

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Flexibility to contingency changes distinguishes habitual and goal-directed strategies in
humans

Julie Lee^{1,#a*}, Mehdi Keramati^{1,#b}

¹Gatsby Computational Neuroscience Unit, University College London, London, UK

^{#a}Current Address: Institute of Ophthalmology, University College London, London, UK

^{#b}Current Address: Max Planck University College London Centre for Computational
Psychiatry and Ageing Research, University College London, London, UK

* Corresponding author

E-mail: julie.lee.15@ucl.ac.uk (JL)

1 **Abstract**

2 Decision-making in the real world presents the challenge of requiring flexible yet
3 prompt behavior, a balance that has been characterized in terms of a trade-off between a
4 slower, prospective goal-directed model-based (MB) strategy and a fast, retrospective
5 habitual model-free (MF) strategy. Theory predicts that flexibility to changes in both reward
6 values and transition contingencies can determine the relative influence of the two systems in
7 reinforcement learning, but few studies have manipulated the latter. Therefore, we developed
8 a novel two-level contingency change task in which transition contingencies between states
9 change every few trials; MB and MF control predict different responses following these
10 contingency changes, allowing their relative influence to be inferred. Additionally, we
11 manipulated the rate of contingency changes in order to determine whether contingency
12 change volatility would play a role in shifting subjects between a MB or MF strategy. We
13 found that human subjects employed a hybrid MB/MF strategy on the task, corroborating the
14 parallel contribution of MB and MF systems in reinforcement learning. Further, subjects did
15 not remain at one level of MB/MF behaviour but rather displayed a shift towards more MB
16 behavior over the first two blocks that was not attributable to the rate of contingency changes
17 but was rather a more general effect of block order. The extent to which each subject used
18 MB control was also related to reward earned, with a correlation between MB weight and
19 reward rate. We demonstrate that flexibility to contingency changes can distinguish MB and
20 MF strategies, with human subjects utilising a hybrid strategy that shifts towards more MB
21 behavior over blocks, consequently corresponding to a higher payoff.

22

23 **Introduction**

24 To make optimal decisions, humans must learn to associate the choices they make
25 with the outcomes that arise from them. Classical learning theories suggest that this problem

26 is addressed by habitual or goal-directed strategies for reinforcement learning [1, 2]. These
27 strategies differ in that habitual behavior seeks simply to reinforce responses based on
28 environmental cues, whereas goal-directed behavior considers action-outcome relationships –
29 that is, contingencies – in the environment. Habitual and goal-directed strategies have been
30 implemented in model-based (MB) and model-free (MF) reinforcement learning algorithms,
31 respectively. Both algorithms make decisions by estimating action values and choosing the
32 actions that maximize reward in the long term [3, 4]. The MF system achieves this
33 retrospectively, caching past rewards using a reward prediction error signal [5] whereas the
34 MB system achieves this prospectively, planning using a learned internal model of the state
35 transitions and rewards in the environment.

36 Recent studies have emphasized that MB and MF systems work in parallel rather than
37 in isolation [4, 6-8]. Early studies discerned MB and MF contributions using manipulations
38 of reward values, such as in reward devaluation paradigms, but did not seek to quantify their
39 relative contributions [1]. A recent study [6] addressed this by developing the hallmark “two-
40 step” task in which, using reward value changes, each trial was informative of the MB/MF
41 tradeoff, thereby permitting model-fitting analyses to quantify their relative influence in
42 decision-making. Human subjects showed a hybrid MB/MF strategy in the task, a result that
43 has been widely replicated under different manipulations [9, 10] and extended to the non-
44 human animal literature (Groman et al. Soc. Neurosci. Abstracts 2014, 558.19, Miranda et al.
45 Soc. Neurosci. Abstracts 2014 756.09, Akam et al. Cosyne Abstracts 2015, II-15; Hasz &
46 Redish, Soc. Neurosci. Abstracts 2016 638.08; [11]).

47 Theory predicts that flexibility to transition contingency changes can – like flexibility
48 to reward value changes – determine the relative influence of MB and MF strategies [4, 12].
49 The advantage of manipulating transitions, rather than reward values, is apparent when
50 contrasting the model-based system to a successor representation (SR) [13]. The successor

51 representation caches transitions in a model-free fashion, but learns reward values in a model-
52 based fashion; thus, changes to reward values cannot distinguish MB and SR representations.
53 In contrast, transition changes ensure that consequent choices only can be explained by an
54 MB system. Two studies have examined the flexibility of MB and MF systems to global
55 contingency changes [14, 15]. However, quantification of the MB/MF tradeoff was limited as
56 these studies manipulated contingency and tested flexibility to the change of contingency in
57 separate phases; at these timescales, it becomes difficult to exclude the effect of adaptation on
58 MB/MF weights. Therefore, we developed a novel two-level contingency change task
59 containing multiple, frequent and interleaved transition contingency changes that elicit
60 different consequent actions by the MB and MF systems. Our design, like the two-step task
61 [6] and its variants, therefore permits model-fitting analyses to robustly determine the relative
62 influence of the MB/MF systems. The contingency change task is structured such that actions
63 following frequent contingency changes are distinctly attributed to either a MB or MF
64 strategy; this then permits quantification of the degree to which each system is in control.

65 On top of a hybrid MB/MF strategy, subjects may not remain at one level of MB/MF
66 control but instead shift their relative weight in accordance with environmental factors. In
67 general, animals show habit formation with time, a robust effect reported since early reward
68 devaluation studies [16] in which extensive training stamped in habits, resulting in
69 insensitivity to reward devaluation; in contrast, limited training retained goal-directed
70 behavior. Sensitivity to contingency degradation (the omission of a previously-learned
71 contingency between actions and outcomes) also decreases with overtraining, likewise
72 reflecting a trend towards habitization with time [17]. In the original two-step task, the
73 MB/MF trade-off was designed to be stable [6], but will shift under manipulations such as
74 limited time [8] or cognitive load [18]. However, habits are not guaranteed to form with time;
75 even after extended training, rats can show residual responding following outcome

76 devaluation, indicating that they retained goal-directed behavior despite overtraining [19]. In
77 another study using the two-step task [20], the level of MB/MF control in fact increased in
78 favour of more MB control (i.e. towards less habitual behavior) over three days of training.
79 However, general shifts in MB/MF control should be disentangled from the effects of
80 environmental volatility, which are known to affect the MB/MF balance [21]. Thus, in this
81 study, we examined whether the MB/MF relationship is affected by environmental stability,
82 or whether it shifts more generally over time.

83 We found that human subjects indeed showed a hybrid strategy in reacting to
84 contingency changes in our task, with an increased influence of MB control over the first two
85 blocks. However, relative MB/MF control did not significantly differ across rates of
86 contingency changes; thus, the increase in MB control may be a more global effect of “anti-
87 habitization” over time. The increased reliance on the MB system was associated with a
88 higher proportion of highly rewarded actions and consequently a higher reward rate,
89 indicating that as subjects proceeded through the session, they became more proficient at
90 exploiting their learned internal model of the task structure to maximize their reward.

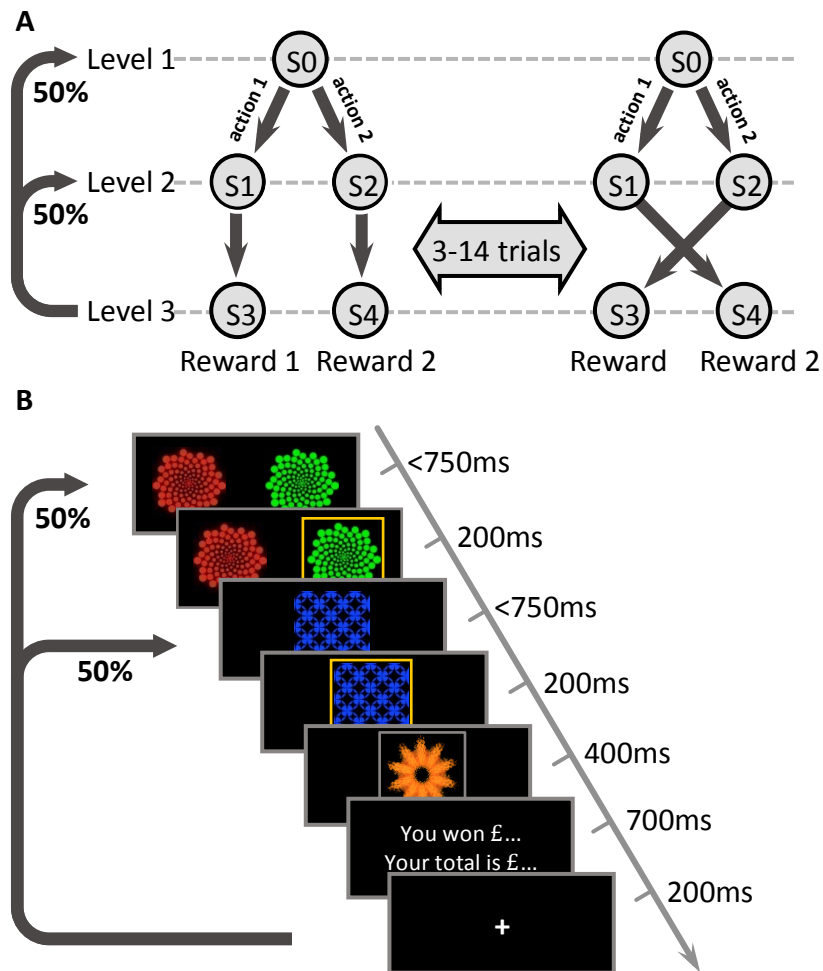
91

92 **Results**

93 Subjects (N=16) performed a two-level contingency change task which consisted of
94 600 trials (Fig 1). Each trial began at either the first level (S0) with 50% probability, or the
95 second level with 50% probability – 25% for each of the two states at this level (S1 or S2). If
96 a trial started at the first level, a two-alternative choice was possible between two abstract
97 stimuli. Each first-level action deterministically always led to the same second-level state, i.e.
98 A1 to S1 and A2 to S2. Critically however, transitions from the second-level states to the
99 terminal states flipped between two contingencies every 3-14 trials. Each of the two terminal
100 states was then associated with either high or low reward, with the exact reward values

101 drifting across trials (see Methods for details). Thus, flexibility to contingency changes was
 102 essential for maximizing reward.

103



104

105 **Fig 1. Schematic of the experimental design.** (A) Each trial started from either the first-
 106 level state (S0), with 50% probability, or one of the two second-level states (S1 or S2), each
 107 with 25% probability. While two choices were available at S0, only a forced choice was
 108 available at the second-level states. The transition structure from the second-level states to
 109 terminal states repeatedly flipped after a random number of trials (every 3- 14), in an
 110 unsignalled fashion. One of the two terminal states (S3 or S4) was associated with a high
 111 reward outcome and the other with a low reward outcome. (B) Timeline of the task for one
 112 example trial.

113

114 If a contingency change occurred, subjects always experienced the new transition
115 structure regardless of whether they started at the first or second level, as contingency could
116 only change between second-level and terminal states. Therefore, provided that an action was
117 possible at the next trial (i.e. that the next trial started at the first level) the MB system would
118 plan using the updated causal structure and thus would take the action that led under the new
119 transition contingencies to the high reward terminal state. However, if a contingency change
120 trial started from the second level, the MF system would not choose the optimal action on the
121 next trial, as neither the received reward nor the new contingency would update the cached
122 values of first-level actions, simply because no first-level action was experienced on those
123 trials. As a result, the relative contribution of MB and MF systems can be measured by the
124 degree of behavioral flexibility on first-level trials following contingency change trials
125 starting from the second level.

126 To examine the effect of environmental volatility on the contribution of the two
127 systems, the frequency of contingency changes was varied – from 3-6 trials for 200 trials, to
128 7-10 trials for the next 200 trials, and then 11-14 trials for the final 200 trials. The order of
129 fast and medium contingency changes was counterbalanced across two subject groups (n=8
130 each). Every 40 trials, assignment of the high and low reward states also flipped to prevent
131 formation of habits over an extended state representation, which could masquerade MF as
132 MB behavior [22].

133 Simulated choices on the task were implemented according to MB and MF
134 reinforcement learning algorithms (see Methods for details). For each system, we measured a
135 “stay probability” index which followed the logic of contingency change trials described
136 above. This index differs from classic stay probabilities [6] as trials starting from the second
137 level do not have any choices to “stay”. Instead, stay probability in our task was defined as

138 the probability of choosing the first-level action that results in the same second-level state as
139 the previous trial. Since first-level to second-level contingencies were fixed, this modified
140 measure provided stay probabilities on any trial, regardless of whether it started at the first or
141 second level. Stay probability was measured for four different conditions: whether the reward
142 received in the previous trial was “high” or “low”, and whether the transition experienced in
143 the previous trial, relative to the trial before that, was “changed” or remained “fixed”. In all
144 cases, analyses were restricted to trials starting from the first level, following a contingency
145 change trial starting at the second level, since only these could distinguish MB and MF
146 strategies.

147 Across these conditions, MB and MF systems showed different stay probability
148 patterns. The MF system, having no experience of the action that led to the new contingency,
149 was more likely to stay on the action leading to the high reward state, and shift on the action
150 leading to the low reward state, under “fixed” than “changed” conditions ($p < 0.01$),
151 indicating it was not flexible to changes in contingencies (Fig 2A). However, the MB system
152 could immediately adapt with the correct next action, staying on the action if it would lead to
153 the high-reward state but shifting if it would lead to the low-reward state, with a main effect
154 of reward ($p < 0.01$) regardless of contingency condition (Fig 2B). As expected, for
155 contingency changes from the first level, MB and MF systems did not differ in stay
156 probability patterns, as the MF system was able to update its action values accordingly, given
157 that it directly experienced the action leading to the new contingency (S1 Fig). In addition to
158 pure MF and pure MB strategies, we simulated a hybrid model that linearly weights MB and
159 MF action values according to a parameter w_{MB} . The stay probability pattern produced by
160 this hybrid system reflected a mixture of the effects observed for the pure MF and MB stay
161 probabilities – that is, showing a main effect of reward ($p < 0.01$), but also an interaction
162 between reward and contingency ($p < 0.01$) (Fig 2C).

163



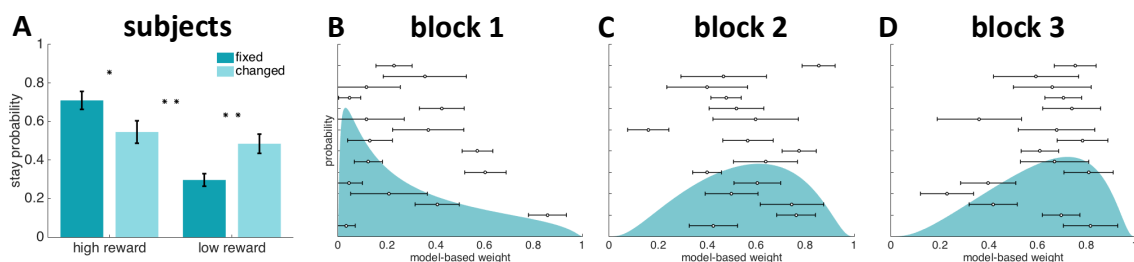
164

165 **Fig 2. Stay probability patterns predicted by simulating model-based (A), model-free**
166 **(B), and hybrid (C) reinforcement learning algorithms.** Stay probability measures the
167 probability of choosing the first-level action that results in the same second-level state as the
168 previous trial. This index was measured when the reward received in the previous trial was
169 “high” or “low”, and when the transition experienced in the previous trial (relative to the trial
170 before that) had its contingency “changed” or remained “fixed”. Stay probabilities are plotted
171 for trials following a change trial that started at the second level, as these distinguish model-
172 based and model-free strategies. * $p < 0.05$, ** $p < 0.01$

173

174 Subjects showed hallmarks of both MB and MF strategies in reacting to contingency
175 changes (Fig 3A), showing a main effect of reward, $F(1,60) = 24.65$, $p < 0.01$, as well as a
176 reward/contingency interaction, $F(1,60) = 13.60$, $p < 0.01$. Therefore, subjects did not solely
177 use a MB or MF strategy when reacting to contingency changes, but rather displayed a hybrid
178 MB/MF strategy.

179



180

181 **Fig 3. Experimental results.** (A) Stay probability pattern from human subjects (N=16)
182 showed significant effects of both model-based ($p < 0.01$) and model-free ($p < 0.01$) strategies.
183 * $p < 0.05$, ** $p < 0.01$ (B) Probability density function over the model-based weight
184 parameter, estimated in three different blocks of the first, middle and last 200 trials (out of
185 600 trials). Overlaid are the individual subjects' model-based weight parameter estimates for
186 each block type. Error bars represent standard deviation.

187

188 To characterize the effect of contingency changes over multiple consecutive trials,
189 lagged logistic regression was performed (S2 Fig). This analysis computes the influence of
190 reward and contingency conditions (predictors) from past trials (lags) on choice probabilities
191 [22, 23]. MB and MF systems differed in the extent to which past predictors influenced the
192 current choice, with the MB system showing more flexibility to recent changes – and less
193 influence of past predictors – than the MF system; this was evidenced by a smoother
194 predictive weight over lags for the MF system than the MB system (S2 Fig). As expected, the
195 pattern of the predictive weights for subjects' choices and the simulated hybrid model
196 reflected a mixture of the MF and MB systems' patterns.

197 While stay probabilities excluded a purely MB or purely MF strategy, this measure
198 could not quantify the degree to which subjects used the hybrid strategy; therefore, we used a
199 hierarchical Bayesian method to fit candidate models of behavior to the subjects' data, to
200 determine which model best explained subjects' choices and to obtain parameter estimates for
201 the MB/MF weighting used by the subjects. The models tested included a pure MB model, a
202 pure MF model, a hybrid model with one constant weight w_{MB} across the session, a hybrid
203 model with three separate w_{MB} weights for the three experimental blocks (which differed in
204 terms of frequency of contingency changes: fast, medium, or slow), and a hybrid model with
205 three separate w_{MB} weights for each range of contingency changes rates. The last two

206 models served to test whether the relative contribution of the two systems depended on
207 volatility of transition structure, or instead more generally trial order. Model-fitting was
208 confirmed to be able to recover true parameter values, as median estimated parameter values
209 from model-fitting (see Methods for details) were well-correlated to known parameter values
210 from simulations, $r \geq 0.99$, $p < 0.01$.

211 Model-fitting results supported the existence of a hybrid MB/MF strategy in our task.
212 Candidate models were compared using two criteria – integrated Bayesian Information
213 Criterion which controls for number of parameters (iBIC) [24] and exceedance probabilities
214 [25] (S2 Table). The hybrid model with three wMB weights over blocks outperformed the
215 other candidate models on both criteria, with the lowest iBIC and a probability of 89.4% that
216 it was the most common of the four models across subjects. Thus, from here we only discuss
217 the results of best-fit model, the three-block hybrid model.

218 The median fitted wMB weights in the three-block hybrid increased across the three
219 blocks (Fig 3B-D), indicating some extent of “anti-habitization” rather than habit formation.
220 The increase of wMB from block 1 to block 2, but not the increase from block 2 to block 3,
221 was significant according to permutation tests, $p < 0.01$. Stay probability analyses were not
222 conducted on the three separate blocks, as slower contingency changes meant that the later
223 blocks had fewer samples of contingency changes for comparison. The increase in wMB
224 across blocks was not attributable to differences in quality of fit from the model-fitting
225 procedure, as the log-likelihood of parameter estimates did not differ significantly across
226 blocks, $F(2,45) = 1.42$, $p > 0.05$. Strength of correlations between fitted and simulated wMB
227 weights were also similar across blocks (block 1: $r = 0.99$, block 2: $r = 1.00$, block 3: $r =$
228 0.99 ; $p < 0.01$ for all blocks). Therefore, the significant increase in wMB from the first to
229 second block was not caused by differences in quality of model fit.

230 To confirm that the increase in model-based weight was not due to differences in the
231 rate of contingency changes, we further analysed the fitted weights from the three-frequency
232 hybrid model, which had a different wMB assigned to each range of contingency change
233 rates, i.e. fast (every 3-6 trials), medium (every 7-10 trials) and slow (every 11-14 trials)
234 contingency change blocks. The estimated wMB weights (S3 Fig) were not significantly
235 different between fast vs. medium, or medium vs. slow frequency of contingency change
236 blocks in permutation tests, $p > 0.05$. Thus, the increase in wMB in our study seemed to be an
237 effect of block order rather than environmental volatility from differences in contingency
238 change rates. In summary, subjects became more model-based across the first two blocks but
239 did not differ in MB influence between different rates of contingency changes; therefore, it
240 seems that block order, but not contingency change volatility, affects wMB in our task.

241 As subjects became more model-based, high reward choices and consequently reward
242 rate also increased. Choice probabilities for the high reward action differed over blocks,
243 $F(2,45) = 5.77$, $p < 0.01$, with post-hoc tests finding a significant increase between the first
244 and third blocks ($p < 0.01$) and the second and third blocks ($p < 0.05$). Additionally, there
245 was a significant difference in reward rate across blocks, $F(2,45) = 3.83$, $p < 0.05$, increasing
246 between the first and third blocks ($p < 0.05$). Mean reaction time and number of missed trials
247 due to timeout did not significantly change across blocks, $p < 0.05$; therefore, the increase in
248 high reward choices over blocks was not necessarily because subjects were worse at the task
249 to begin with. Two analyses were performed to rule out the possibility of practice effects
250 driving the association between reward rate and model-based weight. Within each block,
251 there was a significant correlation of each subject's median wMB and reward rate (block 1: r
252 $= 0.66$, $p < 0.01$, block 2: $r = 0.65$, $p < 0.01$, block 3: $r = 0.56$, $p < 0.05$), indicating that on an
253 individual subject basis, the extent of MB control was related to reward earned. Since these
254 analyses were conducted within blocks, the association with reward rate could not be

255 accounted for by block order. Additionally, the hybrid model was simulated using a range of
256 MB weights (0, 0.2, 0.4, 0.6, 0.8 and 1) using the one-weight hybrid model for simplicity.
257 There was a significant effect of MB weight on reward rate, $F(5,90) = 8.5$, $p < 0.01$. In all,
258 these findings suggest that MB influence in this task truly corresponded to a better “payoff”
259 in terms of reward gained.

260

261 **Discussion**

262 We developed a novel two-level contingency change task in which flexibility to
263 frequently-changing transition contingencies between states could determine whether subjects
264 were using a model-based or a model-free strategy. Subjects showed a hybrid strategy when
265 reacting to contingency changes, corroborating recent evidence of the parallel contribution of
266 MB and MF systems in reward-guided decision-making. Importantly, this finding confirmed
267 that changes to transition contingencies can elicit a balance of MB and MF behavior akin to
268 changes to reward values. Model-fitting analyses indicated that a hybrid model with three
269 MB weights best explained subjects’ choices, with relative MB control increasing over
270 blocks. The rate of contingency changes did not significantly shift the MB/MF balance;
271 rather, MB control increased over the first two blocks of trials. This increase in MB control
272 was concurrent with an increased proportion of high reward choices and consequently
273 increased reward rate; individually, each subject’s MB was also correlated with reward
274 gained in the same block.

275 In all, these results illustrated that not only do subjects use a mixed MB/MF strategy,
276 but within this hybrid strategy, the trade-off shifts towards “anti-habitization” across the first
277 two blocks. This agrees with a previous study [20] that used the two-step task over three
278 days, reporting that their subjects’ MB weight increased across days. One distinction between
279 our findings is that in [20], subjects started relatively model-based (i.e. median $w_{MB} > 0.5$)

280 whereas in our case, subjects began relatively model-free (i.e. median $w_{MB} < 0.5$). This
281 difference in starting MB weight simply may be due to individual differences, which is
282 evident even within our subject pool. Alternatively, differences could be accounted for by the
283 relatively short reaction time limit in our task compared to theirs (750ms in ours vs. 2000ms).
284 A shorter reaction time limit is known to provide a depth-of-planning pressure and favor
285 more MF control [8]. Hence, our subjects may have started more model-free and only
286 become more model-based once they mastered prospective planning of the task structure.
287 This is supported by the lack of significant changes in reaction time across blocks, suggesting
288 that subjects may have used the full extent of their time and eventually learned more efficient
289 planning under time pressure, therefore showing increased MB influence over blocks.

290 These findings of an increase in MB control over blocks, however, goes against
291 another study [26] using a similar task to the two-step task, that found an exponential decay
292 in MB weight over the experimental session, or habit formation. This difference in findings is
293 likely because they used a fixed rather than drifting amount of reward; in stationary
294 environments such as these, habit formation can occur from overtraining, manifesting in an
295 increase in MF rather than MB behavior [21]. Thus, these results point to the importance of
296 maintaining a changing environment, as subjects can otherwise adapt to the change and
297 become habitized.

298 Manipulations of the rate of contingency changes did not seem to affect MB/MF
299 control. While it has been shown that environmental volatility can influence MB/MF levels in
300 the context of reward value changes [21], in our case, the kind and range of contingency
301 change volatility did not elicit a significant difference in relative MB/MF control. Further
302 work is certainly needed to definitively rule out the possibility that environmental volatility in
303 the form of the rates of contingency changes does not affect MB weight, but in the present

304 study, we find that subjects did not change their use of MB control with contingency change
305 volatility, but rather increased MB influence more generally with block order.

306 In conclusion, in a two-level contingency change task, subjects showed a hybrid
307 MB/MF strategy, emphasizing their parallel contribution in reacting to changes in transition
308 contingencies. The inclusion of multiple, frequent changes allowed us to perform model-
309 fitting; by doing so, we found an increase in MB control over the first two blocks, a result not
310 detectable in model-agnostic analyses alone. Our results build on the literature of a hybrid
311 MB/MF strategy in reacting to changes in reward values, demonstrating a mixture of
312 strategies in reacting to multiple, frequent contingency changes that has yet been unexplored.
313 This novel paradigm therefore provides another avenue for exploring the relationship
314 between MB and MF control for future studies in neuropsychiatric disorders that may
315 differentially implicate this balance between changes in transition contingencies and changes
316 in reward values.

317

318 **Methods**

319 **Subjects**

320 Sixteen subjects (nine males, mean age 24 years) took part. The study was approved
321 by the University College London Research Ethics Committee (Project ID 3450/002). All
322 subjects provided written informed consent.

323

324 **Experimental procedure**

325 Subjects performed 600 trials of three blocks (200 each) which differed in frequency
326 of contingency changes: fast (every 3-6 trials), medium (contingency change every 7-10
327 trials) or slow (every 11-14 trials). Each subject was assigned to one of two groups (n=8
328 each), which differed by the order of presentation of fast and medium contingency change

329 blocks, i.e. half of the subjects had fast, medium, then slow contingency changes, and the
330 other half started with medium, fast, then slow frequency of contingency changes.

331 To ensure subjects understood the task structure, they were first trained with practice
332 stimuli (35 trials) then trained on novel test stimuli without reward (55 trials) before starting
333 the experimental session. Subjects were informed that contingency changes would occur, but
334 did not know the frequency of changes nor that those rates would vary across the session.

335 At the first level, subjects had a two-alternative forced choice between two actions
336 (pressing ‘S’ for the action available on the left side of the screen, ‘L’ for the right) with the
337 presentation of stimuli randomized for the left/right side of the screen. To ensure that subjects
338 recognized second-level states, they had to press ‘D’ if they encountered one of these states,
339 and ‘K’ for the other. Both responses had a time limit of 750ms, following which the trial
340 would end with no reward. Missed trials were not repeated.

341 Payoff at the high-reward terminal state varied with a drift rate of 0.2 and offset of
342 0.15 to the bound of £1, with payoff at the low-reward terminal state being £1 minus the
343 reward of the high-reward terminal state. Subjects received a fixed proportion of their total
344 reward gained, with payoff bounded between £5 and £25. To make the task adequately
345 difficult and prevent formation of complex state-space representations [22], high- and low-
346 reward assignments switched every 40 trials between the two terminal states. This change
347 was designed never to co-occur with contingency changes.

348

349 **Model**

350 Both model-free and model-based algorithms seek to estimate the values of state-
351 action pairs in order to choose the actions which can maximize expected future rewards. The
352 state space was modelled as having a first-level state s_0 with two actions a_1 and a_2 , two
353 possible second-level states s_1 and s_2 , and two possible terminal states s_3 and s_4 . There was

354 only one action available on second-level and terminal states, as the subject did not have any
 355 choices at these levels.

356 The model-free algorithm updates values of state-action pairs using temporal
 357 difference Q-learning [3, 27]. The reward r_t is used to compute a reward prediction error δ_t
 358 which updates action values for that state s and action a at time t , $Q_{MF}(s_t, a_t)$. At the first
 359 level r_t is set to be 0 as there is no reward at this level.

$$\delta_t = r_t + \max_{a'} [Q_{MF}(s_{t+1}, a')] - Q_{MF}(s_t, a_t)$$

$$Q_{MF}(s_t, a_t) = Q_{MF}(s_t, a_t) + \alpha_{MF} \lambda \delta_t$$

360 The reward prediction error updates existing action values according to a learning rate
 361 α_{MF} and modified by the eligibility parameter λ . Eligibility governs how much credit past
 362 actions were given for outcomes, with $\lambda = 0$ corresponding to a pure TD algorithm whereby
 363 first-stage actions are updated only by the second-level action values, which in turn is
 364 updated by terminal state rewards. In contrast, $\lambda = 1$ means the algorithm updates first-level
 365 actions only using the final reward from the terminal state reached on that trial.

366 The model-based algorithm learns both transition probabilities P_T and reward
 367 probabilities R_T . The transition probabilities track the transition contingencies P_T between
 368 states s and subsequent states s' . Upon encountering a contingency change, the model-based
 369 system always updated its knowledge of both transitions.

$$P_T \left(s_1 \xrightarrow{a} s_3 \right) = \begin{cases} 1, & \text{if } s' = s_3, s = s_1 \\ 0, & \text{otherwise} \end{cases}$$

$$P_T \left(s_2 \xrightarrow{a} s_4 \right) = \begin{cases} 1, & \text{if } s' = s_4, s = s_2 \\ 0, & \text{otherwise} \end{cases}$$

$$P_T \left(s_1 \xrightarrow{a} s_4 \right) = 1 - P_T \left(s_1 \xrightarrow{a} s_3 \right)$$

$$P_T \left(s_2 \xrightarrow{a} s_3 \right) = 1 - P_T \left(s_2 \xrightarrow{a} s_4 \right)$$

370 The reward probabilities R_T use the reward r_t to update its subjective reward R for
371 that state s and action a at time t .

$$R(s_t, a_t) = R(s_t, a_t) + \alpha_{MB}(R(s_t, a_t) - r_t)$$

372 These learned transition and reward functions are then used to update the action
373 values for the model-based system, Q_{MB} .

$$Q_{MB}(s_t, a_t) = P_T(s \xrightarrow{a} s_3) \cdot R(s_3, a) + P_T(s \xrightarrow{a} s_4) \cdot R(s_4, a)$$

374 Other parameters from the simulated models included learning rates for model-based
375 and model-free systems, α_{MB} and α_{MF} , and a stay bias which temporarily increased the action
376 value for the previously-selected action regardless of outcome, to quantify a perseveration
377 bias. These additional parameters improved fit even when controlling for model complexity
378 (S3 Table).

379 For both systems, values for the non-selected action were updated as well, assuming
380 that subjects knew that the reward for the selected action and reward for the non-selected
381 action were negatively related, according to proposals of fictive reward [28]. Action values
382 were updated for both visited and non-visited states, with the action values of non-visited
383 states corresponding to $1 - Q(s_t, a_t)$ of the visited states. The inclusion of fictive reward
384 updates resulted in a better fit to the subjects' choices (S3 Table).

385 The hybrid model weighted MB and MF action values according to a parameter
386 wMB , with $wMB = 1$ indicating fully MB control:

$$Q_{hybrid}(s_t, a_t) = wMB \cdot Q_{MB}(s_t, a_t) + (1 - wMB) \cdot Q_{MF}(s_t, a_t)$$

387 Action selection was then determined for all models according to a "softmax" rule
388 which computes action probabilities as proportional to the exponential of the action values.

$$p(a_t = a_1 | s_t) = \frac{\exp(\beta \cdot Q(s_t, a_1))}{\exp(\beta \cdot Q(s_t, a_1)) + \exp(\beta \cdot Q(s_t, a_2))}$$

389 The inverse temperature β determined the extent to which action selection was
390 stochastic or deterministic from action values, quantifying an exploration/exploitation trade-
391 off.

392

393 **Simulations**

394 To best replicate the subjects' data of 600 trials for 16 subjects, each simulation was
395 run for 16 initializations of 600 trials each. All reported simulations used fitted parameters
396 from the three-block hybrid model for the learning rates α_{MF} and α_{MB} , inverse temperature
397 β , eligibility trace λ and stay bias (S1 Table). wMB values were 1 for pure MB and 0 for pure
398 MF models.

399

400 **Model-fitting**

401 Subjects' data were fit to the models using mixed effects hierarchical model fitting.
402 Estimation-maximisation was used which iteratively generates group-level distributions over
403 individual subject parameter estimates, choosing the parameters that maximizes the
404 likelihood of the data given those estimates. Parameters were estimated by minimizing the
405 negative log-likelihood of parameter estimates using *fminunc* in Matlab (MathWorks).

406 To ensure the efficacy of wMB parameter estimation for the candidate model, each
407 block wMB was simulated for 11 different parameter values: 0, 0.1, 0.2, ... 1. These resulted
408 in a total of 33 parameter settings for $wMB1$, $wMB2$, $wMB3$, with 16 iterations per setting. All
409 other parameters in the simulations were set constant as the median parameter estimates taken
410 from the hybrid three-block model from model-fitting on the subjects' data. The same model-
411 fitting procedure was performed on the simulated data and estimated parameter values were
412 extracted.

413 The integrated Bayesian information criterion (iBIC) [24] was used to compare the
414 fits of candidate models to the data, with lower scores indicating better fit; this criterion
415 penalizes more complex models. Finally, Bayesian model selection [25] was used to examine
416 the prevalence of each model in the participant population. This quantifies an exceedance
417 probability, the probability that each model is the most common in the subject pool.

418

419 **Permutation tests**

420 Permutation tests were run to evaluate the probability that wMB could differ across
421 blocks by chance. Subjects' blocks were randomly permuted such that each "block"
422 contained a mixture of true first, second and third blocks. Model-fitting was run on each
423 permutation to extract parameter estimates of wMB for each new "block". The probabilities
424 $p(wMB_{block\ 2} > wMB_{block\ 1})$, and $p(wMB_{block\ 3} > wMB_{block\ 2})$ were then evaluated for each
425 permutation. The occurrences of the random permutations which had a smaller $p(wMB_{block\ 2} >$
426 $wMB_{block\ 1})$, and $p(wMB_{block\ 3} > wMB_{block\ 2})$ than the true permutation were then tallied.

427 Likewise, to evaluate the effect of frequency of contingency changes, permutation
428 tests were run to compare wMB for fast, medium and slow contingency change blocks. Each
429 subject was randomly assigned to one of the two groups (which differed in the order of fast
430 and medium contingency change blocks) then wMB of each frequency block was computed
431 for each permutation. Both the aforementioned one-tailed permutation test and a two-tailed
432 Hellinger distance permutation test were used.

433

434 **Acknowledgements**

435 Thanks to Peter Dayan for supervision and comments, and Thomas Akam for
436 comments on the manuscript. JL is supported by a Wellcome Trust doctoral fellowship. MK
437 is supported by the Gatsby Charitable Foundation.

439 **References**

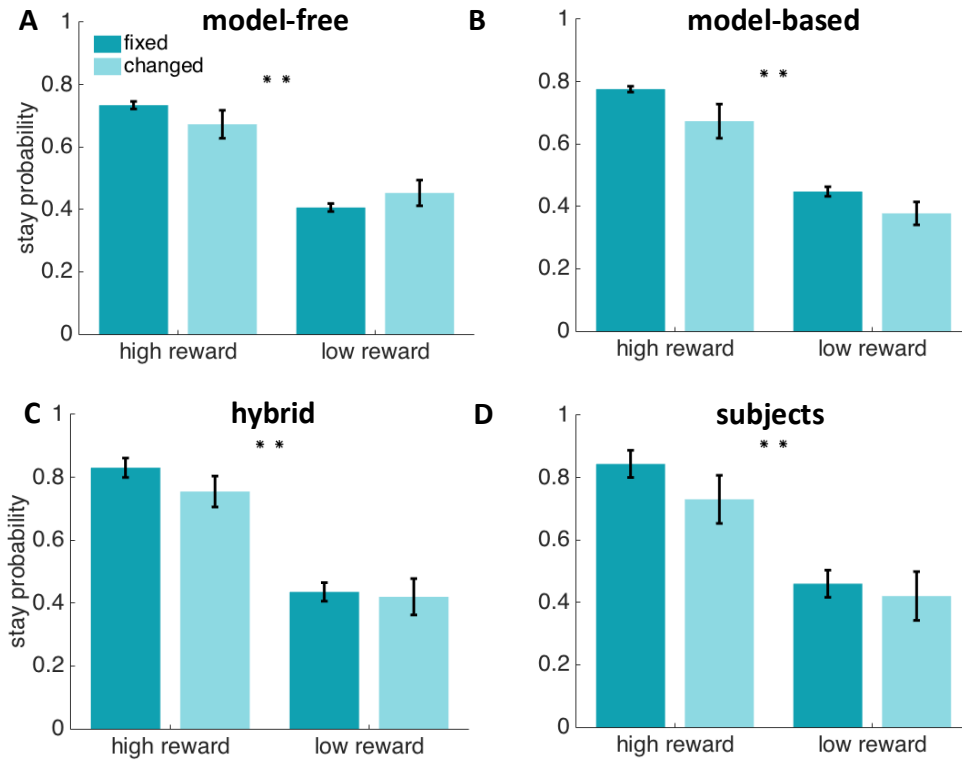
- 440 1. Adams CD, Dickinson A. Instrumental responding following reinforcer devaluation.
441 The Quarterly journal of experimental psychology. 1981;33(2):109-21.
- 442 2. Dickinson A, Balleine B. Motivational control of goal-directed action. *Animal*
443 *Learning & Behavior*. 1994;22(1):1-18.
- 444 3. Sutton RS, Barto AG. Reinforcement learning: An introduction: MIT press
445 Cambridge; 1998.
- 446 4. Daw ND, Niv Y, Dayan P. Uncertainty-based competition between prefrontal and
447 dorsolateral striatal systems for behavioral control. *Nature Neuroscience*.
448 2005;8:1704-11. doi: 10.1038/nn1560.
- 449 5. Schultz W, Dayan P, Montague PR. A neural substrate of prediction and reward.
450 *Science*. 1997;275(5306):1593-9.
- 451 6. Daw ND, Gershman SJ, Seymour B, Dayan P, Dolan RJ. Model-Based Influences on
452 Humans' Choices and Striatal Prediction Errors. *Neuron*. 2011;69:1204-15. doi:
453 10.1016/j.neuron.2011.02.027.
- 454 7. Keramati M, Dezfouli A, Piray P. Speed/accuracy trade-off between the habitual and
455 the goal-directed processes. *PLoS computational biology*. 2011;7:e1002055. doi:
456 10.1371/journal.pcbi.1002055. PubMed PMID: 21637741.
- 457 8. Keramati M, Smittenaar P, Dolan RJ, Dayan P. Adaptive integration of habits into
458 depth-limited planning defines a habitual-goal-directed spectrum. *PNAS*. 2016:1-16.
459 doi: 10.1073/pnas.1609094113. PubMed PMID: 27791110.
- 460 9. Otto AR, Skatova A, Madlon-Kay S, Daw ND. Cognitive control predicts use of
461 model-based reinforcement learning. *Journal of cognitive neuroscience*. 2014.

- 462 10. Wunderlich K, Smittenaar P, Dolan RJ. Dopamine Enhances Model-Based over
463 Model-Free Choice Behavior. *Neuron*. 2012;75:418-24. doi:
464 10.1016/j.neuron.2012.03.042. PubMed PMID: 22884326.
- 465 11. Miller KJ, Botvinick MM, Brody CD. Dorsal hippocampus plays a causal role in
466 model-based planning. *bioRxiv*. 2016. doi: 10.1101/096594.
- 467 12. Balleine BW, O'Doherty JP. Human and Rodent Homologies in Action Control:
468 Corticostriatal Determinants of Goal-Directed and Habitual Action.
469 *Neuropsychopharmacology*. 2010;35:48-69. doi: 10.1038/npp.2009.131.
- 470 13. Dayan P. Improving generalization for temporal difference learning: The successor
471 representation. *Neural Computation*. 1993;5(4):613-24.
- 472 14. Momennejad I, Russek EM, Cheong JH, Botvinick MM, Daw N, Gershman SJ. The
473 successor representation in human reinforcement learning. *bioRxiv*. 2016:1-27. doi:
474 10.1101/083824.
- 475 15. Gershman SJ, Markman AB, Otto AR. Retrospective revaluation in sequential
476 decision making: A tale of two systems. *Journal of Experimental Psychology:*
477 *General*. 2014;143(1):182.
- 478 16. Adams CD. Variations in the sensitivity of instrumental responding to reinforcer
479 devaluation. *The Quarterly Journal of Experimental Psychology*. 1982;34(2):77-98.
- 480 17. Dickinson A. Omission learning after instrumental pretraining. *The Quarterly Journal*
481 *of Experimental Psychology: Section B*. 1998;51(3):271-86.
- 482 18. Otto AR, Gershman SJ, Markman AB, Daw ND. The curse of planning: dissecting
483 multiple reinforcement-learning systems by taxing the central executive.
484 *Psychological science*. 2013;24(5):751-61.

- 485 19. Colwill RM, Rescorla RA. Instrumental responding remains sensitive to reinforcer
486 devaluation after extensive training. *Journal of Experimental Psychology: Animal*
487 *Behavior Processes*. 1985;11(4):520.
- 488 20. Economides M, Kurth-Nelson Z, Lübbert A, Guitart-Masip M, Dolan RJ. Model-
489 Based Reasoning in Humans Becomes Automatic with Training. *PLoS computational*
490 *biology*. 2015:1-19. doi: 10.1371/journal.pcbi.1004463.
- 491 21. Simon DA, Daw ND. Environmental statistics and the trade-off between model-based
492 and TD learning in humans. *Advances in Neural Information Processing Systems*
493 (NIPS). 2011:1-9.
- 494 22. Akam T, Costa R, Dayan P. Simple Plans or Sophisticated Habits? State, Transition
495 and Learning Interactions in the Two-step Task. *PLoS computational biology*.
496 2015:021428. doi: 10.1101/021428.
- 497 23. Miller KJ, Brody CD, Botvinick MM. Identifying Model-Based and Model-Free
498 Patterns in Behavior on Multi-Step Tasks. *bioRxiv*. 2016:096339.
- 499 24. Huys QJ, Eshel N, O'Nions E, Sheridan L, Dayan P, Roiser JP. Bonsai trees in your
500 head: how the Pavlovian system sculpts goal-directed choices by pruning decision
501 trees. *PLoS Comput Biol*. 2012;8(3):e1002410.
- 502 25. Stephan KE, Penny WD, Daunizeau J, Moran RJ, Friston KJ. Bayesian model
503 selection for group studies. *NeuroImage*. 2009;46:1004-17. doi:
504 10.1016/j.neuroimage.2009.03.025.
- 505 26. Gläscher J, Daw Nathaniel D, Dayan P, O'Doherty JP. States versus Rewards:
506 Dissociable Neural Prediction Error Signals Underlying Model-Based and Model-
507 Free Reinforcement Learning. *Neuron*. 2010;66:585-95. doi:
508 10.1016/j.neuron.2010.04.016.
- 509 27. Watkins CJ, Dayan P. Q-learning. *Machine learning*. 1992;8(3-4):279-92.

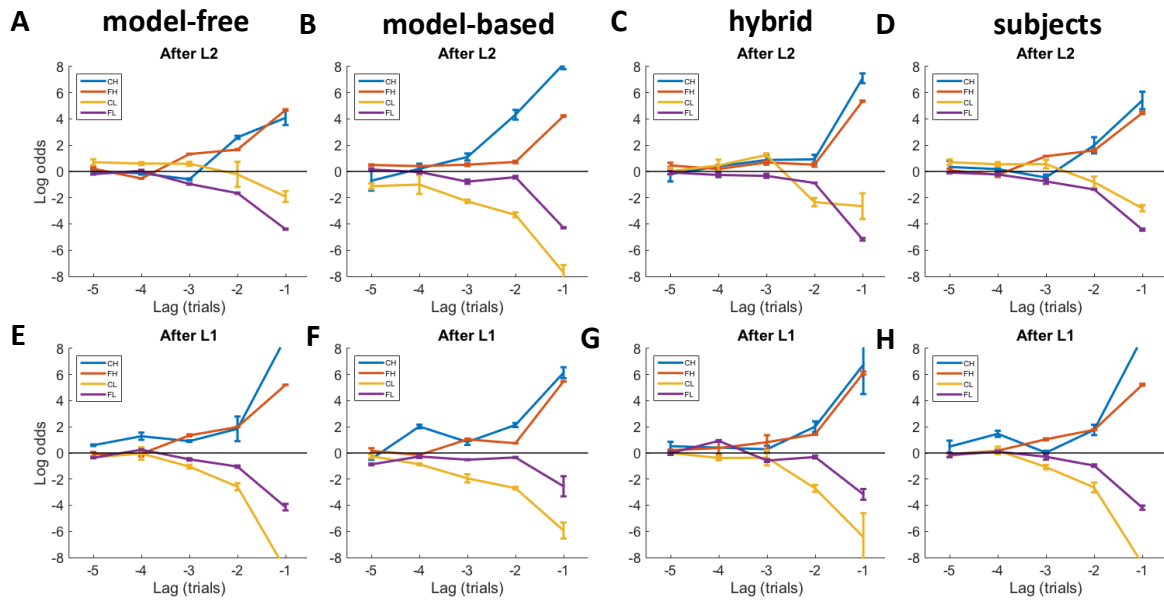
510 28. Lohrenz T, McCabe K, Camerer CF, Montague PR. Neural signature of fictive
511 learning signals in a sequential investment task. Proceedings of the National Academy
512 of Sciences. 2007;104(22):9493-8.

513 **Supporting Information**



514
515 **S1 Fig. Stay probability patterns after first level contingency changes predicted by**
516 **simulating model-based (A), model-free (B), and hybrid (C) reinforcement learning**
517 **algorithms, along with experimental results (D).** Stay-probability measures the probability
518 of choosing the first-level action that results in the same second-level state as the previous
519 trial, following a trial that started at the first level. For each system, this index was measured
520 under four different conditions: when the reward received in the previous trial was “high” or
521 “low”, and when the transition experienced in the previous trial (relative to the trial before
522 that) “changed” or remained “fixed”. * $p < 0.05$, ** $p < 0.01$

523



524

525 **S2 Fig. Predictive weights from lagged logistic regression after second level (A-D) or**

526 **first level (E-H) contingency changes predicted by simulating model-based (A, E),**

527 **model-free (B, F), and hybrid (C, G) reinforcement learning algorithms, along with**

528 **experimental results (D, H). Lagged logistic regression measures the influence of different**

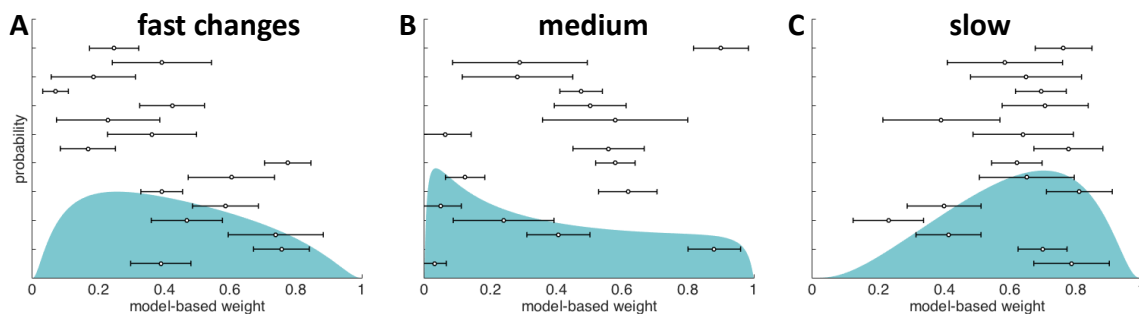
529 **predictors over several trials in the past, in this case, five trials. For each system, this index**

530 **was measured under four different conditions: when the reward received in the previous trial**

531 **was high (H) or low (L), and when the transition experienced in the previous trial (relative to**

532 **the trial before that) “changed” (C) or remained “fixed” (F).**

533



534

535 **S3 Fig. Model-based weights for different frequencies of contingency changes.**

536 **Probability density function over the model-based weight parameters estimated from model-**

537 **fitting, for the blocks of fast (every 3-6 trials), medium (every 7-10 trials) and slow (every**

538 11-14 trials) frequency of contingency changes. Overlaid are the individual subjects'
 539 parameter estimates for each block type. Error bars represent standard deviation.

540

541 **S1 Table. Median Plus Quartile Group-level Parameter Estimates.** Best-fitting parameter
 542 estimates over the subjects from model-fitting.

	β	Stay bias	α_{MB}	α_{MF}	λ	wMB (block 1)	wMB (block 2)	wMB (block 3)
1st quartile	1.88	0.04	0.55	0.03	0.30	0.10	0.40	0.48
Median	2.99	0.10	0.70	0.30	0.47	0.23	0.57	0.63
3rd quartile	4.73	0.22	0.81	0.85	0.65	0.46	0.71	0.76

543

544 **S2 Table. Model Comparison of Candidate Models.** Integrated Bayesian Information
 545 Criterion (iBIC) and negative log-likelihood of all candidate models from model-fitting. The
 546 models tested were: pure model-free (“MF”), pure model-based (“MB”), hybrid MB/MF
 547 (“hybrid”), hybrid MB/MF with different weights fitted for each of the three 200-trial blocks
 548 (“three-block hybrid”), and a hybrid model with different weights fitted for each frequency of
 549 contingency changes (“three-frequency hybrid”). The winning model was the three-block
 550 hybrid, highlighted in gray, according to iBIC and Bayesian model selection [25].

Model	Model-free	Model-based	Hybrid	Three-block hybrid	Three-frequency hybrid
Parameters	4	3	6	8	8
iBIC	9051	8091	7938	7687	7711
Negative Log Likelihood	4489	4018	3914	3770	3782

551

552 **S3 Table. Model Comparison of Additional Parameters.** Integrated Bayesian Information
 553 Criterion (iBIC) and negative log-likelihood of the winning three-block hybrid model with

554 different weights fitted for each of the three 200-trial blocks and the same model without stay
 555 bias, with $\lambda = 1$, with $\lambda = 0$, with only one learning rate for both MF and MB systems, and
 556 without updating fictive reward. The full model fit better to the data than the same model
 557 without each of the aforementioned parameters, even when controlling for model complexity
 558 in the iBIC.

Model	Full model	No stay bias	$\lambda = 1$	$\lambda = 0$	One learning rate	No fictive reward
Parameters	8	7	7	7	7	8
iBIC	7687	8047	7702	7754	7774	8024
Negative Log Likelihood	3770	3959	3787	3813	3823	3939

559