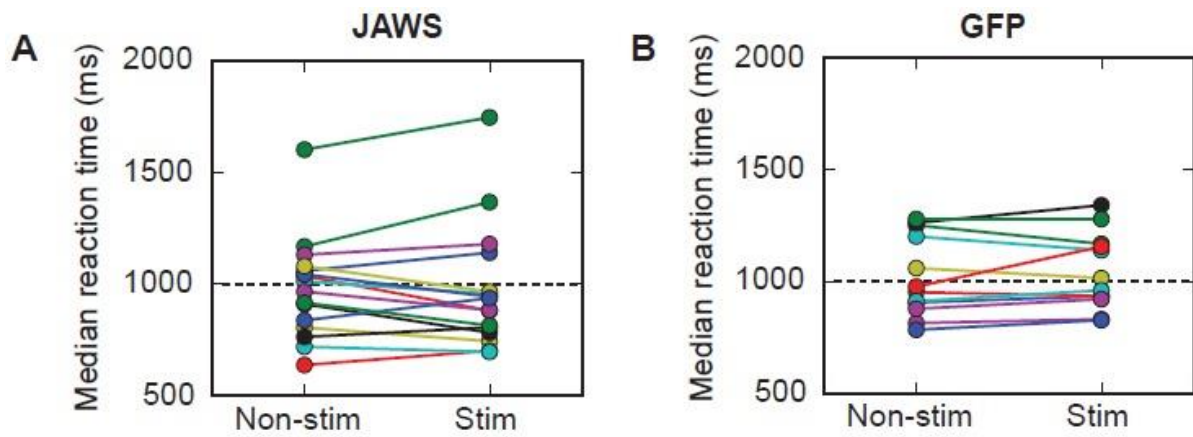
1014

1015

1016

1017

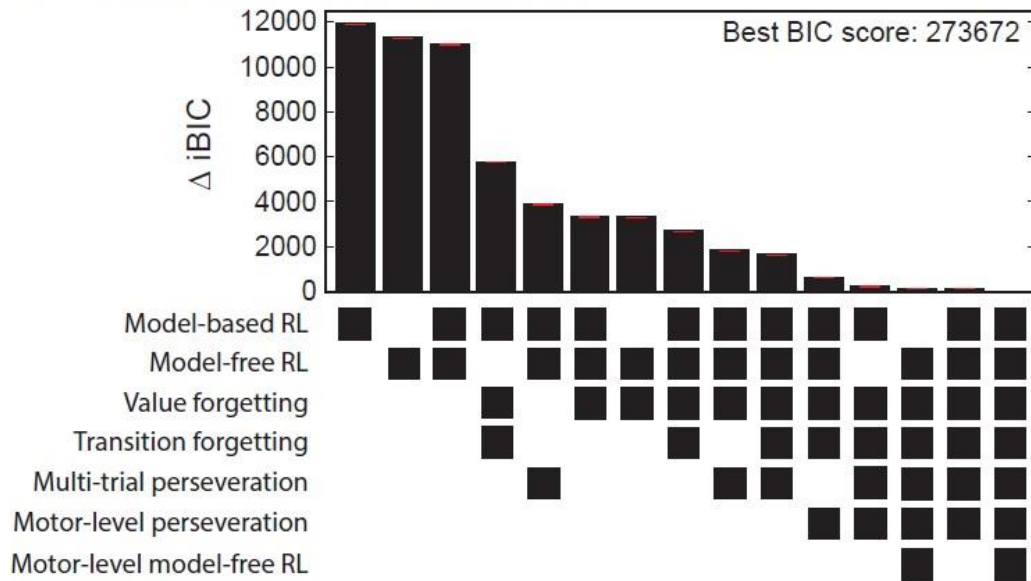
**Figure 4 - figure supplement 2. ACC inhibition full logistic regression model fits.** Fit of the logistic regression model to the JAWS ACC inhibition (A) and GFP controls (B) showing loadings for all 7 parameters. Bars indicate  $\pm 1$  standard deviation of the population level distributions, dots indicate maximum a posteriori session fits.



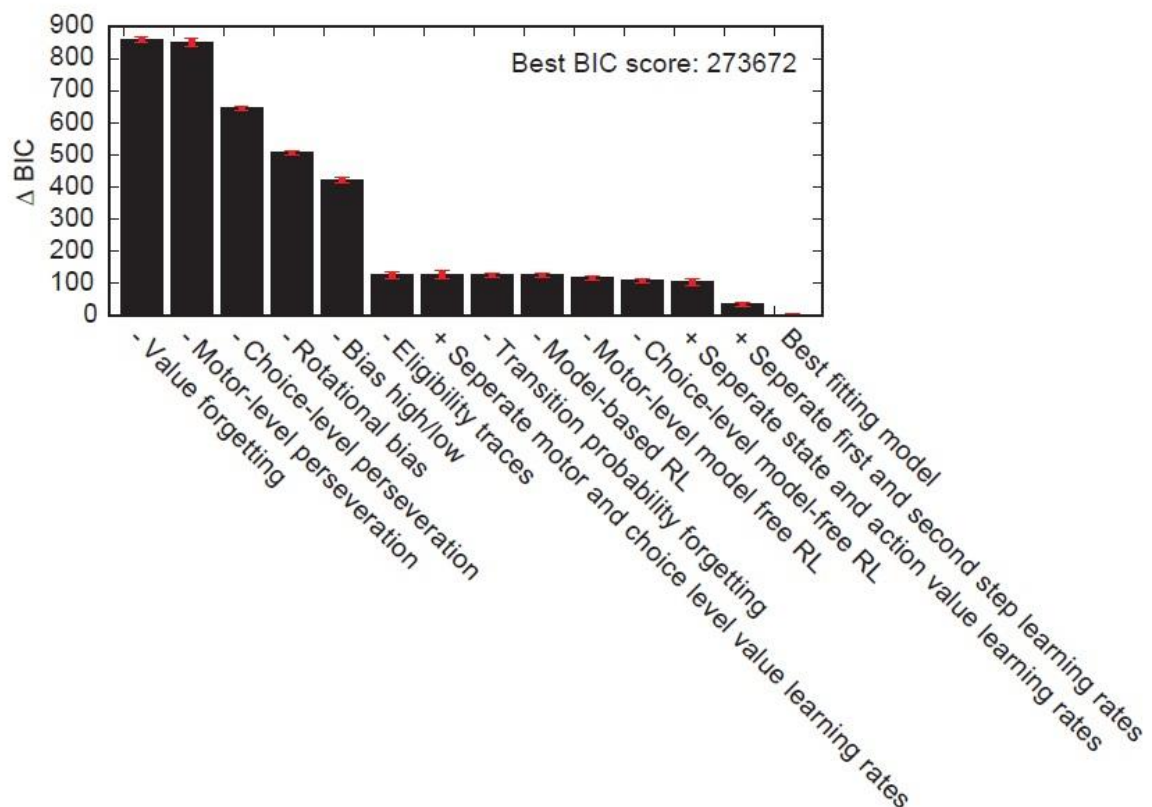
1018  
1019  
1020  
1021  
1022

**Figure 4 - figure supplement 3. ACC inhibition reaction times.** Reaction times for first-step choice on stimulated and non-stimulated trials. Reaction time is measured from the start of the ITI when the subject exits the side poke at the end of the previous trial, until the next high or low poke. The dashed line indicates the end of the ITI at which point the high and low pokes become active.

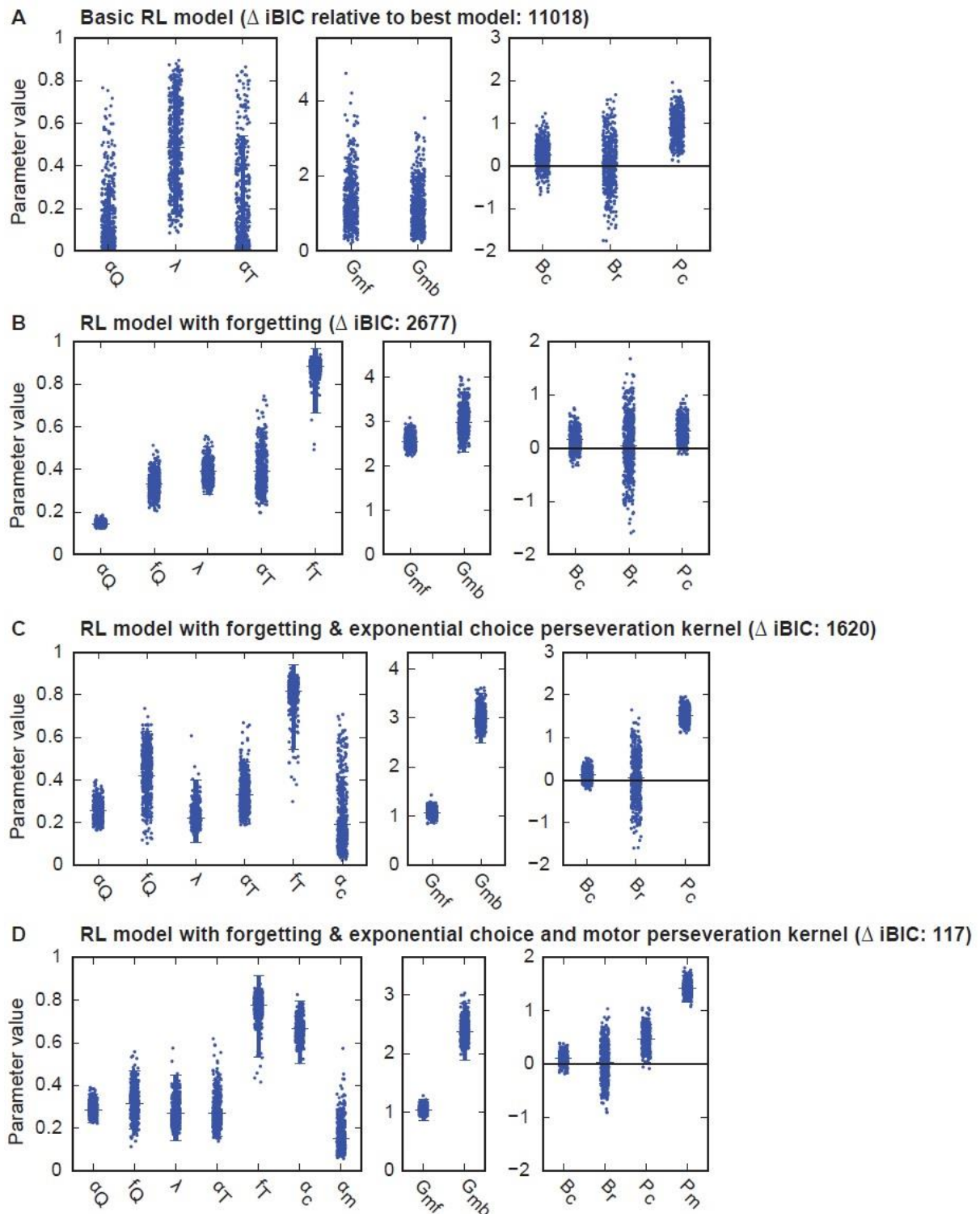
**A** Adding components to basic model.



**B** Adding or removing single components from best model



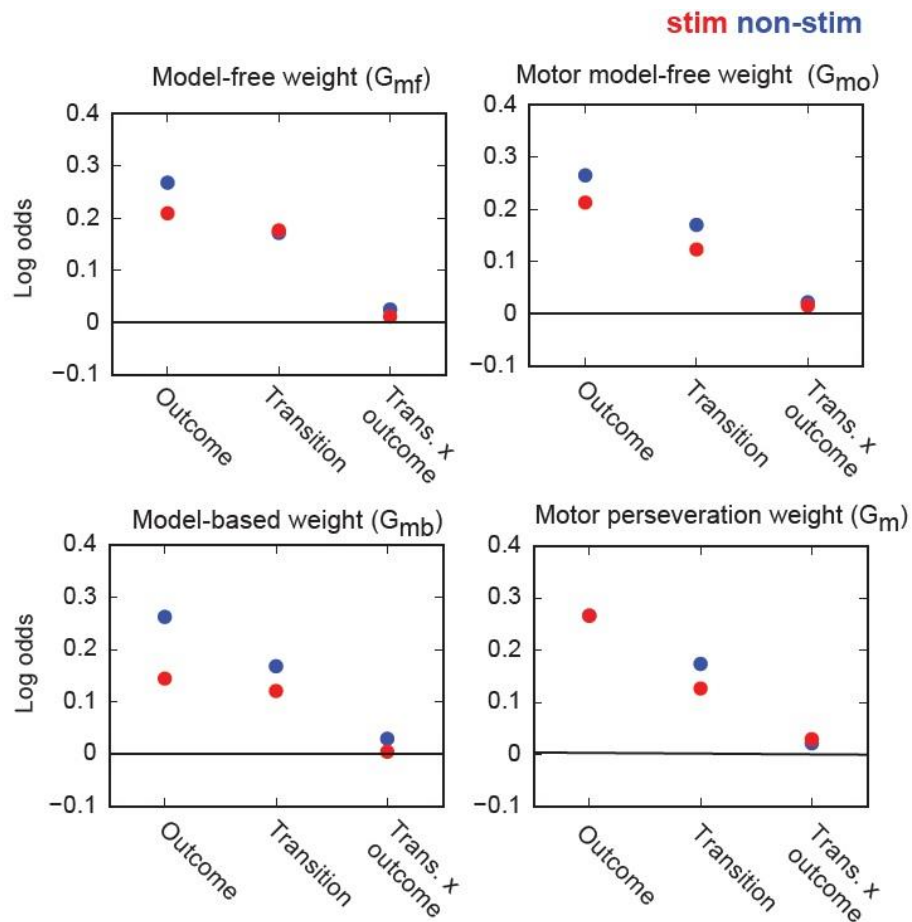
1023 **Figure 5 - figure supplement 1. Baseline dataset BIC score model comparison.** **A)** iBIC score comparison for  
 1024 set of RL models on baseline behavioural dataset. The set of models was constructed as described in  
 1025 supplementary results by iteratively adding features to the RL model. The grid below the plot indicates which  
 1026 features were included in each model. **B)** iBIC score comparison on the baseline dataset for set of RL models  
 1027 created by adding or removing a single feature at a time from the best fitting model. The text below each bar  
 1028 indicates what feature has been added or removed. Error-bars indicate the bootstrap 95% confidence interval  
 1029 on the BIC score.  
 1030



$\alpha_Q$ : Value learning rate,  $f_Q$ : Value forgetting rate,  $\lambda$ : Eligibility trace,  $\alpha_T$ : Transition learning rate,  $f_T$ : Transition forgetting rate,  $\alpha_c$ : Choice perseveration learning rate,  $\alpha_m$ : Motor perseveration learning rate,  $G_{mf}$ : Model-free weight,  $G_{mb}$ : Model-based weight,  $B_c$ : Choice Bias,  $B_r$ : Rotational bias,  $P_c$ : Choice perseveration weight,  $P_m$ : Motor perseveration weight.

1031  
1032 **Figure 5 - figure supplement 2. Alternative RL model fits.** Fit of Reinforcement learning models of different  
1033 levels of complexity. Model complexity increases from **A** to **D** as features are added to the basic RL model. For  
1034 each fit, bars indicate  $\pm 1$  standard deviation of the population level distributions, dots indicate maximum a  
1035 posteriori session fits. For each model the difference in iBIC score between this model and the best fitting  
1036 model is reported.





1037 **Figure 5 - figure supplement 3. Simulating effects of stimulation:** Simulation of the effects of lesioning  
1038 different components of the best fitting RL model on stimulation trials. Model lesioning was implemented by  
1039 setting individual parameters to zero on stimulation trials. Panels show logistic regression loadings for  
1040 stimulated and non-stimulated trials. For each panel the title indicates which model-parameter was set to  
1041 zero on stimulation trials.

1042

1043 Supplementary Material:

1044

1045 *Model comparison:*

1046 The starting point for our model comparison process was the RL agent used in the original Daw two-  
1047 step task (Daw et al., 2011). As the action-state transition probabilities in our task were not fixed,  
1048 we modified the model-based component of the agent to update its estimate of the transition  
1049 probabilities for the chosen action on each trial using an error driven learning rule. As in the original  
1050 Daw agent we included a perseveration parameter which promoted repeating the previous choice.  
1051 Based on the evidence for response biases from the logistic regression, we additionally included in  
1052 the RL agent two parameters capturing a bias towards the high/low poke and the rotational bias  
1053 described in the results section. We compared the goodness of fit of a pure model-free agent, a  
1054 pure model-based agent, and an agent which used a mixture of both strategies. The mixture agent  
1055 provided a better fit to the data than either the pure model-free ( $\Delta$  iBIC = 264, Figure 3B) or pure  
1056 model-based agent ( $\Delta$  iBIC = 888), and the mixture model fit suggested an approximately equal  
1057 contribution of model-based and model-free control (Figure 5 – figure supplement 2A). As the task  
1058 is novel and hence we do not know what features may be present in the behaviour, we performed  
1059 an exploratory process of model comparison to better understand whether the RL model was  
1060 providing a good description of the behaviour. This identified a number of additional features which  
1061 greatly improved fit quality when added to the model.

1062 RL models typically assume that action values of options that are not chosen remain unchanged.  
1063 However, it has been reported that model-fits in some rodent decision making tasks are  
1064 substantially improved by including forgetting about the value of not chosen actions, typically  
1065 implemented as action value decay towards zero (Ito and Doya, 2009, 2015). Including such action  
1066 value forgetting in the mixture agent produced a dramatic improvement in iBIC score for our data ( $\Delta$   
1067 iBIC = 7698). Including forgetting about action-state transition probabilities, implemented as a decay  
1068 of transition probabilities for the not chosen action towards a uniform distribution, further improved  
1069 the goodness of fit ( $\Delta$  iBIC = 643). The mixture agent including value and transition probability  
1070 forgetting again showed approximately equal weighting of the model-based and model-free action  
1071 values in controlling behaviour (Figure 5 – figure supplement 2B). When forgetting was included for  
1072 each agent the mixture agent provided a better fit to the data than either a pure model-free ( $\Delta$  iBIC  
1073 = 612) or pure model-based ( $\Delta$  iBIC = 3066) agent.

1074 Forgetting decreases the value of not chosen relative to chosen options, and therefore promotes  
1075 perseveration of choice. It is therefore possible that if subjects are in fact strongly perseverative,  
1076 this could be mistakenly identified as forgetting in the RL fit. Though the model included a



1077 perseverance parameter for repeating the previous choice, several studies have reported  
1078 perseverance effects spanning multiple trials, even in tasks where decisions optimally should be  
1079 treated as independent (Gold et al., 2008; Akaishi et al., 2014). We therefore tested whether  
1080 goodness of fit was improved by an exponential choice kernel through which prior choices directly  
1081 influenced the current choice with exponentially decreasing weight at increasing lag (Figure 5 –  
1082 figure supplement 2C). This is equivalent to the decision inertia model of Akaishi et al. (2014) in  
1083 which choice is influenced by a variable they term the choice estimate  $CE$ , an average of previous  
1084 choices updated following each decision using the error driven learning rule  $CE_{n+1} = CE_n +$   
1085  $\alpha (C_n - CE_n)$ , where  $C_n$  is the choice on trial  $n$  and  $\alpha$  is a learning rate. The addition of this  
1086 exponential choice kernel dramatically improved fit quality when added to the mixture agent  
1087 without forgetting ( $\Delta$  iBIC = 7133). However even with the exponential choice kernel included, value  
1088 forgetting substantially improved goodness of fit ( $\Delta$  iBIC = 2071), and transition probability forgetting  
1089 further increased goodness of fit ( $\Delta$  iBIC = 194). These results indicate that forgetting about values  
1090 and transitions for not chosen options is a genuine feature of the behaviour and not an artefact due  
1091 to a tendency to perseverate. They further indicate that subjects do in fact show a strong tendency  
1092 to perseverate over multiple trials, which is not captured even by forgetting RL models, presumably  
1093 because it is independent of the recent reinforcement history. Forgetting may be a heuristic used in  
1094 dynamic environments where evidence becomes less reliable with the passage of time due to state  
1095 of the world changing. Alternatively, forgetting may occur due to limitations of the learning systems  
1096 involved, perhaps due to differences between the rapidly changing reward statistics in the task and  
1097 those typical of natural environments.

1098 The choice kernel assumes that perseveration occurs at the level of the decision between the high  
1099 and low pokes, however it is also possible that the perseverative tendency is at the lower level of  
1100 motor actions. In the current task, a given choice (high or low) entails a different motor action  
1101 depending on which side (left or right) the previous trial ended on. We therefore considered a  
1102 model with perseveration at the motor level such that the choice on a given trial only increased the  
1103 probability of repeating that same motor action in future, e.g. a choice taken by moving from the left  
1104 to high poke only increased the probability of choosing high in future following trials which ended on  
1105 the left side (Figure 5 – figure supplement 2D). Motor perseveration was modelled by maintaining  
1106 separate moving averages of choices following trials that ended on the left and right, updated using  
1107 the error driven learning rule described above, which each influenced choices following trials ending  
1108 on their respective sides. Replacing the exponential choice kernel with this motor perseveration  
1109 substantially improved fit quality ( $\Delta$  iBIC = 1004). However, including perseveration both at the  
1110 level of choice, (high vs low, independent of motor action), and at the motor level, further improved

1111 fit quality ( $\Delta$  iBIC = 499), indicating that subjects exhibit perseverative tendencies at both the choice  
1112 and motor level that are not predicted by the RL component of the model. These data support the  
1113 existence of mechanisms which reinforce selected behaviours in a reward-independent fashion, i.e.  
1114 simply choosing to execute a behaviour increases the chance that behaviour will be executed in  
1115 future. This is consistent with previous reports from perceptual (Gold et al., 2008; Akaishi et al.,  
1116 2014) and reward-guided decision making tasks (Miller et al., 2016a), and we think is a parsimonious  
1117 explanation for our results. Such perseveration is somewhat puzzling from a normative perspective  
1118 but may be a signature of a mechanism for automatizing behaviour by reinforcing chosen actions.  
1119 Thorndike proposed such a ‘law of exercise’ (1911) and the idea has recently been revisited by Miller  
1120 et al. (2016a) who suggest that habit formation occurs through outcome-independent reinforcement  
1121 of chosen actions. This framework views habit formation as a supervised learning process in which  
1122 behaviour generated by value sensitive systems, i.e. model-free and model-based RL, is used to train  
1123 value-independent learning systems. Such a mechanism could account for the perseveration  
1124 observed in our data assuming it operated both on actions represented at the level of the choice  
1125 they represent and the level of motor actions. An alternative mechanism which could give rise to  
1126 perseveration would be subjects sampling an option multiple times between choices, which may be  
1127 adaptive if the decision process is costly in time or effort. However, this explanation does not  
1128 account for the observation in our data that perseveration occurred at the level both of choices and  
1129 of motor actions, with different timescales for each (see respective learning rates, Figure 5).

1130 Evidence that perseveration occurred both at the level of choice and motor action raises the  
1131 question of whether reward driven learning also occurs at both levels of representation. This might  
1132 be expected from the architecture of parallel cortical-basal ganglia loops, with circuits linking  
1133 somatosensory and motor cortices to dorsolateral striatum learning values over low level motor  
1134 representations, and circuits linking higher level cortical regions to medial and ventral striatum  
1135 learning values over more abstract state and action representations. We therefore tested an agent  
1136 in which model-free action values were learned in parallel for actions represented both in terms of  
1137 choice (high/low) and motor action (e.g. left→high). This improved goodness of fit ( $\Delta$  iBIC = 117)  
1138 and the resulting model fit indicated that motor-level model-free values had a somewhat stronger  
1139 influence on behaviour than the choice level model-free values (Figure 3a). With the perseveration  
1140 kernels and motor level representations included in each model, the mixture agent again provided a  
1141 better fit to the data than either a pure model-free ( $\Delta$  iBIC = 127) or pure model-based ( $\Delta$  iBIC = 227)  
1142 agent. We tested a number of other modifications to the model including separate learning rates at  
1143 the first and second step, but did not find further improvements in fit quality (Figure 5 – figure  
1144 supplement 1A). Finally, as adding features to the model may make other features which previously

1145 improved the fit unnecessary, we tested whether removing any individual component from the  
1146 model improved fit quality but again did not find further improvements (Figure 5 – figure  
1147 supplement 1B).

#### 1148 *Lesioning Full RL model*

1149 The simulations presented in Figure 3b indicated that data simulated from a model-based RL agent  
1150 showed loading on the transition and outcome predictors while data simulated from a model-free RL  
1151 agent showed loading only on outcome. This suggests that reduced influence of model-based and  
1152 increased influence of model-free RL could produce the observed effect of ACC inhibition. However,  
1153 the full RL model arrived at in the model comparison process included additional features not  
1154 included in those simulations which may complicate the relationship between behavioural strategy  
1155 and regression loadings. Specifically, we were concerned that perseveration or model-free RL for  
1156 actions represented at the motor level (i.e. as a movement from left to high poke, rather than as a  
1157 choice of the high poke irrespective of where the movement started) could produce loading on the  
1158 transition predictor. This is because the state transition determines which second-step state the  
1159 subject ends up in, and hence which motor action they must take to make a given choice on the next  
1160 trial. We therefore performed a set of simulations where we set the influence on choice of different  
1161 components of the model to zero on stimulation trials, which we term lesioning a model component  
1162 (Figure 5 – figure supplement 3). This confirmed that consistent with Figure 2B, lesioning the choice-  
1163 level model-free system selectively reduced loading on the outcome predictor, while lesioning the  
1164 model-based system reduced loading on outcome and transition, and to a lesser extent on the  
1165 interaction predictor. However, lesioning the motor-level model free system (which learned model-  
1166 free action values for individual motor actions such as left→high), also reduced loading on the  
1167 outcome and transition predictors, while lesioning motor-level perseveration reduced loading only  
1168 on the transition predictor. These simulations suggest that the reinforcing effect of experiencing a  
1169 common transition is mediated in part by the use of model-based RL but also in part by  
1170 perseveration and model-free RL occurring at the level of motor actions.