

1 **Human complex exploration strategies are enriched by noradrenaline-modulated heuristics**

2 Dubois M^{1,2}, Habicht J^{1,2}, Michely J^{1,2}, Moran R^{1,2}, Dolan RJ^{1,2} & Hauser TU^{1,2}

3 ¹Max Planck UCL Centre for Computational Psychiatry and Ageing Research, London WC1B
4 5EH, United Kingdom.

5 ²Wellcome Centre for Human Neuroimaging, University College London, London WC1N 3BG,
6 United Kingdom.

7 **Corresponding author**

8 Tobias U. Hauser

9 Max Planck UCL Centre for Computational Psychiatry and Ageing Research

10 University College London

11 10-12 Russell Square

12 London WC1B 5EH

13 United Kingdom

14 Phone: +44 / 207 679 5264

15 Email: t.hauser@ucl.ac.uk

16 Number of pages: 70

17 Number of Figures: 5

18 Number of Tables: 0

19 Abstract: 123 words

20 Introduction: 1036 words

21 Discussion: 2514 words

22 **Data and materials availability:** All necessary resources are publicly available at:

23 <https://github.com/MagDub>.

24
25
26
27
28
29
30
31
32
33
34
35

Abstract

An exploration-exploitation trade-off, the arbitration between sampling a lesser-known against a known rich option, is thought to be solved using computationally demanding exploration algorithms. Given known limitations in human cognitive resources, we hypothesised the presence of additional cheaper strategies. We examined for such heuristics in choice behaviour where we show this involves a value-free random exploration, that ignores all prior knowledge, and a novelty exploration that targets novel options alone. In a double-blind, placebo-controlled drug study, assessing contributions of dopamine (400mg amisulpride) and noradrenaline (40mg propranolol), we show that value-free random exploration is attenuated under the influence of propranolol, but not under amisulpride. Our findings demonstrate that humans deploy distinct computationally cheap exploration strategies and where value-free random exploration is under noradrenergic control.

36 **Introduction**

37 Chocolate, Toblerone, spinach or hibiscus ice-cream? Do you go for the flavour you like
38 the most (chocolate), or another one? In such an exploration-exploitation dilemma, you need to
39 decide whether to go for the option with the highest known subjective value (exploitation) or opt
40 instead for less known or valued options (exploration) so as to not miss out on possibly even
41 higher rewards. In the latter case, you can opt to either chose an option that you have previously
42 enjoyed (Toblerone), an option you are curious about because you do not know what to expect
43 (hibiscus), or even an option that you have disliked in the past (spinach). Depending on your
44 exploration strategy, you may end up with a highly disappointing ice cream encounter, or a life-
45 changing gustatory epiphany.

46 A common approach to the study of complex decision making, for example an
47 exploration-exploitation trade-off, is to take computational algorithms developed in the field of
48 artificial intelligence and test whether key signatures of these are evident in human behaviour.
49 This approach has revealed humans use strategies that reflect an implementation of
50 computationally demanding exploration algorithms (1, 2). One such strategy, directed
51 exploration, involves awarding an ‘information bonus’ to choice options, a bonus that scales with
52 uncertainty. This is captured in algorithms such as the Upper Confidence Bound (UCB) (3, 4)
53 and leads to an exploration of choice options the agent knows little about (1, 5) (e.g. the
54 hibiscus ice-cream). An alternative strategy, sometimes termed ‘random’ exploration, is to
55 induce stochasticity after value computations in the decision process. This can be realised using a
56 fixed parameter as a source of stochasticity, such as a softmax temperature parameter (6, 7),
57 which can be combined with the UCB algorithm (1). Alternatively, one can use a dynamic
58 source of stochasticity, such as in Thompson sampling (8), where stochasticity adapts to an

59 uncertainty about choice options. This exploration is essentially a more sophisticated,
60 uncertainty-driven, version of a softmax. By accounting for stochasticity when comparing choice
61 options' expected values, in effect choosing based on both uncertainty and value, these
62 exploration strategies increase the likelihood of choosing 'good' options that are only slightly
63 less valuable than the best (e.g. the Toblerone ice-cream if you are a chocolate lover).

64 The above processes are computationally demanding, especially when facing real-life
65 multiple-alternative decision problems (6, 9, 10). Human cognitive resources are constrained by
66 capacity limitations (11), metabolic consumption (12), but also because of resource allocation to
67 parallel tasks (e.g. (13, 14)). This directly relates to an agents' motivation to perform a given task
68 (11, 15, 16), as increasing an information demand in one process automatically reduces its
69 availability for others (12). In real-world highly dynamic environments, this arbitration is critical
70 as humans need to maintain resources for alternative opportunities (i.e. flexibility; (11, 17, 18)).
71 This accords with previous studies showing humans are demand-avoidant (17, 19) and suggests
72 that exploration computations tend to be minimised. Here, we examine the explanatory power of
73 two additional computationally less costly forms of exploration, namely value-free random
74 exploration and novelty exploration.

75 Computationally, the least resource demanding way to explore is to ignore all prior
76 information and to choose entirely randomly, de facto assigning the same probability to all
77 options. Such 'value-free' random exploration, as opposed to the two previously considered
78 'value-based' random explorations (for simulations comparing their effects cf. Figure 1 – Figure
79 supplement 2) that add stochasticity during choice value computation, forgoes any costly
80 computation (i.e. value mean and uncertainty), known as an ϵ -greedy algorithmic strategy in
81 reinforcement learning (20). Computational efficiency, however, comes at the cost of sub-

82 optimality due to occasional selection of options of low expected value (e.g. the repulsive
83 spinach ice cream).

84 Despite its sub-optimality, value-free random exploration has neurobiological
85 plausibility. Of relevance in this context is a view that exploration strategies depend on
86 dissociable neural mechanisms (21). Influences from noradrenaline and dopamine are plausible
87 candidates in this regard based on prior evidence (9, 22). Amongst other roles (such as memory
88 (23), or energisation of behaviour (24, 25)), the neuromodulator noradrenaline has been ascribed
89 a function of indexing uncertainty (26–28) or as acting as a ‘reset button’ that interrupts ongoing
90 information processing (29–31). Prior experimental work in rats shows boosting noradrenaline
91 leads to more value-free-random-like random behaviour (32), while pharmacological
92 manipulations in monkeys indicates reducing noradrenergic activity increases choice consistency
93 (33).

94 In human pharmacological studies, interpreting the specific function of noradrenaline on
95 exploration strategies is problematic as many drugs, such as atomoxetine (e.g. (34)), impact
96 multiple neurotransmitter systems. Here, to avoid this issue, we chose the highly specific β -
97 adrenoceptor antagonist propranolol, which has only minimal impact on other neurotransmitter
98 systems (35–37). Using this neuromodulator, we examine whether signatures of value-free
99 random exploration are impacted by administration of propranolol.

100 An alternative computationally efficient exploration heuristic to random exploration is to
101 simply choose an option not encountered previously, which we term novelty exploration.
102 Humans often show novelty seeking (38–41), and this strategy can be used in exploration as
103 implemented by a low-cost version of the UCB algorithm. Here a novelty bonus (42) is added if

104 a choice option has not been seen previously (i.e. it does not have to rely on precise uncertainty
105 estimates). The neuromodulator dopamine is implicated not only in exploration in general (43),
106 but also in signalling such types of novelty bonuses, where evidence indicates a role in
107 processing and exploring novel and salient states (39, 44–47). Although pharmacological
108 dopaminergic studies in humans have demonstrated effects on exploration as a whole (48), they
109 have not identified specific exploration strategies. Here, we used the highly specific D2/D3
110 antagonist, amisulpride, to disentangle the specific role of dopamine and noradrenaline on
111 different exploration strategies.

112 Thus, in the current study, we examine the contributions of value-free random
113 exploration and novelty exploration in human choice behaviour. We developed a novel
114 exploration task combined with computational modeling to probe the contributions of
115 noradrenaline and dopamine. Under double-blind, placebo-controlled, conditions we assessed the
116 impact of two antagonists with high affinity and specificity for either dopamine (amisulpride) or
117 noradrenaline (propranolol) respectively. Our results provide evidence that both exploration
118 heuristics supplement computationally more demanding exploration strategies, and that value-
119 free random exploration is particularly sensitive to noradrenergic modulation, with no effect of
120 amisulpride.

121 **Results**

122 *Probing the contributions of heuristic exploration strategies*

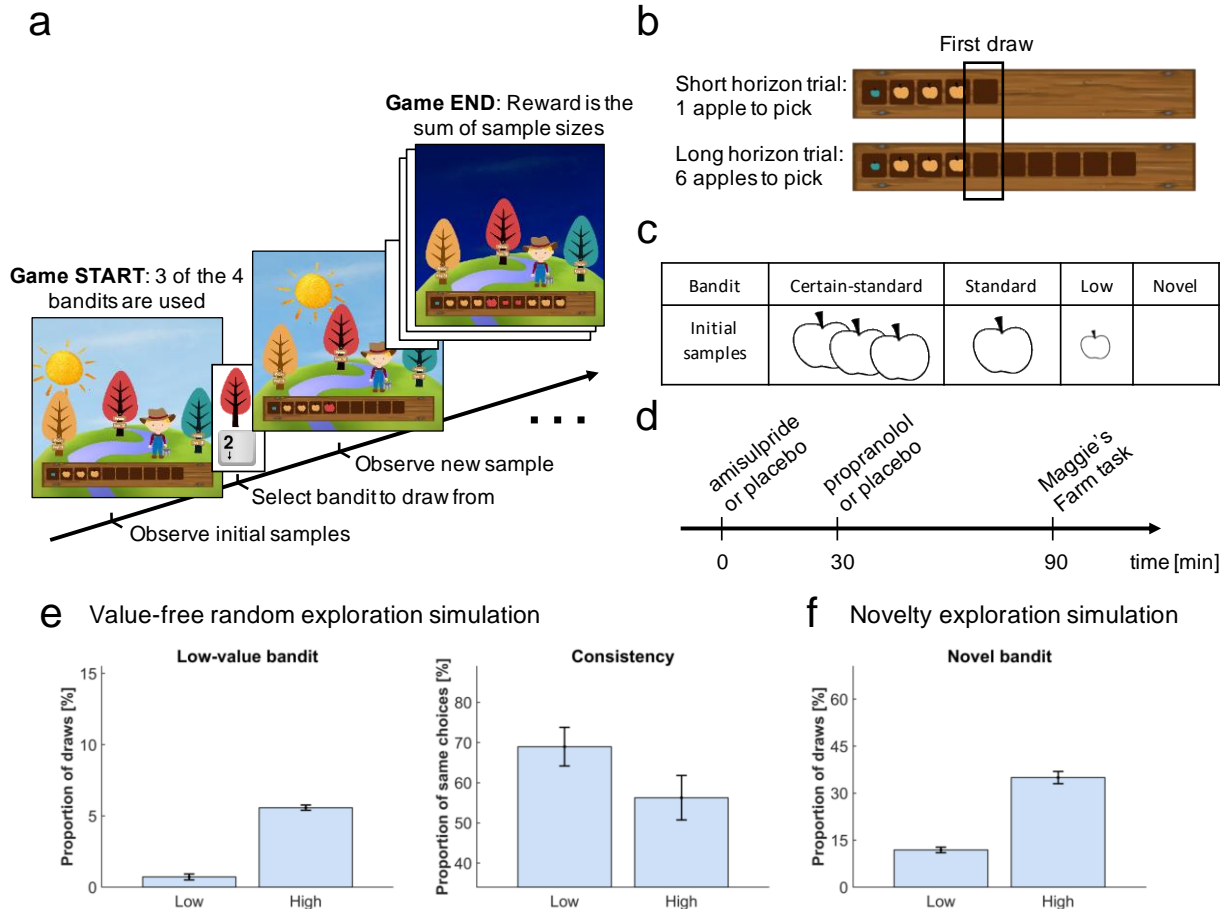
123 We developed a novel multi-round three-armed bandit task (Figure 1; bandits depicted as
124 trees), enabling us to assess the contributions of value-free random exploration and novelty
125 exploration in addition to Thompson sampling and UCB (combined with a softmax). In
126 particular, we exploited the fact that both heuristic strategies make specific predictions about
127 choice patterns. The novelty exploration assigns a ‘novelty bonus’ only to bandits for which
128 subjects have no prior information, but not to other bandits. This can be seen as a low-resolution
129 version of UCB, which assigns a bonus to all choice options proportionally to how informative
130 they are, in effect a graded bonus which scales to each bandits’ uncertainty. Thus, to capture this
131 heuristic, we manipulated the amount of prior information with bandits carrying only little
132 information (i.e. 1 vs 3 initial samples) or no information (0 initial samples). A high novelty
133 exploration predicts a higher frequency of selecting the novel option (Figure 1f). This is in
134 contrast to high exploration using other strategies which does not predict such a strong effect on
135 the novel option (cf. Figure 1 - Figure supplement 5).

136 Value-free random exploration, captured here by ϵ -greedy, predicts that all prior
137 information is discarded entirely and that there is equal probability attached to all choice options.
138 This strategy is distinct from other exploration strategies as it is likely to choose bandits known
139 to be substantially worse than the other bandits. Thus, a high value-free random exploration
140 predicts a higher frequency of selecting the low-value option (Figure 1e), whereas high
141 exploration using other strategies does not predict such effect (cf. Figure 1 - Figure supplement
142 3). A second prediction is that choice consistency, across repeated trials, is directly affected by
143 value-free random exploration, in particular by comparison to other more deterministic

144 exploration strategies (e.g. directed exploration) that are value-guided and thus will consistently
145 the most informative and valuable options. Given that value-free random exploration splits its
146 choice probability equally (i.e. 33.3% of choosing any bandit out of the three displayed), an
147 increase in such exploration predicts a lower likelihood of choosing the same bandit again, even
148 under identical choice options (Figure 1e). This contrasts to other strategies that make consistent
149 exploration predictions (e.g. UCB would consistently explore the choice option that carries a
150 high information bonus; Figure 1 - Figure supplement 4).

151 We generated bandits from four different generative processes (Figure 1c) with distinct
152 sample means (but a fixed sampling variance) and number of initial samples (i.e. samples shown
153 at the beginning of a trial for this specific bandit). Subjects were exposed to these bandits before
154 making their first draw. The ‘certain-standard bandit’ and the (less certain) ‘standard bandit’
155 were bandits with comparable means but varying levels of uncertainty, providing either three or
156 one initial samples (depicted as apples; similar to the horizon task (7)). The ‘low-value bandit’
157 was a bandit with one initial sample from a substantially lower generative mean, thus appealing
158 to a value-free random exploration strategy alone. The last bandit, with a mean comparable with
159 that of the standard bandits, was a ‘novel bandit’ for which no initial sample was shown,
160 primarily appealing to a novelty exploration strategy (cf. Materials and Methods for a full
161 description of bandit generative processes). To assess choice consistency, all trials were repeated
162 once. In the pilot experiments (data not shown), we noted some exploration strategies tended to
163 overshadow other strategies. To effectively assess all exploration strategies, we opted to present
164 only three of the four different bandit types on each trial, as different bandit triples allow
165 different explorations to manifest. Lastly, to assess whether subjects’ behaviour captured
166 exploration, we manipulated the degree to which subjects could interact with the same bandits.

167 Similar to previous studies (7), subjects could perform either one draw, encouraging exploitation
 168 (short horizon condition) or six draws encouraging more substantial explorative behaviour (long
 169 horizon condition) (7, 34).



170
 171 **Figure 1.** Study design. In the Maggie's farm task, subjects had to choose from three bandits
 172 (depicted as trees) to maximise an outcome (sum of reward). The rewards (apple size) of each
 173 bandit followed a normal distribution with a fixed sampling variance. (a) At the beginning of
 174 each trial, subjects were provided with some initial samples on the wooden crate at the bottom of
 175 the screen and had to select which bandit they wanted to sample from next. (b) Depending on the
 176 condition, they could either perform one draw (short horizon) or six draws (long horizon). The
 177 empty spaces on the wooden crate (and the suns' position) indicated how many draws they had
 178 left. The first draw in both conditions was the main focus of the analysis. (c) In each trial, three
 179 bandits were displayed, selected from four possible bandits, with different generative processes
 180 that varied in terms of their sample mean and number of initial samples (i.e. samples shown at
 181 the beginning of a trial). The 'certain-standard bandit' and the 'standard bandit' had comparable
 182 means but different levels of uncertainty about their expected mean: they provided three and one
 183 initial sample respectively; the 'low-value bandit' had a low mean and displayed one initial
 184 sample; the 'novel bandit' did not show any initial sample and its mean was comparable with

185 that of the standard bandits. (d) Prior to the task, subjects were administered different drugs:
186 400mg amisulpride that blocks dopaminergic D2/D3 receptors, 40mg propranolol to block
187 noradrenergic β -receptors, and inert substances for the placebo group. Different administration
188 times were chosen to comply with the different drug pharmacokinetics (placebo matching the
189 other groups' administration schedule). (e) Simulating value-free random behaviour with a low
190 vs high model parameter (ϵ) in this task shows that in a high regime, agents choose the low-value
191 bandit more often (left panel; mean \pm SD) and are less consistent in their choices when facing
192 identical choice options (right panel). (f) Novelty exploration exclusively promotes choosing
193 choice options for which subjects have no prior information, captured by the 'novel bandit' in
194 our task. For details about simulations cf. Materials and Methods. For details about the task
195 display cf. Figure 1 – Figure supplement 1. For simulations of different exploration strategies
196 and their impact of different bandits cf. Figure 1 – Figure supplement 2-5.

197

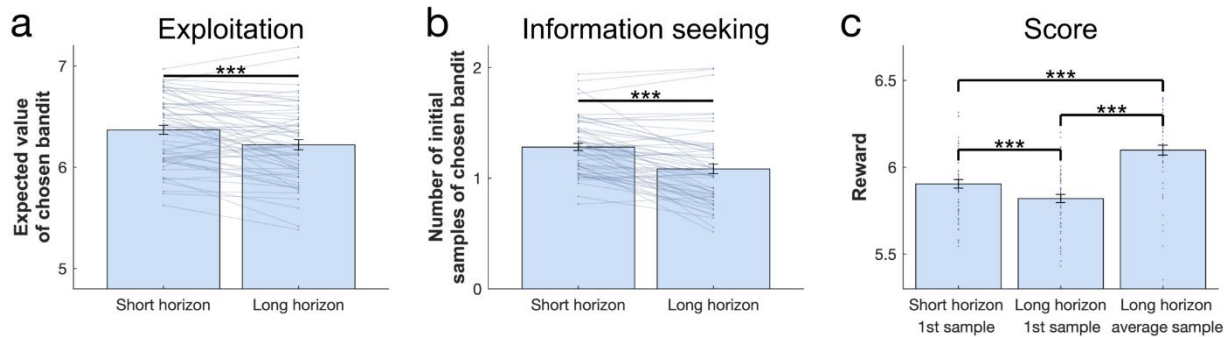
198 *Testing the role of catecholamines noradrenaline and dopamine*

199 In a double-blind, placebo-controlled, between-subjects, study design we assigned
200 subjects (N=60) randomly to one of three experimental groups: amisulpride, propranolol or
201 placebo. The first group received 40mg of the β -adrenoceptor antagonist propranolol to alter
202 noradrenaline function, while the second group was administered 400mg of the D2/D3
203 antagonist amisulpride that alters dopamine function. Because of different pharmacokinetic
204 properties, these drugs were administered at different times (Figure 1d) and compared to a
205 placebo group that received a placebo at both drug times to match the corresponding antagonists'
206 time. One subject (amisulpride group) was excluded from the analysis due to a lack of
207 engagement with the task. Reported findings were corrected for IQ and mood, as drug groups
208 differed marginally in those measures (cf. Appendix 2 Table 1), by adding WASI (49) and
209 PANAS (50) negative scores as covariates in each ANOVA. Similar results were obtained in an
210 analysis that corrected for physiological effects as from the analysis without covariates (cf.
211 Appendix 1).

212 *Increased exploration when information can subsequently be exploited*

213 Our task embodied two decision-horizon conditions, a short and a long. To assess whether
214 subjects explored more in a long horizon condition, in which additional information can inform
215 later choices, we examined which bandit subjects chose in their first draw (in accordance with
216 the horizon task (7)), irrespective of their drug group. A marker of exploration here is evident if
217 subjects chose bandits with lower expected values, computed as the mean value of their initial
218 samples shown (trials where the novel bandit was chosen were excluded). As expected, subjects
219 chose bandits with a lower expected value in the long compared to the short horizon (repeated-
220 measures ANOVA for the expected value: $F(1, 56)=19.457$, $p<.001$, $\eta^2=.258$; Figure 2a). To
221 confirm that this was a consequence of increased exploration, we analysed the proportion of how
222 often the high-value option was chosen (i.e. the bandit with the highest expected reward based on
223 its initial samples) and we found that subjects (especially those with higher IQ) sampled from it
224 more in the short compared to the long horizon, (WASI-by-horizon interaction: $F(1,54)=13.304$,
225 $p=.001$, $\eta^2=.198$; horizon main effect: $F(1, 54)=3.909$, $p=.053$, $\eta^2=.068$; Figure 3a), confirming
226 a reduction in exploitation when this information could be subsequently used. Interestingly, this
227 frequency seemed to be marginally higher in the amisulpride group, suggesting an overall higher
228 tendency to exploitation following dopamine blockade (cf. Appendix 1). This horizon-specific
229 behaviour resulted in a lower reward on the 1st sample in the long compared to the short horizon
230 ($F(1, 56)=23.922$, $p<.001$, $\eta^2=.299$; Figure 2c). When we tested whether subjects were more
231 likely to choose options they knew less about (computed as the mean number of initial samples
232 shown), we found that subjects chose less known (i.e. more informative) bandits more often in
233 the long horizon compared to the short horizon ($F(1, 56)=58.78$, $p<.001$, $\eta^2=.512$; Figure 2b).

234 Next, to evaluate whether subjects used the additional information beneficially in the long
235 horizon condition, we compared the average reward (across six draws) obtained in the long
236 compared to short horizon (one draw). We found that the average reward was higher in the long
237 horizon ($F(1, 56)=103.759$, $p<.001$, $\eta^2=.649$; Figure 2c), indicating that subjects tended to
238 choose less optimal bandits at first but subsequently learnt to appropriately exploit the harvested
239 information to guide choices of better bandits in the long run. Additionally, when looking
240 specifically at the long horizon condition, we found that subjects earned more when their first
241 draw was explorative versus exploitative (Figure 2 - Figure supplement 1c-d; cf. Appendix 2 for
242 details).



244 **Figure 2.** Benefits of exploration. To investigate the effect of information on performance we
 245 collapsed subjects over all three treatment groups. (a) The expected value (average of its initial
 246 samples) of the first chosen bandit as a function of horizon. Subjects chose bandits with a lower
 247 expected value (i.e. they explored more) in the long horizon compared to the short horizon. (b)
 248 The mean number of samples for the first chosen bandit as a function of horizon. Subjects chose
 249 less known (i.e. more informative) bandits more in the long compared to the short horizon. (c)
 250 The first draw in the long horizon led to a lower reward than the first draw in the short horizon,
 251 indicating that subjects sacrificed larger initial outcomes for the benefit of more information.
 252 This additional information helped making better decisions in the long run, leading to a higher
 253 earning over all draws in the long horizon. For values and statistics cf. Appendix 2 Table 3. For
 254 response times and details about all long horizons' samples cf. Figure 2 – Figure supplement 1.
 255 *** = $p < .001$. Data are shown as mean \pm SEM and each dot/line represent a subject.
 256

257

258 *Subjects demonstrate value-free random behaviour*

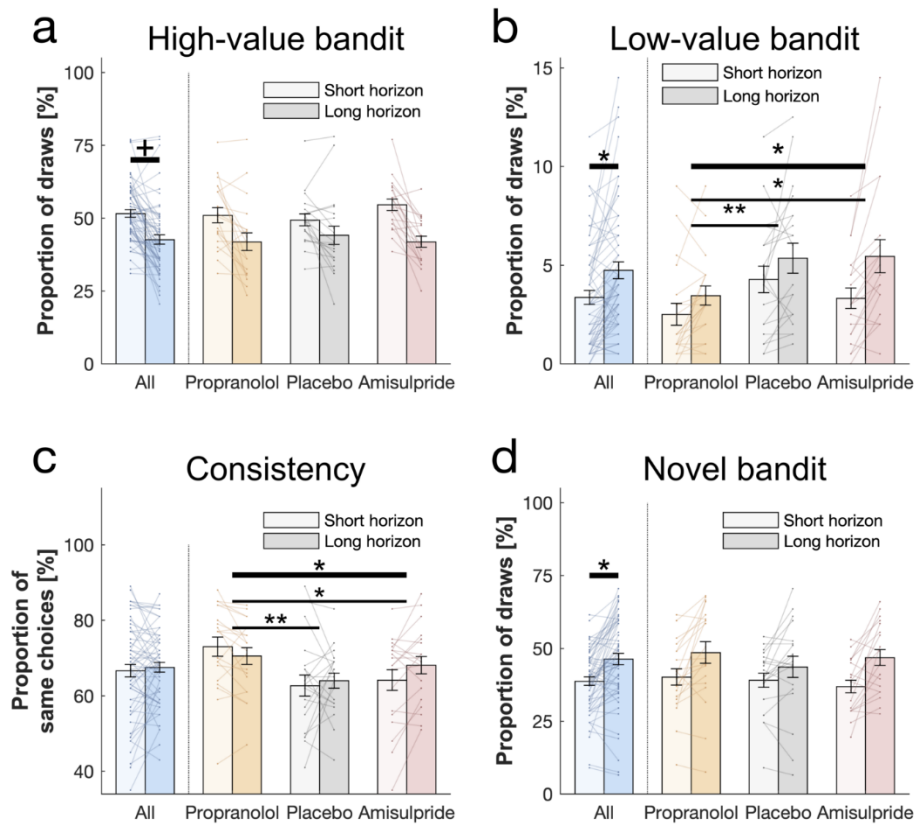
259 Value-free random exploration (analogue to ϵ -greedy) predicts that ϵ % of the time each
 260 option will have an equal probability of being chosen. In such a regime (compared to more
 261 complex strategies that would favour options with a higher expected value with a similar
 262 uncertainty), the probability of choosing bandits with a low expected value (here the low-value
 263 bandit; Fig. 1e) will be higher (cf. Figure 1 – Figure supplement 3). We investigated whether the
 264 frequency of picking the low-value bandit was increased in the long horizon condition across all
 265 subjects (i.e. when exploration is useful), and we found a significant main effect of horizon ($F(1,$
 266 $54)=4.069$, $p=.049$, $\eta^2=.07$; Figure 3b). This demonstrates that value-free random exploration is
 267 utilised more when exploration is beneficial.

268 *Value-free random behaviour is modulated by noradrenaline function*

269 When we tested whether value-free random exploration was sensitive to neuromodulatory
270 influences, we found a difference in how often drug groups sampled from the low-value option
271 (drug main effect: $F(2, 54)=7.003$, $p=.002$, $\eta^2=.206$; drug-by-horizon interaction: $F(2,$
272 $54)=2.154$, $p=.126$, $\eta^2=.074$; Figure 3b). This was driven by the propranolol group choosing the
273 low-value option significantly less often than the other two groups (placebo vs propranolol:
274 $t(40)=2.923$, $p=.005$, $d=.654$; amisulpride vs propranolol: $t(38)=2.171$, $p=.034$, $d=.496$) with no
275 difference between amisulpride and placebo: ($t(38)=-0.587$, $p=.559$, $d=.133$). These findings
276 demonstrate that a key feature of value-free random exploration, the frequency of choosing low-
277 value bandits, is sensitive to influences from noradrenaline.

278 To further examine drug effects on value-free random exploration, we assessed a second
279 prediction, namely choice consistency. Because value-free random exploration ignores all prior
280 information and chooses randomly, it should result in a decreased choice consistency when
281 presented identical choice options (cf. Figure 1 – Figure supplement 2 & 4, compared to more
282 complex strategies which are always biased towards the rewarding or the information providing
283 bandit for example). To this end, each trial was duplicated in our task, allowing us to compute
284 the consistency as the percentage of time subjects sampled from an identical bandit when facing
285 the exact same choice options. In line with the above analysis, we found a difference in
286 consistency by which drug groups sampled from different option (drug main effect: $F(2,$
287 $54)=7.154$, $p=.002$, $\eta^2=.209$; horizon main effect: $F(1, 54)=1.333$, $p=.253$, $\eta^2=.024$; drug-by-
288 horizon interaction: $F(2, 54)=3.352$, $p=.042$, $\eta^2=.11$; Figure 3c), driven by the fact that the
289 propranolol group chose significantly more consistently than the other two groups (pairwise
290 comparisons: placebo vs propranolol: $t(40)=-3.525$, $p=.001$, $d=.788$; amisulpride vs placebo:

291 $t(38)=1.107$, $p=.272$, $d=.251$; amisulpride vs propranolol: $t(38)=-2.267$, $p=.026$, $d=.514$). Please
 292 see Appendix 1 for further discussion and analysis of the drug-by-horizon interaction. Taken
 293 together, these results indicate that value-free random exploration depends critically on
 294 noradrenaline functioning, such that an attenuation of noradrenaline leads to a reduction in value-
 295 free random exploration.



296

297 **Figure 3.** Behavioural horizon and drug effects. Choice patterns in the first draw for each
 298 horizon and drug group (propranolol, placebo and amisulpride). (a) Subjects sampled from the
 299 high-value bandit (i.e. bandit with the highest average reward of initial samples) more in the
 300 short horizon compared to the long horizon indicating reduced exploitation. (b) Subjects sampled
 301 from the low-value bandit more in the long horizon compared to the short horizon indicating
 302 value-free random exploration, but subjects in the propranolol group sampled less from it
 303 overall, and (c) were more consistent in their choices overall, indicating that noradrenaline
 304 blockade reduces value-free random exploration. (d) Subjects sampled from the novel bandit
 305 more in the long horizon compared to the short horizon indicating novelty exploration. Please
 306 note that some horizon effects were modulated by subjects' intellectual abilities when
 307 additionally controlling for them (cf. Appendix 2 Table 4). Horizontal bars represent rm-
 308 ANOVA (thick) and pairwise comparisons (thin). † = $p < .07$, * = $p < .05$, ** = $p < .01$. Data are

309 shown as mean \pm SEM and each line represent one subject. For values and statistics cf. Appendix
310 2 Table 4. For response times and frequencies specific to the displayed bandits cf. Figure 3 –
311 Figure supplement 1-2.

312

313 *Novelty exploration is unaffected by catecholaminergic drugs*

314 Next, we examined whether subjects show evidence for novelty exploration by choosing
315 the novel bandit for which there was no prior information (i.e. no initial samples), as predicted
316 by model simulations (Figure 1f). We found a significant main effect of horizon ($F(1, 54)=5.593$,
317 $p=.022$, $\eta^2=.094$; WASI-by-horizon interaction: $F(1, 54) =13.897$, $p<.001$, $\eta^2=.205$; Figure 3d)
318 indicating that subjects explored the novel bandit significantly more often in the long horizon
319 condition, and this was particularly strong for subjects with a higher IQ. We next assessed
320 whether novelty exploration was sensitive to our drug manipulation, but found no drug effects on
321 the novel bandit ($F(2, 54)=1.498$, $p=.233$, $\eta^2=.053$; drug-by-horizon interaction: $F(2, 54)=.542$,
322 $p=.584$, $\eta^2=.02$; Figure 3d). Thus, there was no evidence that an attenuation of dopamine or
323 noradrenaline function impact novelty exploration in this task.

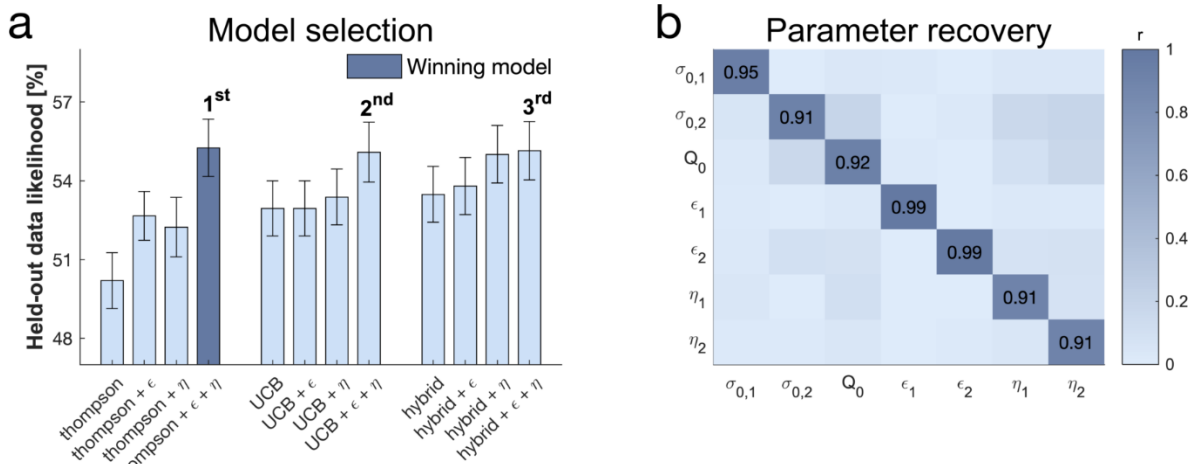
324 *Subjects combine computationally demanding strategies and exploration heuristics*

325 To examine the contributions of different exploration strategies to choice behaviour, we
326 fitted a set of computational models to subjects' behaviour, building on models developed in
327 previous studies (1). In particular, we compared models incorporating UCB, Thompson
328 sampling, an ϵ -greedy algorithm and the novelty bonus (cf. Materials and Methods). Essentially,
329 each model makes different exploration predictions. In the Thompson model, Thompson
330 sampling (8, 51) leads to an uncertainty-driven value-based random exploration, where both
331 expected value and uncertainty contribute to choice. In this model higher uncertainty leads to
332 more exploration such that instead of selecting a bandit with the highest mean, bandits are

333 chosen relative to how often a random sample would yield the highest outcome, thus accounting
334 for uncertainty (2). The UCB model (3, 4), capturing directed exploration, predicts that each
335 bandit is chosen according to a mixture of expected value and an additional expected information
336 gain (2). This is realised by adding a bonus to the expected value of each option, proportional to
337 how informative it would be to select this option (i.e. the higher the uncertainty in the options'
338 value, the higher the information gain). This computation is then passed through a softmax
339 decision model, capturing value-based random exploration. Novelty exploration is a simplified
340 version of the information bonus in the UCB algorithm, which only applies to entirely novel
341 options. It defines the intrinsic value of selecting a bandit about which nothing is known, and
342 thus saves demanding computations of uncertainty for each bandit. Lastly, the value-free random
343 ϵ -greedy algorithm selects any bandit ϵ % of the time, irrespective of the prior information of this
344 bandit. For additional models cf. Appendix 1.

345 We used cross-validation for model selection (Figure 4a) by comparing the likelihood of
346 held-out data across different models, an approach that adequately arbitrates between model
347 accuracy and complexity. The winning model encompasses uncertainty-driven value-based
348 random exploration (Thompson sampling) with value-free random exploration (ϵ -greedy
349 parameter) and novelty exploration (novelty bonus parameter η). The winning model predicted
350 held-out data with a 55.25% accuracy (SD=8.36%; chance level =33.33%). Similarly to previous
351 studies (1), the hybrid model combining UCB and Thompson sampling explained the data better
352 than each of those processes alone, but this was no longer the case when accounting for novelty
353 and value-free random exploration (Figure 4a). The winning model further revealed that all
354 parameter estimates could be accurately recovered (Figure 4b; Figure 4 – Figure supplement 3).
355 Interestingly, although the 2nd and 3rd place models made different prediction about the complex

356 exploration strategy, using a directed exploration with value-based random exploration (UCB) or
 357 a combination of complex strategies (hybrid) respectively, they share the characteristic of
 358 benefitting from value-free random and novelty exploration. This highlights that subjects used a
 359 mixture of computationally demanding and heuristic exploration strategies.



360
 361 **Figure 4.** Subjects use a mixture of exploration strategies. (a) A 10-fold cross-validation of the
 362 likelihood of held-out data was used for model selection (chance level =33.3%; for model
 363 selection at the individual level cf. Figure 4 – Figure supplement 1). The Thompson model with
 364 both the ϵ -greedy parameter and the novelty bonus η best predicted held-out data (b) Model
 365 simulation with 4^7 simulations predicted good recoverability of model parameters (for
 366 correlations between behaviour and model parameters cf. Figure 4 – Figure supplement 2); σ_0
 367 is the prior variance and Q_0 is the prior mean (for parameter recovery correlation plots cf. Figure 4
 368 – Figure supplement 3). 1 stands for short horizon-, and 2 for long horizon-specific parameters.
 369 For values and parameter details cf. Appendix 2 Table 5.

370

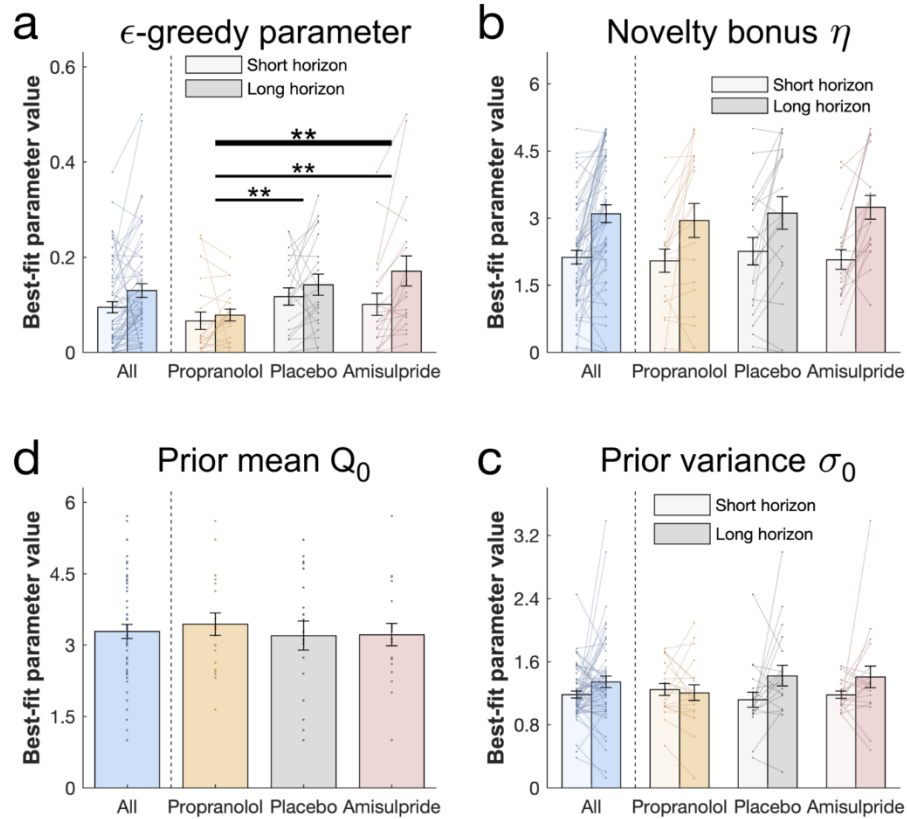
371 *Noradrenaline controls value-free random exploration*

372 To more formally compare the impact of catecholaminergic drugs on different exploration
 373 strategies, we assessed the free parameters of the winning model between drug groups (Figure 5,
 374 cf. Appendix 2 Table 6 for exact values). First, we examined the ϵ -greedy parameter that
 375 captures the contribution of value-free random exploration to choice behaviour. We assessed
 376 how this value-free random exploration differed between drug groups. A significant drug main
 377 effect (drug main effect: $F(2, 54)=6.722, p=.002, \eta^2=.199$; drug-by-horizon interaction: $F(2,$

378 54)=1.305, $p=.28$, $\eta^2=.046$; Figure 5a) demonstrates that the drug groups differ in how strongly
379 they deploy this exploration strategy. Post-hoc analysis revealed that subjects with reduced
380 noradrenaline functioning had the lowest values of ϵ (pairwise comparisons: placebo vs
381 propranolol: $t(40)=3.177$, $p=.002$, $d=.71$; amisulpride vs propranolol: $t(38)=2.723$, $p=.009$,
382 $d=.626$) with no significant difference between amisulpride vs placebo: ($t(38)=.251$, $p=.802$,
383 $d=.057$). Critically, the effect on ϵ was also significant when the complex exploration strategy
384 was a directed exploration with value-based random exploration (2nd place model) and,
385 marginally significant, when it was a combination of the above (3rd place model; cf. Appendix
386 1).

387 The ϵ -greedy parameter was also closely linked to the above behavioural metrics
388 (correlation between the ϵ -greedy parameter with draws from the low-value bandit:
389 $R_{Pearson}=.828$, $p<.001$; and with choice consistency: $R_{Pearson}=-.596$, $p<.001$; Figure 4 – Figure
390 supplement 2), and showed a similar horizon effect (horizon main effect: $F(1, 54)=1.968$,
391 $p=.166$, $\eta^2=.035$; WASI-by-horizon interaction: $F(1, 54)=6.08$, $p=.017$, $\eta^2=.101$; Figure 5a). Our
392 findings thus accord with the model-free analyses and demonstrate that noradrenaline blockade
393 reduces value-free random exploration.

394



396

397 **Figure 5.** Drug effects on model parameters. The winning model's parameters were fitted to
 398 each subject's first draw (for model simulations cf. Figure 5 – Figure supplement 1). (a) Subjects
 399 had higher values of ϵ (value-free random exploration) in the long compared to the short horizon.
 400 Notably, subjects in the propranolol group had lower values of ϵ overall, indicating that
 401 attenuation of noradrenaline functioning reduces value-free random exploration. Subjects from
 402 all groups (b) assigned a similar value to novelty, captured by the novelty bonus η , which was
 403 higher (more novelty exploration) in the long compared to the short horizon. (c) The groups had
 404 similar beliefs Q_0 about a bandits' mean before seeing any initial samples and (d) were similarly
 405 uncertain σ_0 about it (for gender effects cf. Figure 5 – Figure supplement 2). Please note that
 406 some horizon effects were modulated by subjects' intellectual abilities when additionally
 407 controlling for them (cf. Appendix 2 Table 6). ** = $p < .01$. Data are shown as mean \pm SEM and
 408 each dot/line represent one subject. For parameter values and statistics cf. Appendix 2 Table 6.

409

410 *No drug effects on other parameters*

411 The novelty bonus η captures the intrinsic reward of selecting a novel option. In line with
 412 the model-free behavioural findings, there was no difference between drug groups in terms of

413 this effect ($F(2, 54)=.249$, $p=.78$, $\eta^2=.009$; drug-by-horizon interaction: $F(2, 54)=.03$, $p=.971$,
414 $\eta^2=.001$). There was also a close alignment between model-based and model-agnostic analyses
415 (correlation between the novelty bonus η with draws from the novel bandit: $R_{Pearson}=.683$,
416 $p<.001$; Figure 4 – Figure supplement 2), and we found a similarly increased novelty bonus
417 effect in the long horizon in subjects with a higher IQ (WASI-by-horizon interaction: $F(1, 54)$
418 $=8.416$, $p=.005$, $\eta^2=.135$; horizon main effect: $F(1, 54)=1.839$, $p=.181$, $\eta^2=.033$; Figure 5b).

419 When analysing the additional model parameter, we found that subjects had similar prior
420 beliefs about bandits, given by the initial estimate of a bandit's mean (prior mean Q_0 : $F(2,$
421 $54)=.118$, $p=.889$, $\eta^2=.004$; Figure 5c) and their uncertainty about it (prior variance σ_0 : horizon
422 main effect: $F(1, 54)=.129$, $p=.721$, $\eta^2=.002$; drug main effect: $F(2, 54)=.06$, $p=.942$, $\eta^2=.002$;
423 drug-by-horizon interaction: $F(2, 54)=2.162$, $p=.125$, $\eta^2=.074$; WASI-by-horizon interaction:
424 $F(1, 54)=.022$, $p=.882$, $\eta^2<.001$; Figure 5d). Interestingly, our dopamine manipulation seemed to
425 affect this uncertainty in a gender-specific manner, with female subjects having larger values of
426 σ_0 compared to males in the placebo group, and with the opposite being true in the amisulpride
427 group (cf. Appendix 1). Taken together, these findings show that value-free random exploration
428 was most sensitive to our drug manipulations.

429 Discussion

430 Solving the exploration-exploitation problem is non trivial, and one suggestion is that
431 humans solve it using computationally demanding exploration strategies (1, 2), taking account of
432 the uncertainty (variance) as well as the expected reward (mean) of each choice. Although
433 tracking the distribution of summary statistics (e.g. mean and variance) is less resource costly
434 than keeping track of full distributions (52), it nevertheless carries considerable costs when one
435 has to keep track of multiple options, as in exploration. Indeed, in a three-bandit task such as that
436 considered here, this results in a necessity to compute 6 key-statistics, drastically limiting
437 computational resources when selecting among choice options (10). Real-life decisions often
438 comprise an unlimited range of options, which results in a tracking of a multitude of key-
439 statistics, potentially mandating a deployment of alternative more efficient strategies. Here, we
440 demonstrate that two additional, less resource-hungry heuristics are at play during human
441 decision-making, value-free random exploration and novelty exploration.

442 By assigning intrinsic value (novelty bonus (42)) to an option not encountered before
443 (53), a novelty bonus can be seen as an efficient simplification of demanding algorithms, such as
444 UCB (3, 4). It is interesting to note that our winning model did not include UCB, but instead
445 novelty exploration. This indicates humans might use such a novelty shortcut to explore unseen,
446 or rarely visited, states to conserve computational costs when such a strategy is possible. A
447 second exploration heuristic that also requires minimal computational resources, value-free
448 random exploration, also plays a role in our task. Even though less optimal, its simplicity and
449 neural plausibility renders it a viable strategy. Indeed, we observe an increase in performance in
450 each model after adding ϵ , supporting the notion that this strategy is a relevant additional human
451 exploration heuristic. Interestingly, the benefit of ϵ is somewhat smaller in a simple UCB model

452 (without novelty bonus), which probably arises because value-based random exploration
453 partially captures some of the increased noisiness. We show through converging behavioural and
454 modelling measures that both value-free random and novelty exploration were deployed in a
455 goal-directed manner, coupled with increased levels of exploration when this was strategically
456 useful. Importantly, these heuristics were observed in all best models (1st, 2nd and 3rd position)
457 even though each incorporated different exploration strategies. This suggests that the complex
458 models make similar predictions in our task. This is also observed in our simulations, and
459 demonstrates that value-free random exploration is at play even when accounting for other value-
460 based forms of random exploration (1, 7), whether fixed or uncertainty-driven.

461 Exploration was captured in a similar manner to previous studies (7), by comparing in the
462 same setting (i.e. same prior information) the first choice in a long decision horizon, where
463 reward can be increased in the long term through information gain, and in a short decision
464 horizon where information cannot subsequently be put to use. This means that by changing the
465 opportunity to benefit from the information gained for the first sample, the long horizon invites
466 extended exploration (7), what we find also in our study. This experimental manipulation is a
467 well-established means for altering exploration and has been used extensively in previous studies
468 (7, 21, 34, 54). Nevertheless, there remains a possibility that a longer horizon may also affect the
469 psychological nature of the task. In our task, reward outcomes were presented immediately after
470 every draw, rendering it unlikely that perception of reward delays (i.e. delay discounting) is
471 impacted. Moreover, a monetary bonus was given only at the end of the task, and thus did not
472 impact a horizon manipulation. We also consider our manipulation was unlikely to change effort
473 in each horizon, because the reward (i.e. size of the apple) remains the same at every draw,

474 resulting in an equivalent reward-effort ratio (55–58). However, this issue can be addressed in
475 further studies, for example, by equating the amount of button presses across both conditions.

476 Value-free random exploration might reflect other influences, such as attentional lapses
477 or impulsive motor responses. We consider these as unlikely to a significant factor at play here.
478 Indeed, there are two key features that would signify such effects. Firstly, these influences would
479 be independent of task condition. Secondly, they would be expected to lead to shorter, or more
480 variable, response latencies. In our data, we observe an increase in value-free exploration in the
481 long horizon condition in both behavioural measures and model parameters, speaking against an
482 explanation based upon simple mistakes. Moreover, we did not observe a difference in response
483 latency for choices that were related to value-free random exploration (cf. Appendix 1), further
484 arguing against mistakes. Lastly, the sensitivity of value-free random exploration to propranolol
485 supports this being a separate process, and previous studies using the same drug did not find an
486 effect on task mistakes (e.g. on accuracy (59); (33, 58–60)). However, future studies could
487 explore these exploration strategies in more detail including by reference to subjects' own self-
488 reports.

489 It is still unclear how exploration strategies are implemented neurobiologically.
490 Noradrenaline inputs, arising from the locus coeruleus (63) (LC) are thought to modulate
491 exploration (2, 64, 65), though empirical data on its precise mechanisms and means of action
492 remains limited. In this study, we found that noradrenaline impacted value-free random
493 exploration, in contrast to novelty exploration and complex exploration. This might suggest that
494 noradrenaline influences ongoing valuation or choice processes that discards prior information.
495 Importantly, this effect was observed whether the complex exploration was an uncertainty-driven
496 value-based random exploration (winning model), a directed exploration with value-based

497 random exploration (2nd place model) or a combination of the above (3rd place model; cf.
498 Appendix 1). This is consistent with findings in rodents where enhanced anterior cingulate
499 noradrenaline release leads to more random behaviour (32). It is also consistent with
500 pharmacological findings in monkeys that show enhanced choice consistency after reducing LC
501 noradrenaline firing rates (33). It would be interesting for future studies to determine, in more
502 detail, whether value-free random exploration is corrupting a value computation itself, or
503 whether it exclusively biases the choice process.

504 We note that pupil diameter has been used as an indirect marker of noradrenaline activity
505 (66), although the link between the two is not always straightforward (36). Because the effect of
506 pharmacologically induced changes of noradrenaline levels on pupil size remains poorly
507 understood (36, 67), including the fact that previous studies found no effect of propranolol on
508 pupil diameter (36, 68), we opted against using pupillometry in this study. However, our current
509 findings align with previous human studies that show an association between this indirect marker
510 and exploration, but that study did not dissociate between the different potential exploration
511 strategies that subjects could deploy (69). Future studies might usefully include indirect
512 measures of noradrenaline activity, for example pupillometry, to examine a potential link
513 between natural variations in noradrenaline levels and a propensity towards value-free random
514 exploration.

515 The LC has two known modes of synaptic signalling (63), tonic and phasic, thought to
516 have complementary roles (31). Phasic noradrenaline is thought to act as a reset button (31),
517 rendering an agent agnostic to all previously accumulated information, a de facto signature of
518 value-free random exploration. Tonic noradrenaline has been associated, although not
519 consistently (70), with increased exploration (64, 71), decision noise in rats (72) and more

520 specifically with random as opposed to directed exploration strategies (34). This later study
521 unexpectedly found that boosting noradrenaline decreased (rather than increased) random
522 exploration, which the authors speculated was due to an interplay with phasic signalling.
523 Importantly, the drug used in that study also affects dopamine function making it difficult to
524 assign a precise interpretation to the finding. A consideration of this study influenced our
525 decision to opt for drugs with high specificity for either dopamine or noradrenaline (59),
526 enabling us to reveal highly specific effects on value-free random exploration. Although the
527 contributions of tonic and phasic noradrenaline signalling cannot be disentangled in our study,
528 our findings align with theoretical accounts and non-primate animal findings, indicating that
529 phasic noradrenaline promotes value-free random exploration.

530 Aside from this ‘reset signal’ role, noradrenaline has been assigned other roles, including
531 a role in memory function (23, 73, 74). To minimise a possible memory-related impact, we
532 designed the task such that all necessary information was visible on the screen at all times. This
533 means subjects did not have to memorise values for a given trial, rendering the task less
534 susceptible to forgetting or other memory effects. Another role for noradrenaline relates to
535 volatility and uncertainty estimation (26–28), as well as the energisation of behaviour (24, 25).
536 Non-human primates studies demonstrate a higher LC activation for high effort choices,
537 suggesting that noradrenaline release facilitates energy mobilisation (24). Theoretical models
538 also suggest that the LC is involved in the control of effort exertion. Thus, it is thought to
539 contribute to trading off between effortful actions leading to large rewards and “effortless”
540 actions leading to small rewards by modulating “raw” reward values as a function of the required
541 effort (25). Our task can be interpreted as encapsulating such a trade-off: complex exploration
542 strategies are effortful but optimal in terms of reward gain, while value-free random exploration

543 requires little effort while occasionally leading to low reward. Applying this model, a
544 noradrenaline boost could optimise cognitive effort allocation for high reward gain (25), thereby
545 facilitating complex exploration strategies compared to value-free random exploration. In such a
546 framework, blocking noradrenaline release should decrease usage of complex exploration
547 strategies, leading to an increase of value-free random exploration which is the opposite of what
548 we observed in our data. Another interpretation of an effort-facilitation model of noradrenaline is
549 that a boost would help overcoming cost, i.e. the lack of immediate reward when selecting the
550 low-value bandit, essentially providing a significant increase to the value of information gain. In
551 line with our results, a decrease would interrupt this boost in valuation, removing an incentive to
552 choose the low-value option. However, this theory is currently limited by the absence of
553 empirical evidence for noradrenaline boosting valuation.

554 Noradrenaline blockade by propranolol has been shown previously to enhance
555 metacognition (75), decrease information gathering (59), and attenuate arousal-induced boosts in
556 incidental memory (36). All of these findings, including a decrease in value-free random
557 exploration found here, suggests propranolol may influence how neural noise affects information
558 processing. In particular, the results indicate that under propranolol behaviour is less stochastic
559 and less influenced by ‘task-irrelevant’ distractions. This aligns with theoretical ideas, as well as
560 recent optogenetic evidence (32), that proposes noradrenaline infuses noise in a temporally
561 targeted way (31). It also accords with studies implicating noradrenaline in attention shifts (for a
562 review cf. (76)). Other gain-modulation theories of noradrenaline/catecholamine function have
563 proposed an effect on stochasticity (64, 65), although a hypothesized direction of effect is
564 different (i.e. noradrenaline decreases stochasticity). Several aspects of noradrenaline functioning
565 may explain the contradictory accounts of its link with stochasticity. For example, they might be

566 capturing different aspects of an assumed U-shaped noradrenaline functioning curve, and/or
567 distinct activity modes of noradrenaline (i.e. tonic and phasic firing) (64). Further studies can
568 shed light on how different modes of activity affect value-free random exploration. This idea can
569 be extended also to tasks where propranolol has been shown to attenuate a discrimination
570 between different levels of loss (with no effect on the value-based exploration parameter,
571 referred to in these studies as consistency) (62) and a reduction in loss aversion (60). This hints
572 at additional roles for noradrenaline on prior information and task-distractibility during
573 exploration in loss frame environments. Future studies investigating exploration in loss contexts
574 might provide important additional information on these questions.

575 It is important to mention here that β -adrenergic receptors, the primary target of
576 propranolol, have been shown (unlike α -adrenergic receptors) to increase synaptic inhibition
577 within rat cortex (77), specifically through inhibitory GABA-mediated transmission (78).
578 Additionally β -adrenergic receptors are more concentrated in the intermediate layers in the
579 prefrontal area (79), within which inhibition is favoured (80). Thus inhibitory mechanisms might
580 account for noradrenaline-related task-distractibility and randomness, or the role of β -adrenergic
581 receptors in executive function impairments (81). This raises the question of whether blocking β -
582 adrenergic receptors might lead to an accumulation of synaptic noradrenaline, and therefore act
583 via α -adrenergic receptors. To the best of our knowledge, evidence for such an effect is limited.
584 A second question is whether the observed effects are a pure consequence of propranolol's
585 impact on the brain, or whether they reflect peripheral effects of propranolol. When we examined
586 peripheral markers (i.e. heart rate) and behaviour we found no evidence for an effect on any of
587 our findings, rendering such influences unlikely. However, future studies using drugs that

588 exclusively targets peripheral, but not central, noradrenaline receptors (e.g. (82)) are needed to
589 answer this question conclusively.

590 Dopamine has been ascribed multiple functions besides reward learning (83), such as
591 novelty seeking (46, 84, 85) or exploration in general (43). In fact, studies have demonstrated
592 that there are different types of dopaminergic neurons in the ventral tegmental area, and that
593 some contribute to non-reward signals, such as saliency and novelty (44). This suggests a role in
594 novelty exploration. Moreover, dopamine has been suggested as important in an exploration-
595 exploitation arbitration (21, 86, 87), although its precise role remains unclear, given reported
596 effects on random exploration (88), on directed exploration (45, 89), or no effects at all (90). A
597 recent study found no effect following dopamine blockade using haloperidol (87), which
598 interestingly also affects noradrenaline function (e.g. (91, 92)). Our results did not demonstrate
599 any main effect of dopamine manipulation on exploration strategies, even though blocking
600 dopamine was associated with a trend level increase in exploitation (cf. Appendix 1). We believe
601 it unlikely this reflects an ineffective drug dose as previous studies have found neurocognitive
602 effects with the same dose (36, 59, 93, 94).

603 One possible reason for an absence of significant findings is that our dopaminergic
604 blockade targets D2/D3 receptors rather than D1 receptors, a limitation due a lack of available
605 specific D1 receptor blockers for use in humans. An expectation of greater D1 involvement
606 arises out of theoretical models (95) and a prefrontal hypothesis of exploration (89).
607 Interestingly, we observed a weak gender-specific differential drug effect on subjects'
608 uncertainty about an expected reward, with women being more uncertain than men in the
609 placebo setting, but more certain in the dopamine blockade setting (cf. Appendix 1). This might
610 be meaningful as other studies using the same drug have also found behavioural gender-specific

611 drug effects (96). Upcoming, novel drugs (97) might be able help unravel a D1 contribution to
612 different forms of exploration. Additionally, future studies could use approved D2/D3 agonists
613 (e.g. ropinirole) in a similar design to probe further whether enhancing dopamine leads to a
614 general increase in exploration.

615 In conclusion, humans supplement computationally expensive exploration strategies with
616 less resource demanding exploration heuristics, and as shown here the latter include value-free
617 random and novelty exploration. Our finding that noradrenaline specifically influences value-
618 free random exploration demonstrates that distinct exploration strategies may be under specific
619 neuromodulator influence. Our current findings may also be relevant to enabling a richer
620 understanding of disorders of exploration, such as attention-deficit/hyperactivity disorder (22,
621 98) including how aberrant catecholamine function might contribute to its core behavioural
622 impairments.

623 **Materials and Methods**

624 *Subjects*

625 Sixty healthy volunteers aged 18 to 35 (mean =23.22, SD =3.615) participated in a
626 double-blind, placebo-controlled, between-subjects study. The sample size was determined using
627 power calculation taking effect sizes from our prior studies that used the same drug
628 manipulations (36, 59, 75). Each subject was randomly allocated to one of three drug groups,
629 controlling for an equal gender balance across all groups (cf. Appendix 1). Candidate subjects
630 with a history of neurological or psychiatric disorders, current health issues, regular medications
631 (except contraceptives), or prior allergic reactions to drugs were excluded from the study.
632 Subjects had (self-reported) normal or corrected-to-normal vision. The groups consisted of 20
633 subjects each matched (cf. Appendix 2 Table 1) for gender and age. To evaluate peripheral drug
634 effects, heart rate, systolic and diastolic blood pressure were collected to at three different time-
635 points: ‘at arrival’, ‘pre-task’ and ‘post-task’, cf. Appendix 1 for details. At 50 minutes after
636 administering the 2nd drug, subjects were filled in the PANAS questionnaires (50) and
637 completed the WASI Matrix Reasoning subtest (49). Subjects differed in mood (PANAS
638 negative affect, cf. Appendix 1 for details) and marginally in intellectual abilities (WASI), and so
639 we control for these potential confounders in our analyses (cf. Appendix 1 for uncorrected
640 results). Subjects were reimbursed for their participation on an hourly basis and received a bonus
641 according to their performance (proportional to the sum of all the collected apples’ size). One
642 subject from the amisulpride group was excluded due to not engaging in the task and performing
643 at chance level. The study was approved by the UCL research ethics committee and all subjects
644 provided written informed consent.

645 *Pharmacological manipulation*

646 To reduce noradrenaline functioning, we administered 40mg of the non-selective β -
647 adrenoceptor antagonist propranolol 60 minutes before the task (Fig 1D). To reduce dopamine
648 functioning, we administered 400mg of the selective D2/D3 antagonist amisulpride 90 minutes
649 before the task. Because of different pharmacokinetic properties, drugs were administered at
650 different times. Each drug group received the drug on its corresponding time point and a placebo
651 at the other time point. The placebo group received placebo at both time points, in line with our
652 previous studies (36, 59, 75).

653 *Experimental paradigm*

654 To quantify different exploration strategies, we developed a multi-armed bandit task
655 implemented using Cogent (<http://www.vislab.ucl.ac.uk/cogent.php>) for MATLAB (R2018a).
656 Subjects had to choose between bandits (i.e. trees) that produced samples (i.e. apples) with
657 varying reward (i.e. size) in two different horizon conditions (Figure 1a-b). Bandits were
658 displayed during the entire duration of a trial and there was no time limit for sampling from
659 (choosing) the bandits. The sizes of apples they collected were summed and converted to an
660 amount of juice (feedback), which was displayed during 2000 ms at the end of each trial.
661 Subjects were instructed to endeavour to make the most juice and that they would receive a cash
662 bonus proportional to their performance. Overall subjects received £10 per hour and a mean
663 bonus of £1.12 (std: £0.06).

664 Similar to the horizon task (7), to induce different extents of exploration, we manipulated
665 the horizon (i.e. number of apples to be picked: 1 in the short horizon, 6 in the long horizon)
666 between trials. This horizon-manipulation, which has been extensively used to modulate
667 exploratory behaviour (21, 34, 54, 99), promotes exploration in the long horizon condition as
668 there are more opportunities to gather reward.

669 Within a single trial, each bandit had a different mean reward μ (i.e. apple size) and
670 associated uncertainty as captured by the number of initial samples (i.e. number of apples shown
671 at the beginning of the trial). Each bandit (i.e. tree) i was from one of four generative processes
672 (Figure 1c) characterised by different means μ_i and number of initial samples. The rewards
673 (apple sizes) for each bandit were sampled from a normal distribution with mean μ_i , specific to
674 the bandit, and with a fixed variance, $S^2=0.8$. The rewards were those sampled values rounded to
675 the closest integer. Each distribution was truncated to $[2, 10]$, meaning that rewards with values
676 above or below this interval were excluded, resulting in a total of 9 possible rewards (i.e. 9
677 different apple sizes; cf. Figure 1 - Figure supplement 1 for a representation). The ‘certain
678 standard bandit’ provided three initial samples and on every trial its mean μ_{cs} was sampled from
679 a normal distribution: $\mu_{cs} \sim N(5.5, 1.4)$. The ‘standard bandit’ provided one initial sample and to
680 make sure that its mean μ_s was comparable to μ_{cs} , the trials were split equally between the four
681 following: $\{\mu_s = \mu_{cs} + 1; \mu_s = \mu_{cs} - 1; \mu_s = \mu_{cs} + 2; \mu_s = \mu_{cs} - 2\}$. The ‘novel bandit’
682 provided no initial samples and its mean μ_n was comparable to both μ_{cs} and μ_s by splitting the
683 trials equally between the eight following: $\{\mu_n = \mu_{cs} + 1; \mu_n = \mu_{cs} - 1; \mu_n = \mu_{cs} + 2; \mu_n =$
684 $\mu_{cs} - 2; \mu_n = \mu_s + 1; \mu_n = \mu_s - 1; \mu_n = \mu_s + 2; \mu_n = \mu_s - 2\}$. The ‘low bandit’ provided one
685 initial sample which was smaller than all the other bandits’ means on that trial: $\mu_l =$
686 $\min(\mu_{cs}, \mu_s, \mu_n) - 1$. We ensured that the initial sample from the low-value bandit was the
687 smallest by resampling from each bandit in the trials were that was not the case. To make sure
688 that our task captures heuristic exploration strategies, we simulated behaviour (cf. Figure 1).
689 Additionally, in each trial, to avoid that some exploration strategies overshadow other ones, only
690 three of the four different groups were available to choose from. Based on the mean of the initial

691 samples, we identified the high-value option (i.e. the bandit with the highest expected reward) in
692 trials where both the certain-standard and the standard bandit were present.

693 There were 25 trials of each of the four three-bandit combination making it a total of 100
694 different trials. They were then duplicated to measure choice consistency, defined as the
695 frequency of making the same choice on identical trials (in contrast to a previous propranolol
696 study where consistency was defined in terms of a value-based exploration parameter (60)). Each
697 subject played these 200 trials both in a short and in a long horizon setting, resulting in a total of
698 400 trials. The trials were randomly assigned to one of four blocks and subjects were given a
699 short break at the end of each of them. To prevent learning, the bandits' positions (left, middle or
700 right) as well as their colour (8 sets of 3 different colours) were shuffled between trials. To
701 ensure subjects distinguished different apple sizes and understood that apples from the same tree
702 were always of similar size (generated following a normal distribution), they needed to undergo
703 training prior to the main experiment. In training, based on three displayed apples of similar size,
704 they were tasked to guess between two options, namely which apple was most likely to come
705 from the same tree and then received feedback about their choice.

706 *Statistical analyses*

707 All statistical analyses were performed using the R Statistical Software (100). For
708 computing ANOVA tests and pairwise comparisons the 'rstatix' package was used, and for
709 computing effect sizes the 'lsr' package (101) was used. To ensure consistent performance across
710 all subjects, we excluded one outlier subject (belonging to the amisulpride group) from our
711 analysis due to not engaging in the task and performing at chance level (defined as randomly
712 sampling from one out of three bandits, i.e. 33%). Each bandits' selection frequency for a
713 horizon condition was computed over all 200 trials and not only over the trials where this

714 specific bandit was present (i.e. $3/4$ of $200 = 150$ trials). In all the analysis comparing horizon
715 conditions, except when looking at score values (Figure 2c), only the 1st draw of the long horizon
716 was used. We compared behavioural measures and model parameters using (paired-samples) t-
717 tests and repeated-measures (rm-) ANOVAs with a between-subject factor of drug group
718 (propranolol group, amisulpride group, placebo group) and a within-subject factor horizon (long,
719 short). Information seeking, expected values and scores were analysed using rm-ANOVAS with
720 a within-subject factor horizon. Measures that were horizon-independent (e.g. prior mean), were
721 analysed using one-way ANOVAs with a between-subject factor drug group. As drug groups
722 differed in negative affect (cf. Appendix 2 Table 1), which, through its relationship to anxiety
723 (*I02*) is thought to affect cognition (*I03*) and potentially exploration (*I04*). We corrected for
724 negative affect (PANAS) and IQ (WASI) in each analysis by adding those two measures as
725 covariates in each ANOVA mentioned above (cf. Appendix 1 for analysis without covariates and
726 analysis with physiological effect as an additional covariates). We report effect sizes using
727 partial eta squared (η^2) for ANOVAs and Cohen's d (d) for t-tests (*I05*).

728 *Computational modelling*

729 We adapted a set of Bayesian generative models from previous studies (*I*), where each
730 model assumed that different characteristics account for subjects' behaviour. The binary
731 indicators (c_{tr}, c_n) indicate which components (value-free random and novelty exploration
732 respectively) were included in the different models. The value of each bandit is represented as a
733 distribution $N(Q, S)$ with $S = 0.8$, the sampling variance fixed to its generative value. Subjects
734 have prior beliefs about bandits' values which we assume to be Gaussian with mean Q_0 and
735 uncertainty σ_0 . The subjects' initial estimate of a bandit's mean (Q_0 ; prior mean) and its
736 uncertainty about it (σ_0 ; prior variance) are free parameters.

737 These beliefs are updated according to Bayes rule (detailed below) for each initial sample (note
 738 that there are no updates for the novel bandit).

739 *Mean and variance update rules*

740 At each time point t , in which a sample m , of one of the bandits is presented, the
 741 expected mean Q and precision $\tau = \frac{1}{\sigma^2}$ of the corresponding bandit i are updated as follows:

$$Q_{i,t+1} = \frac{\tau_{i,t} * Q_{i,t} + \tau_{samp} * m}{\tau_{i,t} + \tau_{samp}}$$

$$\tau_{t+1}^i = \tau_{samp} + \tau_t^i$$

742 where $\tau_{samp} = \frac{1}{S^2}$ is the sampling precision, with the sampling variance $S = 0.8$ fixed. Those
 743 update rules are equivalent to using a Kalman filter (106) in stationary bandits.

744 We examined three base models: the UCB model, the Thompson model and the hybrid
 745 model. The UCB model encompasses the UCB algorithm (captures directed exploration) and a
 746 softmax choice function (captures a value-based random exploration). The Thompson model
 747 reflects Thompson sampling (captures an uncertainty-driven value-based random exploration).
 748 The hybrid model captures the contribution of the UCB model and the Thompson model,
 749 essentially a mixture of the above. We computed three extensions of each model by either adding
 750 value-free random exploration $(c_{vf}, c_n) = (1,0)$, novelty exploration $(c_{vf}, c_n) = (0,1)$ or both
 751 heuristics $(c_{vf}, c_n) = (1,1)$, leading to a total of 12 models (see the labels on the x-axis in Figure
 752 4a; $(c_{vf}, c_n) = (0,0)$ is the model with no extension). For additional models cf. Appendix 1. A
 753 coefficient $c_{vf}=1$ indicates that a ϵ -greedy component was added to the decision rule, ensuring
 754 that once in a while (every ϵ % of the time), another option than the predicted one is selected. A
 755 coefficient $c_n=1$ indicates that the novelty bonus η is added to the computation of the value of
 756 novel bandits and the Kronecker delta δ in front of this bonus ensures that it is only applied to

757 the novel bandit. The models and their free parameters (summarised in Appendix 2 Table 5) are
 758 described in detail below.

759 *Choice rules*

760 *UCB model.* In this model, an information bonus γ is added to the expected reward of each
 761 option, scaling with the option's uncertainty (UCB). The value of each bandit i at timepoint t is:

$$V_{i,t} = Q_{i,t} + \gamma\sigma_{i,t} + c_n\eta\delta_{[i=new]}$$

762 The probability of choosing bandit i was given by passing this into the softmax decision
 763 function:

$$P(c_t = i) = \frac{e^{\beta V_{i,t}}}{\sum_x e^{\beta V_{x,t}}} * (1 - c_{vf}\epsilon) + c_{vf}\frac{\epsilon}{3}$$

764 where β is the inverse temperature of the softmax (lower values producing more value-
 765 based random exploration), and the coefficient c_{vf} adds the value-free random exploration
 766 component.

767 *Thompson model.* In this model, based on Thompson sampling, the overall uncertainty can be
 768 seen as a more refined version of a decision temperature (I). The value of each bandit i is as
 769 before:

$$V_{i,t} = Q_{i,t} + c_n\eta\delta_{[i=new]}$$

770 A sample $x_{i,t} \sim N(V_{i,t}, \sigma_{i,t}^2)$ is taken from each bandit. The probability of choosing a
 771 bandit i depends on the probability that all pairwise differences between the sample from bandit i
 772 and the other bandits $j \neq i$ were greater or equal to 0 (see the probability of maximum utility
 773 choice rule (107)). In our task, because three bandits were present, two pairwise differences
 774 scores (contained in the two-dimensional vector u) were computed for each bandit. The
 775 probability of choosing bandit i is:

$$\begin{aligned}
P(c_t = i) &= P(\forall j: x_{i,t} > x_{j,t}) * (1 - c_{vf}\epsilon) + c_{vf}\frac{\epsilon}{3} \\
&= \int_0^\infty \int_0^\infty \phi(u; M_{i,t}, C_{i,t}) du * (1 - c_{vf}\epsilon) + c_{vf}\frac{\epsilon}{3}
\end{aligned}$$

776 where ϕ is the multivariate Normal density function with mean vector

$$M_{i,t} = A_i \begin{pmatrix} V_{1,t} \\ V_{2,t} \\ V_{3,t} \end{pmatrix}$$

777 and covariance matrix

$$C_{i,t} = A_i \begin{pmatrix} \sigma_{1,t} & 0 & 0 \\ 0 & \sigma_{2,t} & 0 \\ 0 & 0 & \sigma_{3,t} \end{pmatrix} A_i^T$$

778 Where the matrix A_i computes the pairwise differences between bandit i and the other bandits.

779 For example, for bandit $i = 1$:

$$A_1 = \begin{pmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix}$$

780 *Hybrid model.* This model allows a combination of the UCB model and the Thompson model.

781 The probability of choosing bandit i is:

$$P(c_t = i) = \left(wP_{UCB}(c_t = i) + (1 - w)P_{Thompson}(c_t = i) \right) * (1 - c_{vf}\epsilon) + c_{vf}\frac{\epsilon}{3}$$

782 where w specifies the contribution of each of the two models. P_{UCB} and $P_{Thompson}$ are

783 calculated for $c_{vf}=0$. If $w=1$, only the UCB model is used while if $w=0$ only the Thompson

784 model is used. In between values indicate a mixture of the two models.

785 All the parameters besides Q_0 and w were free to vary as a function of the horizon (cf.
786 Appendix 2 Table 5) as they capture different exploration forms: directed exploration
787 (information bonus γ ; UCB model), novelty exploration (novelty bonus η), value-based random
788 exploration (inverse temperature β ; UCB model), uncertainty-directed exploration (prior
789 variance σ_0 ; Thompson model) and value-free random exploration (ϵ -greedy parameter). The
790 prior mean Q_0 was fitted to both horizons together as we do not expect the belief of how good a
791 bandit is to depend on the horizon. The same was done for w as assume the arbitration between
792 the UCB model and the Thompson model does not depend on horizon.

793 *Parameter estimation.*

794 To fit the parameter values, we used the maximum a posteriori probability (MAP) estimate. The
795 optimisation function used was `fmincon` in MATLAB. The parameters could vary within the
796 following bounds: $\sigma_0 = [0.01, 6]$, $Q_0 = [1, 10]$, $\epsilon = [0, 0.5]$, $\eta = [0, 5]$. The prior distribution
797 used for the prior mean parameter Q_0 was the normal distribution: $Q_0 \sim N(5, 2)$ that
798 approximates the generative distributions. For the ϵ -greedy parameter, the novelty bonus η and
799 the prior variance parameter σ_0 , a uniform distribution (of range equal to the specific parameters'
800 bounds) was used, which is equivalent to performing MLE. A summary of the parameter values
801 per group and per horizon can be found in Appendix 2 Table 6.

802 *Model comparison.*

803 We performed a K-fold cross-validation with $K = 10$. We partitioned the data of each subject
804 ($N_{trials} = 400$; 200 in each horizon) into K folds (i.e. subsamples). For model fitting in our
805 model selection, we used maximum likelihood estimation (MLE), where we maximised the
806 likelihood for each subject individually (`fmincon` was ran with 8 randomly chosen starting point

807 to overcome potential local minima). We fitted the model using K-1 folds and validated the
808 model on the remaining fold. We repeated this process K times, so that each of the K fold is used
809 as a validation set once, and averaged the likelihood over held out trials. We did this for each
810 model and each subject and averaged across subjects. The model with the highest likelihood of
811 held-out data (the winning model) was the Thompson sampling with $(c_{tr}, c_n) = \{1,1\}$. It was
812 also the model which accounted best for the largest number of subjects (Figure 4 – Figure
813 supplement 1).

814 *Parameter recovery.*

815 To make sure that the parameters are interpretable, we performed a parameter recovery analysis.
816 For each parameter, we took 4 values, equally spread, within a reasonable parameter range
817 ($\sigma_0 = [0.5, 2.5]$, $Q_0 = [1, 6]$, $\epsilon = [0, 0.5]$, $\eta = [0, 5]$). All parameters but Q_0 were free to vary as
818 a function of the horizon. We simulated behaviour with one artificial agent for each 4^7
819 combinations using a new trial for each. The model was fitted using MAP estimation (cf.
820 Parameter estimation) and analysed how well the generative parameters (generating parameters
821 in Figure 5) correlated with the recovered ones (fitted parameters in Figure 5) using Pearson
822 correlation (summarised in Figure 5c). In addition to the correlation we examined the spread
823 (Figure 4 – Figure supplement 3) of the recovered parameters. Overall the parameters were well
824 recoverable.

825 *Model validation*

826 To validate our model, we used each subjects' fitted parameters to simulate behaviour on our
827 task (4000 trials per agent). The stimulated data (Figure 5 – Figure supplement 1), although not
828 perfect, resembles the real data reasonably well. Additionally, to validate the behavioural
829 indicators of the two different exploration heuristics we stimulated the behaviour of 200 agents

830 using the winning model on one horizon condition (i.e. trials = 200). For the indicators of value-
831 free random exploration, we stimulated behaviour with low ($\epsilon = 0$) and high ($\epsilon = 0.2$) values of
832 the ϵ -greedy parameter. The other parameters were set to the mean parameter fits ($\sigma_0 =$
833 $1.312, \eta = 2.625, Q_0 = 3.2$). This confirms that higher amounts of value-free random
834 exploration are captured by the proportion of low-value bandit selection (Figure 1f) and the
835 choice consistency (Figure 1e). Similarly, for the indicator of novelty exploration, we simulated
836 behaviour with low ($\eta = 0$) and high ($\eta = 2$) values of the novelty bonus η to validate the use of
837 the proportion of the novel-bandit selection (Figure 1g). Again, the remaining parameters were
838 set to the mean parameter fits ($\sigma_0 = 1.312, \epsilon = 0.1, Q_0 = 3.2$). Parameter values for high and
839 low exploration were selected empirically from pilot and task data. Additionally, we simulated
840 the effects of other exploration strategies in short and long horizon conditions (Figure 1 – Figure
841 supplement 3-5). To simulate a long (versus short) horizon condition we increased the overall
842 exploration by increasing other exploration strategies. Details about parameter values can be
843 found in Appendix 2 Table 7.

844 **Conflict of interest**

845 The authors declare no competing financial interests.

846 **Acknowledgements**

847 M.D. is a predoctoral fellow of the International Max Planck Research School on Computational
848 Methods in Psychiatry and Ageing Research. The participating institutions are the Max Planck
849 Institute for Human Development and the University College London (UCL). T.U.H. is
850 supported by a Wellcome Sir Henry Dale Fellowship (211155/Z/18/Z), a grant from the Jacobs
851 Foundation (2017-1261-04), the Medical Research Foundation, a 2018 NARSAD Young
852 Investigator Grant (27023) from the Brain and Behavior Research Foundation, and an ERC
853 Starting Grant. R.J.D. holds a Wellcome Trust Investigator Award (098362/Z/12/Z). The Max
854 Planck UCL Centre is a joint initiative supported by UCL and the Max Planck Society. The
855 Wellcome Centre for Human Neuroimaging is supported by core funding from the Wellcome
856 Trust (203147/Z/16/Z).

857 **References**

- 858 1. S. J. Gershman, Deconstructing the human algorithms for exploration. *Cognition*. **173**, 34–42 (2018).
859 2. E. Schulz, S. J. Gershman, The algorithmic architecture of exploration in the human brain. *Curr. Opin.*
860 *Neurobiol.* **55**, 7–14 (2019).
861 3. P. Auer, Using confidence bounds for exploitation-exploration trade-offs. *J. Mach. Learn. Res.* **3**, 397–422
862 (2003).
863 4. A. Carpentier, A. Lazaric, M. Ghavamzadeh, R. Munos, P. Auer, Upper-confidence-bound algorithms for
864 active learning in multi-armed bandits. *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell.*
865 *Lect. Notes Bioinformatics)*. **6925 LNAI**, 189–203 (2011).
866 5. P. Schwartenbeck, J. Passerker, T. U. Hauser, T. H. FitzGerald, M. Kronbichler, K. J. Friston,
867 Computational mechanisms of curiosity and goal-directed exploration. *Elife* (2019),
868 doi:10.7554/eLife.41703.
869 6. N. D. Daw, J. P. O’Doherty, P. Dayan, B. Seymour, R. J. Dolan, Cortical substrates for exploratory
870 decisions in humans. *Nature*. **441**, 876–879 (2006).
871 7. R. C. Wilson, A. Geana, J. M. White, E. A. Ludvig, J. D. Cohen, Humans use directed and random
872 exploration to solve the explore–exploit dilemma. *J. Exp. Psychol. Gen.* **143**, 2074–2081 (2014).
873 8. W. R. Thompson, On the Likelihood that One Unknown Probability Exceeds Another in View of the
874 Evidence of Two Samples. **25**, 285–294 (1933).
875 9. J. D. Cohen, S. M. McClure, A. J. Yu, Should I stay or should I go? How the human brain manages the
876 trade-off between exploitation and exploration. *Philos. Trans. R. Soc. B Biol. Sci.* **362**, 933–942 (2007).
877 10. I. C. Dezza, A. Cleeremans, W. Alexander, Should we control? The interplay between cognitive control and
878 information integration in the resolution of the exploration-exploitation dilemma. *J. Exp. Psychol. Gen.*
879 (2019), doi:10.1037/xge0000546.
880 11. D. Papadopetraki, M. Froböse, A. Westbrook, B. Zandbelt, R. Cools, Quantifying the cost of cognitive
881 stability and flexibility (2019), doi:10.1101/743120.
882 12. Z. Alexandre, S. Oleg, P. Giovanni, An information-theoretic perspective on the costs of cognition.
883 *Neuropsychologia*. **123**, 5–18 (2019).
884 13. B. Wahn, P. König, Is attentional resource allocation across sensory modalities task-dependent? *Adv. Cogn.*
885 *Psychol.* **13**, 83–96 (2017).
886 14. R. Marois, J. Ivanoff, Capacity limits of information processing in the brain. *Trends Cogn. Sci.* **9**, 296–305
887 (2005).
888 15. M. Botvinick, T. Braver, Motivation and cognitive control: From behavior to neural mechanism. *Annu. Rev.*
889 *Psychol.* (2015), doi:10.1146/annurev-psych-010814-015044.
890 16. M. I. Froböse, A. Westbrook, M. Bloemendaal, E. Aarts, R. Cools, Catecholaminergic modulation of the
891 cost of cognitive control in healthy older adults. *PLoS One*. **15**, 1–26 (2020).

- 892 17. W. Kool, J. T. McGuire, Z. B. Rosen, M. M. Botvinick, Decision Making and the Avoidance of Cognitive
893 Demand. *J. Exp. Psychol. Gen.* **139**, 665–682 (2010).
- 894 18. R. Cools, The cost of dopamine for dynamic cognitive control. *Curr. Opin. Behav. Sci.* (2015), ,
895 doi:10.1016/j.cobeha.2015.05.007.
- 896 19. M. I. Froböse, R. Cools, Chemical neuromodulation of cognitive control avoidance. *Curr. Opin. Behav. Sci.*
897 **22**, 121–127 (2018).
- 898 20. R. S. Sutton, A. G. Barto, Introduction to Reinforcement Learning. *MIT Press Cambridge* (1998),
899 doi:10.1.1.32.7692.
- 900 21. W. K. Zajkowski, M. Kossut, R. C. Wilson, A causal role for right frontopolar cortex in directed, but not
901 random, exploration. *Elife.* **6**, 1–18 (2017).
- 902 22. T. U. Hauser, V. G. Fiore, M. Moutoussis, R. J. Dolan, Computational Psychiatry of ADHD: Neural Gain
903 Impairments across Marrian Levels of Analysis. *Trends Neurosci.* **39**, 63–73 (2016).
- 904 23. S. J. Sara, A. Vankov, A. Hervé, Locus coeruleus-evoked responses in behaving rats: A clue to the role of
905 noradrenaline in memory. *Brain Res. Bull.* **35**, 457–465 (1994).
- 906 24. C. Varazzani, A. San-Galli, S. Gilardeau, S. Bouret, Noradrenaline and dopamine neurons in the
907 reward/effort trade-off: A direct electrophysiological comparison in behaving monkeys. *J. Neurosci.* **35**,
908 7866–7877 (2015).
- 909 25. M. Silvetti, E. Vassena, E. Abrahamse, T. Verguts, *Dorsal anterior cingulate-brainstem ensemble as a*
910 *reinforcement meta-learner* (2018), vol. 14.
- 911 26. M. Silvetti, R. Seurinck, M. E. van Bochove, T. Verguts, The influence of the noradrenergic system on
912 optimal control of neural plasticity. *Front. Behav. Neurosci.* **7**, 1–6 (2013).
- 913 27. A. J. Yu, P. Dayan, Uncertainty, neuromodulation, and attention. *Neuron.* **46**, 681–692 (2005).
- 914 28. M. R. Nassar, K. M. Rumsey, R. C. Wilson, K. Parikh, B. Heasley, J. I. Gold, Rational regulation of learning
915 dynamics by pupil-linked arousal systems. *Nat. Neurosci.* (2012), doi:10.1038/nn.3130.
- 916 29. J. David Johnson, Noradrenergic control of cognition: global attenuation and an interrupt function. *Med.*
917 *Hypotheses* (2003), doi:10.1016/s0306-9877(03)00021-5.
- 918 30. S. Bouret, S. J. Sara, Network reset: A simplified overarching theory of locus coeruleus noradrenaline
919 function. *Trends Neurosci.* (2005), doi:10.1016/j.tins.2005.09.002.
- 920 31. P. Dayan, A. J. Yu, Phasic norepinephrine: A neural interrupt signal for unexpected events. *Netw. Comput.*
921 *Neural Syst.* **17**, 335–350 (2006).
- 922 32. D. G. R. Tervo, M. Proskurin, M. Manakov, M. Kabra, A. Vollmer, K. Branson, A. Y. Karpova, Behavioral
923 variability through stochastic choice and its gating by anterior cingulate cortex. *Cell.* **159**, 21–32 (2014).
- 924 33. C. I. Jahn, S. Gilardeau, C. Varazzani, B. Blain, J. Sallet, M. E. Walton, S. Bouret, Dual contributions of
925 noradrenaline to behavioural flexibility and motivation, 2687–2702 (2018).
- 926 34. C. M. Warren, R. C. Wilson, N. J. Van Der Wee, E. J. Giltay, M. S. Van Noorden, J. D. Cohen, S.
927 Nieuwenhuis, The effect of atomoxetine on random and directed exploration in humans. *PLoS One* (2017),
928 doi:10.1371/journal.pone.0176034.
- 929 35. P. F. Fraundorfer, R. H. Fertel, D. D. Miller, D. R. Feller, Biochemical and pharmacological characterization
930 of high-affinity trimetoquinol analogs on guinea pig and human beta adrenergic receptor subtypes: Evidence
931 for partial agonism. *J. Pharmacol. Exp. Ther.* **270**, 665–674 (1994).
- 932 36. T. U. Hauser, E. Eldar, N. Purg, M. Moutoussis, R. J. Dolan, Distinct roles of dopamine and noradrenaline
933 in incidental memory. *J. Neurosci.* (2019), doi:10.1523/jneurosci.0401-19.2019.
- 934 37. Ki Database, (available at <https://pdsp.unc.edu/databases/pdsp.php>).
- 935 38. N. Bunzeck, C. F. Doeller, R. J. Dolan, E. Duzel, Contextual interaction between novelty and reward
936 processing within the mesolimbic system. *Hum. Brain Mapp.* (2012), doi:10.1002/hbm.21288.
- 937 39. B. C. Wittmann, N. D. Daw, B. Seymour, R. J. Dolan, Striatal Activity Underlies Novelty-Based Choice in
938 Humans. *Neuron* (2008), doi:10.1016/j.neuron.2008.04.027.
- 939 40. S. J. Gershman, Y. Niv, Novelty and Inductive Generalization in Human Reinforcement Learning. *Top.*
940 *Cogn. Sci.* **7**, 391–415 (2015).
- 941 41. H. Stojic, E. Shulz, P. P. Analytis, M. Speekenbrink, It’s new, but is it good? How generalization and
942 uncertainty guide the exploration of novel options. *PsyArXiv* (2018).
- 943 42. R. M. Krebs, B. H. Schott, H. Schütze, E. Düzel, The novelty exploration bonus and its attentional
944 modulation. *Neuropsychologia* (2009), doi:10.1016/j.neuropsychologia.2009.01.015.
- 945 43. M. J. Frank, B. B. Doll, J. Oas-Terpstra, F. Moreno, Prefrontal and striatal dopaminergic genes predict
946 individual differences in exploration and exploitation. *Nat. Neurosci.* **12**, 1062–1068 (2009).
- 947 44. E. S. Bromberg-Martin, M. Matsumoto, O. Hikosaka, Dopamine in Motivational Control: Rewarding,

- 948 Aversive, and Alerting. *Neuron* (2010), , doi:10.1016/j.neuron.2010.11.022.
- 949 45. V. D. Costa, V. L. Tran, J. Turchi, B. B. Averbeck, Dopamine modulates novelty seeking behavior during
950 decision making. *Behav. Neurosci.* (2014), doi:10.1037/a0037128.
- 951 46. E. Düzel, W. D. Penny, N. Burgess, Brain oscillations and memory. *Curr. Opin. Neurobiol.* (2010), ,
952 doi:10.1016/j.conb.2010.01.004.
- 953 47. K. Iigaya, T. U. Hauser, Z. Kurth-Nelson, J. P. O’Doherty, P. Dayan, R. J. Dolan, The value of what’s to
954 come: neural mechanisms coupling prediction error and reward anticipation. *bioRxiv* (2019),
955 doi:10.1101/588699.
- 956 48. A. S. Kayser, J. M. Mitchell, D. Weinstein, M. J. Frank, Dopamine, locus of control, and the exploration-
957 exploitation tradeoff. *Neuropsychopharmacology* (2015), doi:10.1038/npp.2014.193.
- 958 49. D. Wechsler, WASI -II: Wechsler abbreviated scale of intelligence - second edition. *J. Psychoeduc. Assess.*
959 (2013), doi:10.1177/0734282912467756.
- 960 50. D. Watson, L. A. Clark, A. Tellegen, Development and Validation of Brief Measures of Positive and
961 Negative Affect: The PANAS Scales. *J. Pers. Soc. Psychol.* (1988), doi:10.1037/0022-3514.54.6.1063.
- 962 51. S. Agrawal, N. Goyal, Analysis of thompson sampling for the multi-armed bandit problem. *J. Mach. Learn.*
963 *Res.* **23**, 1–26 (2012).
- 964 52. M. D’Acemont, P. Bossaerts, Neurobiological studies of risk assessment: A comparison of expected utility
965 and mean-variance approaches. *Cogn. Affect. Behav. Neurosci.* **8**, 363–374 (2008).
- 966 53. N. C. Foley, D. C. Jangraw, C. Peck, J. Gottlieb, Novelty enhances visual salience independently of reward
967 in the parietal lobe. *J. Neurosci.* (2014), doi:10.1523/JNEUROSCI.4171-13.2014.
- 968 54. C. M. Wu, E. Schulz, M. Speekenbrink, J. D. Nelson, B. Meder, Generalization guides human exploration in
969 vast decision spaces. *Nat. Hum. Behav.* **2**, 915–924 (2018).
- 970 55. V. Skvortsova, S. Palminteri, M. Pessiglione, Learning to minimize efforts versus maximizing rewards:
971 Computational principles and neural correlates. *J. Neurosci.* **34**, 15621–15630 (2014).
- 972 56. T. U. Hauser, E. Eldar, R. J. Dolan, Separate mesocortical and mesolimbic pathways encode effort and
973 reward learning signals. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E7395–E7404 (2017).
- 974 57. M. E. Walton, S. Bouret, What Is the Relationship between Dopamine and Effort? *Trends Neurosci.* **42**, 79–
975 91 (2019).
- 976 58. J. D. Salamone, S. E. Yohn, L. López-Cruz, N. San Miguel, M. Correa, Activational and effort-related
977 aspects of motivation: Neural mechanisms and implications for psychopathology. *Brain.* **139**, 1325–1347
978 (2016).
- 979 59. T. U. Hauser, M. Moutoussis, N. Purg, P. Dayan, R. J. Dolan, Beta-Blocker Propranolol Modulates Decision
980 Urgency During Sequential Information Gathering. *J. Neurosci.* (2018), doi:10.1523/jneurosci.0192-
981 18.2018.
- 982 60. P. Sokol-Hessner, S. F. Lackovic, R. H. Tobe, C. F. Camerer, B. L. Leventhal, E. A. Phelps, Determinants
983 of Propranolol’s Selective Effect on Loss Aversion. *Psychol. Sci.* **26**, 1123–1130 (2015).
- 984 61. D. Campbell-Meiklejohn, J. Wakeley, V. Herbert, J. Cook, P. Scollo, M. K. Ray, S. Selvaraj, R. E.
985 Passingham, P. Cowen, R. D. Rogers, Serotonin and dopamine play complementary roles in gambling to
986 recover losses. *Neuropsychopharmacology.* **36**, 402–410 (2011).
- 987 62. R. D. Rogers, M. Lancaster, J. Wakeley, Z. Bhagwagar, Effects of beta-adrenoceptor blockade on
988 components of human decision-making. *Psychopharmacology (Berl).* **172**, 157–164 (2004).
- 989 63. J. Rajkowski, P. Kubiak, G. Aston-Jones, Locus coeruleus activity in monkey: Phasic and tonic changes are
990 associated with altered vigilance. *Brain Res. Bull.* **35**, 607–616 (1994).
- 991 64. G. Aston-Jones, J. D. Cohen, AN INTEGRATIVE THEORY OF LOCUS COERULEUS-
992 NOREPINEPHRINE FUNCTION: Adaptive Gain and Optimal Performance. *Annu. Rev. Neurosci.* **28**, 403–
993 450 (2005).
- 994 65. D. Servan-Schreiber, H. Printz, J. D. Cohen, A network model of catecholamine effects: Gain, signal-to-
995 noise ratio, and behavior. *Science (80-)*. (1990), doi:10.1126/science.2392679.
- 996 66. S. Joshi, Y. Li, R. M. Kalwani, J. I. Gold, Relationships between Pupil Diameter and Neuronal Activity in
997 the Locus Coeruleus, Colliculi, and Cingulate Cortex. *Neuron* (2016), doi:10.1016/j.neuron.2015.11.028.
- 998 67. S. Joshi, J. I. Gold, Pupil Size as a Window on Neural Substrates of Cognition. *Trends Cogn. Sci.* **24**, 466–
999 480 (2020).
- 1000 68. V. Koudas, A. Nikolaou, E. Hourdaki, S. G. Giakoumaki, P. Roussos, P. Bitsios, Comparison of ketanserin,
1001 buspirone and propranolol on arousal, pupil size and autonomic function in healthy volunteers.
1002 *Psychopharmacology (Berl).* (2009), doi:10.1007/s00213-009-1508-5.
- 1003 69. M. Jepma, S. Nieuwenhuis, Pupil diameter predicts changes in the exploration-exploitation trade-off:

- 1004 Evidence for the adaptive gain theory. *J. Cogn. Neurosci.* (2011), doi:10.1162/jocn.2010.21548.
- 1005 70. Jepma, The role of the noradrenergic system in the exploration-exploitation trade-off: a pharmacological
1006 study. *Front. Hum. Neurosci.* (2010), doi:10.3389/fnhum.2010.00170.
- 1007 71. M. Usher, J. D. Cohen, D. Servan-Schreiber, J. Rajkowski, G. Aston-Jones, The role of locus coeruleus in
1008 the regulation of cognitive performance. *Science (80-.)*. **283**, 549–554 (1999).
- 1009 72. G. A. Kane, E. M. Vazey, R. C. Wilson, A. Shenhav, N. D. Daw, G. Aston-Jones, J. D. Cohen, Increased
1010 locus coeruleus tonic activity causes disengagement from a patch-foraging task. *Cogn. Affect. Behav.*
1011 *Neurosci.* **17**, 1073–1083 (2017).
- 1012 73. Z. L. Rossetti, S. Carboni, Noradrenaline and dopamine elevations in the rat prefrontal cortex in spatial
1013 working memory. *J. Neurosci.* **25**, 2322–2329 (2005).
- 1014 74. M. E. Gibbs, D. S. Hutchinson, R. J. Summers, Noradrenaline release in the locus coeruleus modulates
1015 memory formation and consolidation; roles for α - and β -adrenergic receptors. *Neuroscience.* **170**, 1209–
1016 1222 (2010).
- 1017 75. T. U. Hauser, M. Allen, N. Purg, M. Moutoussis, G. Rees, R. J. Dolan, Noradrenaline blockade specifically
1018 enhances metacognitive performance. *Elife* (2017), doi:10.7554/eLife.24901.
- 1019 76. I. Trofimova, T. W. Robbins, Temperament and arousal systems: A new synthesis of differential psychology
1020 and functional neurochemistry. *Neurosci. Biobehav. Rev.* **64**, 382–402 (2016).
- 1021 77. B. D. Waterhouse, H. C. Moises, H. H. Yeh, D. J. Woodward, Norepinephrine enhancement of inhibitory
1022 synaptic mechanisms in cerebellum and cerebral cortex: Mediation by beta adrenergic receptors. *J.*
1023 *Pharmacol. Exp. Ther.* **221**, 495–506 (1982).
- 1024 78. B. D. Waterhouse, H. C. Moises, H. H. Yeh, H. M. Geller, D. J. Woodward, Comparison of norepinephrine-
1025 and benzodiazepine-induced augmentation of Purkinje cell response to γ -aminobutyric acid (GABA). *J.*
1026 *Pharmacol. Exp. Ther.* **228**, 257–267 (1984).
- 1027 79. P. S. Goldman-Rakic, M. S. Lidow, D. W. Gallager, Overlap of dopaminergic, adrenergic, and
1028 serotonergic receptors and complementarity of their subtypes in primate prefrontal cortex. *J. Neurosci.* **10**,
1029 2125–2138 (1990).
- 1030 80. J. S. Isaacson, M. Scanziani, How Inhibition Shapes Cortical Activity Excitation and inhibition walk hand in
1031 hand. *Neuron.* **72**, 231–243 (2011).
- 1032 81. H. Salgado, M. Treviño, M. Atzori, Layer- and area-specific actions of norepinephrine on cortical synaptic
1033 transmission. *Brain Res.* **1641**, 163–176 (2016).
- 1034 82. B. De Martino, B. A. Strange, R. J. Dolan, Noradrenergic neuromodulation of human attention for emotional
1035 and neutral stimuli. *Psychopharmacology (Berl).* **197**, 127–136 (2008).
- 1036 83. W. Schultz, P. Dayan, P. R. Montague, A neural substrate of prediction and reward. *Science (80-.)*. **275**,
1037 1593–1599 (1997).
- 1038 84. B. C. Wittmann, N. D. Daw, B. Seymour, R. J. Dolan, Striatal Activity Underlies Novelty-Based Choice in
1039 Humans. *Neuron.* **58**, 967–973 (2008).
- 1040 85. V. D. Costa, V. L. Tran, J. Turchi, B. B. Averbeck, Dopamine modulates novelty seeking behavior during
1041 decision making. *Behav. Neurosci.* **128**, 556–566 (2014).
- 1042 86. A. S. Kayser, J. M. Mitchell, D. Weinstein, M. J. Frank, Dopamine, locus of control, and the exploration-
1043 exploitation tradeoff. *Neuropsychopharmacology.* **40**, 454–462 (2015).
- 1044 87. K. Chakroun, D. Mathar, A. Wiehler, F. Ganzer, J. Peters, Dopaminergic modulation of the
1045 exploration/exploitation trade-off in human decision-making. *bioRxiv*, 706176 (2019).
- 1046 88. F. Cinotti, V. Fresno, N. Aklil, E. Coutureau, B. Girard, A. R. Marchand, M. Khamassi, Dopamine blockade
1047 impairs the exploration-exploitation trade-off in rats. *Sci. Rep.* **9**, 1–14 (2019).
- 1048 89. M. J. Frank, B. B. Doll, J. Oas-Terpstra, F. Moreno, Prefrontal and striatal dopaminergic genes predict
1049 individual differences in exploration and exploitation. *Nat. Neurosci.* (2009), doi:10.1038/nn.2342.
- 1050 90. L. K. Krugel, G. Biele, P. N. C. Mohr, S. C. Li, H. R. Heekeren, Genetic variation in dopaminergic
1051 neuromodulation influences the ability to rapidly and flexibly adapt decisions. *Proc. Natl. Acad. Sci. U. S. A.*
1052 (2009), doi:10.1073/pnas.0905191106.
- 1053 91. J. Fang, P. H. Yu, Effect of haloperidol and its metabolites on dopamine and noradrenaline uptake in rat
1054 brain slices. *Psychopharmacology (Berl).* (1995), doi:10.1007/BF02246078.
- 1055 92. M. Toru, M. Takashima, Haloperidol in large doses reduces the cataleptic response and increases
1056 noradrenaline metabolism in the brain of the rat. *Neuropharmacology* (1985), doi:10.1016/0028-
1057 3908(85)90079-6.
- 1058 93. T. Kahnt, S. C. Weber, H. Haker, T. W. Robbins, P. N. Tobler, Dopamine D2-Receptor Blockade Enhances
1059 Decoding of Prefrontal Signals in Humans. *J. Neurosci.* (2015), doi:10.1523/jneurosci.4182-14.2015.

- 1060 94. T. Kahnt, P. N. Tobler, Dopamine Modulates the Functional Organization of the Orbitofrontal Cortex. *J.*
1061 *Neurosci.* (2017), doi:10.1523/jneurosci.2827-16.2016.
- 1062 95. M. D. Humphries, M. Khamassi, K. Gurney, Dopaminergic control of the exploration-exploitation trade-off
1063 via the basal ganglia. *Front. Neurosci.* (2012), doi:10.3389/fnins.2012.00009.
- 1064 96. A. Soutschek, C. J. Burke, A. Raja Beharelle, R. Schreiber, S. C. Weber, I. I. Karipidis, J. Ten Velden, B.
1065 Weber, H. Haker, T. Kalenscher, P. N. Tobler, The dopaminergic reward system underpins gender
1066 differences in social preferences. *Nat. Hum. Behav.* (2017), doi:10.1038/s41562-017-0226-y.
- 1067 97. A. Soutschek, G. Gvozdanic, R. Kozak, S. Duvvuri, N. de Martinis, B. Harel, D. L. Gray, E. Fehr, A.
1068 Jetter, P. N. Tobler, Dopaminergic D1 Receptor Stimulation Affects Effort and Risk Preferences. *Biol.*
1069 *Psychiatry* (2019), doi:10.1016/j.biopsych.2019.09.002.
- 1070 98. T. U. Hauser, R. Iannaccone, J. Ball, C. Mathys, D. Brandeis, S. Walitza, S. Brem, Role of the medial
1071 prefrontal cortex in impaired decision making in juvenile attention-deficit/hyperactivity disorder. *JAMA*
1072 *Psychiatry.* **71**, 1165–1173 (2014).
- 1073 99. D. Guo, A. J. Yu, in *Advances in Neural Information Processing Systems* (2018).
- 1074 100. R. R Development Core Team, *R: A Language and Environment for Statistical Computing* (2011).
- 1075 101. D. Navarro, *Learning statistics with R: A tutorial for psychology students and other beginners. (Version 0.5)*
1076 (2015; <http://ua.edu.au/ccs/teaching/lr>).
- 1077 102. D. Watson, L. A. Clark, G. Carey, Positive and Negative Affectivity and Their Relation to Anxiety and
1078 Depressive Disorders. *J. Abnorm. Psychol.* **97**, 346–353 (1988).
- 1079 103. S. J. Bishop, C. Gagne, Anxiety, Depression, and Decision Making: A Computational Perspective. *Annu.*
1080 *Rev. Neurosci.* **41**, 371–388 (2018).
- 1081 104. L. de Visser, L. J. van der Knaap, A. J. A. E. van de Loo, C. M. M. van der Weerd, F. Ohl, R. van den Bos,
1082 Trait anxiety affects decision-making differently in healthy men and women: Towards gender-specific
1083 endophenotypes of anxiety. *Neuropsychologia.* **48**, 1598–1606 (2010).
- 1084 105. J. T. E. Richardson, Eta squared and partial eta squared as measures of effect size in educational research.
1085 *Educ. Res. Rev.* (2011), , doi:10.1016/j.edurev.2010.12.001.
- 1086 106. C. M. Bishop, in *Information Science and Statistics* (2006).
- 1087 107. M. Speekenbrink, E. Konstantinidis, Uncertainty and exploration in a restless bandit problem. *Top. Cogn.*
1088 *Sci.* (2015), doi:10.1111/tops.12145.
- 1089
1090

1091 **Appendix 1**

1092

1093 **Drug effect on response times**

1094 There were no differences in response times (RT) between drug groups in the one-way ANOVA. Neither in the
1095 mean RT (ANOVA: $F(2, 54)=1.625, p=.206, \eta^2=.057$) nor in its variability (standard deviation; $F(2, 54)=1.85,$
1096 $p=.16, \eta^2=.064$).

1097

1098 **Bandit effect on response times**

1099 There was no difference in response times between bandits in the repeated-measures ANOVA (bandit main effect:
1100 $F(1.78, 99.44)=1.634, p=.203, \eta^2=.028$; Figure 3 – Figure supplement 1).

1101

1102 **Interaction effects on response times**

1103 When looking at the 1st choice in both conditions, no differences were evident in RT in the repeated-measures
1104 ANOVA with a between-subject factor drug group and within-subject factors horizon and bandit (bandit main
1105 effect: $F(1.71,92.46)=1.203, p=.3, \eta^2=.022$; horizon main effect: $F(1,54)=.71, p=.403, \eta^2=.013$; drug main effect:
1106 $F(2,54)=2.299, p=.11, \eta^2=.078$; drug-by-bandit interaction: $F(3.42,92.46)=.431, p=.757, \eta^2=.016$; drug-by-horizon
1107 interaction: $F(2,54)=.204, p=.816, \eta^2=.008$; bandit-by-horizon interaction: $F(1.39,75.01)=.298, p=.662, \eta^2=.005$;
1108 drug-by-bandit-by-horizon interaction: $F(2.78,75.01)=1.015, p=.387, \eta^2=.036$).

1109 In the long horizon, when looking at all 6 samples, no differences were evident in RT between drug group in the
1110 repeated-measures ANOVA with a between-subject factor drug group and within-subject factors bandits and
1111 samples (drug main effect: $F(2,56)=.542, p=.585, \eta^2=.019$). There was an effect of bandit (bandit main effect:
1112 $F(1.61,90.12)=7.137, p=.003, \eta^2=.113$), of sample (sample main effect: $F(1.54,86.15)=427.047, p<.001, \eta^2=.884$)
1113 and an interaction between the two (bandit-by-sample interaction: $F(3.33,186.41)=4.789, p=.002, \eta^2=.079$; drug-by-
1114 bandit interaction: $F(3.22,90.12)=.525, p=.679, \eta^2=.018$; drug-by-sample interaction: $F(3.08,86.15)=1.039, p=.381,$
1115 $\eta^2=.036$; drug-by-bandit-by-sample interaction: $F(6.66,186.41)=.645, p=.71, \eta^2=.023$). Further analysis (not
1116 corrected for multiple comparisons) revealed that the interaction between bandit and sample reflected the fact that
1117 when looking at samples individually, there was a bandit main effect in the 2nd sample (bandit main effect:
1118 $F(1.27,70.88)=27.783, p<.001, \eta^2=.332$; drug main effect: $F(2,56)=.201, p=.819, \eta^2=.007$; drug-by-bandit
1119 interaction: $F(2.53,70.88)=.906, p=.429, \eta^2=.031$) and in the 3rd sample (bandit main effect: $F(1.23,68.93)=21.318,$
1120 $p<.001, \eta^2=.276$; drug main effect: $F(2,56)=.102, p=.903, \eta^2=.004$; drug-by-bandit interaction: $F(2.46,68.93)=.208,$
1121 $p=.855, \eta^2=.007$), but not in the other samples (1st sample: drug main effect: $F(2,56)=1.108, p=.337, \eta^2=.038$; bandit
1122 main effect: $F(2,112)=.339, p=.713, \eta^2=.006$; drug-by-bandit interaction: $F(4,112)=.414, p=.798, \eta^2=.015$; 4th
1123 sample: (drug main effect: $F(2,56)=.43, p=.652, \eta^2=.015$; bandit main effect: $F(1.36,76.22)=1.348, p=.259, \eta^2=.024$;
1124 drug-by-bandit interaction: $F(2.72,76.22)=.396, p=.737, \eta^2=.014$; 5th sample: drug main effect: $F(2,56)=.216,$
1125 $p=.806, \eta^2=.008$; bandit main effect: $F(1.25,69.79)=.218, p=.696, \eta^2=.004$; drug-by-bandit interaction:
1126 $F(2.49,69.79)=.807, p=.474, \eta^2=.028$; 6th sample: drug main effect: $F(2,56)=1.026, p=.365, \eta^2=.035$; bandit main
1127 effect: $F(1.05,58.81)=.614, p=.444, \eta^2=.011$; drug-by-bandit interaction: $F(2.1,58.81)=1.216, p=.305, \eta^2=.042$). In
1128 the 2nd sample, the high-value bandit was chosen faster (high-value bandit vs low-value bandit : $t(59)=-5.736,$
1129 $p<.001, d=.917$; high-value bandit vs novel bandit: $t(59)=-6.24, p<.001, d=.599$) and the low-value bandit was
1130 chosen slower (low-value bandit vs novel bandit: $t(59)=3.756, p<.001, d=.432$). In the 3rd sample, the low-value
1131 bandit was chosen slower (high-value bandit vs low-value bandit : $t(59)=-5.194, p<.001, d=.571$; low-value bandit
1132 vs novel bandit: $t(59)=4.448, p<.001, d=.49$; high-value bandit vs novel bandit: $t(59)=-1.834, p=.072, d=.09$).

1133

1134 **Horizon effect on response times**

1135 There were no differences in RT between horizon conditions in the repeated-measures ANOVA with the between-
1136 subject factor drug group, the within-subject factor horizon condition and the covariates WASI and PANAS
1137 negative score (horizon main effect: $F(1, 54)=1.443, p=.235, \eta^2=.026$; drug main effect: $F(2, 54)=1.625, p=.206,$
1138 $\eta^2=.057$; drug-by-horizon interaction: $F(2, 54)=.431, p=.652, \eta^2=.016$. In the long horizon, the RT decreased with
1139 each sample (sample main effect: $F(1.36, 73.5)=13.626, p<.001, \eta^2=0.201$; Pairwise comparisons: sample 1 vs 2:
1140 $t(59)=20.968, p<.001, d=2.73$; sample 2 vs 3: $t(59)=11.825, p<.001, d=1.539$; sample 3 vs 4: $t(59)=7.862, p<.001,$
1141 $d=1.024$; sample 4 vs 5: $t(59)=4.117, p<.001, d=1.539$; sample 5 vs 6: $t(59)=2.646, p=.01, d=1.024$; Figure 2 –
1142 Figure supplement 1b).

1143

1144 **PANAS**

1145 The Positive Affect and Negative Affect scale (PANAS; (50)) was completed 50 minutes after the 2nd drug
1146 administration and 10 minutes prior to the task. Groups had similar positive affect but differed in negative affect (cf.
1147 Appendix 2 Table 1), driven by a higher score in the placebo group (pairwise comparisons: placebo vs propranolol:
1148 $t(56)=2.801$, $p=.007$, $d=.799$; amisulpride vs placebo: $t(56)=-2.096$, $p=.041$, $d=.557$; amisulpride vs propranolol:
1149 $t(56)=.669$, $p=.506$, $d=.383$). It is unclear whether this difference was driven by the drug manipulation, but similar
1150 studies have not reported such an effect (e.g. (36, 59, 61, 62, 75)). We controlled for a possible influence of these
1151 measures in all our analyses.
1152

1153 **Physiological effects**

1154 Heart rate, systolic and diastolic pressure were obtained at 3 time points: at the beginning of the experiment before
1155 giving the drug ('at arrival'), after giving the drug just before the task ('pre-task'), and after finishing task and
1156 questionnaires ('post-task'). The post-task heart rate was lower for participants who received propranolol compared
1157 to the other 2 groups (1-way ANOVA: $F(2, 55)=7.249$, $p=.002$, $\eta^2=.209$; cf. Appendix 2 Table 2). A two-way
1158 ANOVA with the between-subject factor of drug group and within-subject factor of time (all three time points),
1159 showed a time-dependent decrease in heart rate ($F(1.74, 95.97)=99.341$, $p<.001$, $\eta^2=.644$), in systolic pressure ($F(2,$
1160 $110)=8.967$, $p<.001$, $\eta^2=.14$) and in diastolic pressure ($F(2, 110)=.874$, $p=.42$, $\eta^2=.016$), indicating subjects relaxed
1161 across the course of the study. Those reductions did not differ between drug group (drug main effect: heart rate: $F(2,$
1162 $55)=1.84$, $p=.169$, $\eta^2=.063$; systolic pressure: $F(2, 55)=1.08$, $p=.347$, $\eta^2=.038$; diastolic pressure: $F(2, 55)=.239$,
1163 $p=.788$, $\eta^2=.009$; drug-by-time interaction: heart rate: $F(3.49, 95.97)=1.928$, $p=.121$, $\eta^2=.066$; systolic pressure: $F(4,$
1164 $110)=1.6$, $p=.179$, $\eta^2=.055$; diastolic pressure: $F(4, 110)=.951$, $p=.438$, $\eta^2=.033$).
1165

1166 **Task performance score**

1167 The performance did not differ between drug groups (total score: drug main effect: $F(2, 5)=2.313$, $p=.109$, $\eta^2=.079$)
1168 but it was increased in subjects with higher IQ scores (WASI main effect: $F(1, 54)=17.172$, $p<.001$, $\eta^2=.241$).
1169 In the long horizon, the score increased with each sample (sample main effect: $F(3.12, 174.97)=103.469$, $p<.001$,
1170 $\eta^2=.649$; Pairwise comparisons: sample 1 vs 2: $t(59)=-6.737$, $p<.001$, $d=.877$; sample 2 vs 3: $t(59)=-3.69$, $p<.001$,
1171 $d=.48$; sample 3 vs 4: $t(59)=-5.167$, $p<.001$, $d=.673$; sample 4 vs 5: $t(59)=-2.832$, $p=.006$, $d=.48$; sample 5 vs 6:
1172 $t(59)=-2.344$, $p=.022$, $d=.673$; Figure 2 – Figure supplement 1a). The increase in reward was larger in trials where
1173 the first draw was exploratory (linear regression slope coefficient: $\text{mean}=.118$, $\text{sd}=.038$) compared to when it was
1174 exploitative (linear regression slope coefficient: $\text{mean}=.028$, $\text{sd}=.041$; t-tests for slope coefficients: $t(58)=-12.161$,
1175 $p<.001$, $d=-1.583$; Figure 2 - Figure supplement 1d), suggesting that exploration was used beneficially and subjects
1176 benefitted from their initial exploration.
1177

1178 **Dopamine effect on high-value bandit sampling frequency**

1179 The amisulpride group had a marginal tendency towards selecting the high-value bandit, meaning that they were
1180 disposed to exploit more overall (propranolol group excluded: horizon main effect: $F(1, 35)=3.035$, $p=.09$, $\eta^2=.08$;
1181 drug main effect: $F(1, 35)=3.602$, $p=.066$, $\eta^2=.093$; drug-by-horizon interaction: $F(1, 35)=2.15$, $p=.151$, $\eta^2=.058$).
1182 This trend effect was not observed when all 3 groups were included (horizon main effect: $F(1, 54)=3.909$, $p=.053$,
1183 $\eta^2=.068$; drug main effect: $F(2, 54)=1.388$, $p=.258$, $\eta^2=.049$; drug-by-horizon interaction: $F(2, 54)=.834$, $p=.44$,
1184 $\eta^2=.03$).

1185 **Gender effects**

1186 When adding gender as a between-subjects variable in the repeated-measures ANOVAs, none of the main results
1187 changed. Interestingly, we observed a drug-by-gender interaction in the prior variance σ_0 (drug-by-gender
1188 interaction: $F(2, 51)=5.914$, $p=.005$, $\eta^2=.188$; Figure 5 – Figure supplement 2), driven by the fact that, female
1189 subjects in the placebo group had a larger average σ_0 (across both horizon conditions) compared to males
1190 ($t(20)=2.836$, $p=.011$, $d=1.268$), whereas male subjects have a larger σ_0 compared to females in the amisulpride
1191 group, ($t(19)=-2.466$, $p=.025$, $d=1.124$; propranolol group: $t(20)=-.04$, $p=.969$, $d=.018$). This suggests that in a
1192 placebo setting, females are on average more uncertain about an option's expected value, whereas in a dopamine
1193 blockade setting males are more uncertain. Besides this effect, we observed a trend-level significance in response
1194 times (RT), driven primarily by female subjects tending to have a faster RT in the long horizon compared to male
1195 subjects (gender main effect: $F(1, 51)=3.54$, $p=.066$, $\eta^2=.065$).

1196 **Horizon and drug effects without covariate**

1197 When analysing the results without correcting for IQ (WASI) and negative affect (PANAS), similar results are
1198 obtained. The high-value bandit is picked more in the short-horizon condition indicating exploitation ($F(1,$
1199 $56)=44.844, p<.001, \eta^2=.445$), whereas the opposite phenomenon is observed in the low-value bandit ($F(1,$
1200 $56)=24.24, p<.001, \eta^2=.302$) and the novel bandit (horizon main effect: $F(1, 56)=30.867, p<.001, \eta^2=.355$),
1201 indicating exploration. In line with these results, the model parameters for value-free random exploration (ϵ : $F(1,$
1202 $56)=10.362, p=.002, \eta^2=.156$) and novelty exploration (η : $F(1, 56)=38.103, p<.001, \eta^2=.405$) are larger in the long
1203 compared to the short horizon condition. Additionally, noradrenaline blockade reduces value-free random
1204 exploration as can be seen in the two behavioural signatures, frequency of picking the low-value bandit ($F(2,$
1205 $56)=2.523, p=.089, \eta^2=.083$; Pairwise comparisons: placebo vs propranolol: $t(40)=2.923, p=.005, d=.654$;
1206 amisulpride vs placebo: $t(38)=-.587, p=.559, d=.133$; amisulpride vs propranolol: $t(38)=2.171, p=.034, d=.496$), and
1207 in the consistency ($F(2, 56)=3.596, p=.034, \eta^2=.114$; Pairwise comparisons: placebo vs propranolol: $t(40)=-3.525,$
1208 $p=.001, d=.788$; amisulpride vs placebo: $t(38)=1.107, p=.272, d=.251$; amisulpride vs propranolol: $t(38)=-2.267,$
1209 $p=.026, d=.514$), as well as in the model parameter for value-free random exploration (ϵ : $F(2, 56)=3.205, p=.048,$
1210 $\eta^2=.103$; Pairwise comparisons: placebo vs propranolol: $t(40)=3.177, p=.002, d=.71$; amisulpride vs placebo:
1211 $t(38)=.251, p=.802, d=.057$; amisulpride vs propranolol: $t(38)=2.723, p=.009, d=.626$).

1212 1213 **Horizon and drug effects with heart rate as covariate**

1214 When analysing results but now correcting for the post-experiment heart rate (cf. Appendix 2 Table 1) in addition to
1215 IQ (WASI) and negative affect (PANAS), we obtained similar results. Noradrenaline blockade reduced value-free
1216 random exploration as seen in two behavioural signatures, frequency of picking the low-value bandit ($F(2, 52)=$
1217 $4.014, p=.024, \eta^2=.134$; Pairwise comparisons:(placebo vs propranolol: $t(40)= 2.923, p=.005, d=.654$; amisulpride
1218 vs propranolol: $t(38)= 2.171, p=.034, d=.496$; amisulpride vs placebo: $t(38)= -.587, p=.559, d=.133$), and
1219 consistency ($F(2, 52)= 5.474, p=.007, \eta^2=.174$; Pairwise comparisons: placebo vs propranolol: $t(40)= -3.525,$
1220 $p=.001, d=.788$; amisulpride vs propranolol: $t(38)= -2.267, p=.026, d=.514$; amisulpride vs placebo: $t(38)= 1.107,$
1221 $p=.272, d=.251$), as well as in a model parameter for value-free random exploration (ϵ : $F(2, 52)= 4.493, p=.016,$
1222 $\eta^2=.147$; Pairwise comparisons: placebo vs propranolol: $t(40)= 3.177, p=.002, d=.71$; amisulpride vs propranolol:
1223 $t(38)= 2.723, p=.009, d=.626$; amisulpride vs placebo: $t(38)=.251, p=.802, d=.057$).

1224 1225 **Other model results**

1226 When analysing the fitted parameter values of both the 2nd winning model (UCB + ϵ + η) and 3rd winning model
1227 (hybrid + ϵ + η), similar results pertain. Thus, a value-free random exploration parameter was reduced following
1228 noradrenaline blockade in the 2nd winning model (ϵ : $F(2, 54)=4.503, p=.016, \eta^2=.143$; Pairwise comparisons:
1229 placebo vs propranolol: $t(38)=2.185, p=.033, d=.386$; amisulpride vs propranolol: $t(40)=1.724, p=.089, d=.501$;
1230 amisulpride vs placebo: $t(40)=-.665, p=.508, d=.151$) and was affected at a trend-level significance in the 3rd
1231 winning model (ϵ : $F(2, 54)=3.04, p=.056, \eta^2=.101$). These results highlight our finding that value-free random
1232 exploration is modulated by noradrenaline and additionally demonstrates this is independent of the complex
1233 exploration strategy used as well as the value function.

1234 1235 **Bandit combination effect**

1236 Behavioural results were analysed additionally for each bandit combination separately. The high-value bandit was
1237 chosen more when there was no novel bandit (pairwise comparisons: [certain-standard, standard, low] vs [certain-
1238 standard, standard, novel]: $t(59)=15.122, p<.001, d=1.969$; [certain-standard, standard, low] vs [certain-standard,
1239 novel, low]: $t(59)=12.905, p<.001, d=2.389$; [certain-standard, standard, low] vs [standard, novel, low]:
1240 $t(59)=18.348, p<.001, d=1.68$), and less when its value was less certain ([standard, novel, low] vs [certain-standard,
1241 standard, novel]: $t(59)=-6.986, p<.001, d=.407$; [standard, novel, low] vs [certain-standard, novel, low] : $t(59)=-$
1242 $5.44, p<.001, d=.708$; bandit combination main effect: $F(1.81, 101.33)=237.051, p<.001, \eta^2=.809$; [certain-standard,
1243 standard, novel] vs [certain-standard, novel, low]: $t(59)=.364, p=.717, d=.909$; Figure 3 – Figure supplement 2a).
1244 The novel bandit was chosen most often when the high-value bandit was less certain, then when the high-value
1245 bandit was more certain and was chosen least when both certain and certain standard bandits were present
1246 ([standard, novel, low] vs [certain-standard, novel, low]: $t(59)=5.001, p<.001, d=.651$; [standard, novel, low] vs
1247 [certain-standard, standard, novel]: $t(59)=9.414, p<.001, d=1.226$; [certain-standard, novel, low] vs [certain-
1248 standard, standard, novel]: $t(59)=4.146, p<.001, d=.54$; bandit combination main effect: $F(2, 112)=42.44, p<.001,$
1249 $\eta^2=.431$; Figure 3 – Figure supplement 2b). The low-value bandit was chosen less when the high-value bandit was
1250 more certain ([certain-standard, novel, low] vs [certain-standard, standard, low]: $t(59)=-2.731, p=.008, d=.356$;
1251 [certain-standard, novel, low] vs [standard, novel, low]: $t(59)=-1.958, p=.055, d=.255$; bandit combination main

1252 effect: $F(1.66, 92.74)=4.534$, $p=.019$, $\eta^2=.075$; [certain-standard, standard, low] vs [standard, novel, low]:
1253 $t(59)=1.32$, $p=.192$, $d=.172$; Figure 3 – Figure supplement 2c).

1254

1255 **Other effects on choice consistency**

1256 Our results demonstrate a drug-by-horizon interaction on choice consistency ($F(2, 54)=3.352$, $p=.042$, $\eta^2=.110$;
1257 Figure 3c), mainly driven by the fact that frequency of selecting the same option is increased in the long (compared
1258 to the short) horizon in the amisulpride group, while there is no significant horizon difference in the other two drug
1259 groups (pairwise comparison for horizon effect: amisulpride group: $t(19)=2.482$, $p=.023$, $d=.569$; propranolol group:
1260 $t(20)=-1.91$, $p=.071$, $d=.427$; placebo group: $t(20)=.505$, $p=.619$, $d=.113$). It is not entirely clear why catecholamines
1261 would increase the differentiation between the horizon conditions and this relatively weak effect should be
1262 replicated before interpreting.

1263

1264 **Stand-alone heuristic models**

1265 We also analysed stand-alone heuristic models, in which there is no value computation (value of each bandit i :
1266 $V_i = 0$). The held-out data likelihood for such heuristic model combined with novelty exploration had a mean of
1267 $m=0.367$ ($sd=0.005$). The model in which we added value-free random exploration on top of novelty exploration
1268 had a mean of $m=0.384$ ($sd=0.006$). These models performed poorly, although better than chance level. Importantly,
1269 adding value-free random exploration improved performance. This highlights that subjects' combine complex and
1270 heuristic modules in exploration.

1271 **Appendix 2**

	Propranolol	Placebo	Amisulpride	
Gender (M/F)	10/10	10/10	10/9	
Age	22.80 (3.59)	23.80 (4.23)	23.05 (3.01)	F(2,56)=.404, p=.669, η^2 =.014
Intellectual abilities	22.8 (1.85)	22.6 (3.70)	24.37 (2.45)	F(2,56)=2.337, p=.106, η^2 =.077
Positive affect	24.55 (8.99)	28.90 (7.56)	29.58 (10.21)	F(2,56)=1.832, p=.170, η^2 =.061
Negative affect	10.65 (.81)	12.75 (3.63)	11.16 (1.71)	F(2,56)=4.259, p=.019, η^2 =.132

1272 **Appendix 2 Table 1.**

1273 Characteristics of drug groups. The drug groups did not differ in gender, age, nor in intellectual abilities (adapted
 1274 WASI matrix test). Groups differed in negative affect (PANAS), driven by a higher score in the placebo group
 1275 (pairwise comparisons: placebo vs propranolol: $t(56)=2.801$, $p=.007$, $d=.799$; amisulpride vs placebo: $t(56)=-2.096$,
 1276 $p=.041$, $d=.557$; amisulpride vs propranolol: $t(56)=.669$, $p=.506$, $d=.383$). For more details cf. Appendix 1. Mean
 1277 (SD).

		Propranolol	Placebo	Amisulpride	
Heart rate (BPM)	At arrival	74.9 (10.8)	77,2 (12,6)	77.7 (13.8)	F(2, 55)=.290, p=.749, η^2 =.010
	Pre-task	62,6 (8,5)	65,8 (8,3)	64,6 (9,8)	F(2, 55)=.667, p=.517, η^2 =.024
	Post-task	55,7 (6,7)	64,4 (6,9)	63,4 (10,0)	F(2, 55)=7.249, p=.002, η^2 =.209
Systolic blood pressure	At arrival	117,2 (10,4)	115,0 (9,7)	117,9 (9,7)	F(2, 55)=.438, p=.648, η^2 =.016
	Pre-task	109,4 (9,2)	111,8 (8,6)	114,9 (8,6)	F(2, 55)=1.841, p=.168, η^2 =.063
	Post-task	109,5 (8,2)	113,9 (11,3)	114,6 (9,3)	F(2, 55)=1.584, p=.214, η^2 =.054
Diastolic blood pressure	At arrival	71,5 (7,8)	71,2 (6,7)	72,3 (6,7)	F(2, 55)=.115, p=.891, η^2 =.004
	Pre-task	68,3 (7,0)	71,1 (10,6)	72,0 (5,9)	F(2, 55)=1.111, p=.337, η^2 =.039
	Post-task	70,8 (7,3)	70,9 (8,0)	70,3 (6,6)	F(2, 55)=.037, p=.964, η^2 =.001

1278 **Appendix 2 Table 2.**

1279 Physiological effects on drug groups. The drug groups also differed in post-experiment heart rate, driven by lower
1280 values in the propranolol group (pairwise comparisons: placebo vs propranolol: $t(55)=3.5$, $p=.001$, $d=1.293$;
1281 amisulpride vs placebo: $t(55)= -.394$, $p=.695$, $d=.119$; amisulpride vs propranolol: $t(55)=3.013$, $p=.004$, $d=.921$).
1282 For detailed statistics and analysis accounting for this cf. Appendix 1. Mean (SD).

	Horizon	Mean (sd)	Two-way repeated-measures ANOVA
			Main effect of horizon
Expected value	short	6.368 (0.335)	F(1, 56)=19.457, p<.001, η^2 =.258
	long	6.221 (0.379)	
Initial samples	short	1.282 (0.247)	F(1, 56)=58.78, p<.001, η^2 =.512
	long	1.084 (0.329)	
Score (1 st sample)	short	5.904 (0.192)	F(1, 56)=58.78, p<.001, η^2 =.512
	long	5.82 (0.182)	
Score (average)	short	5.904 (0.192)	F(1, 56)=103.759, p<.001, η^2 =.649
	long	6.098 (0.222)	

1283 **Appendix 2 Table 3.**

1284 Table of statistics and behavioural values of Figure 2. All of those measures were modulated by the horizon
1285 condition.

	Horizon	Mean (sd)			Two-way repeated-measures ANOVA			
		Amisulpride	Placebo	Propranolol	Main effect		Interaction	
High-value bandit	short	54.55 (8.87)	49.38 (9.10)	50.98 (11.4)	D	F(2, 54)=1.388, p=.258, η^2 =.049	DH	F(2, 54)=.834, p=.440, η^2 =.030
	long	41.90 (8.47)	44.10 (13.88)	41.90 (13.57)	H	F(1, 54)=3.909, p=.053, η^2 =.068	HW	F(1, 54)=13.304, p=.001, η^2 =.198
Low-value bandit	short	3.32 (2.33)	4.28 (2.98)	2.50 (2.48)	D	F(2, 54)=7.003, p=.002, η^2 =.206	DH	F(2, 54)=2.154, p=.126, η^2 =.074
	long	5.45 (3.76)	5.35 (3.40)	3.45 (2.18)	H	F(1, 54)=4.069, p=.049, η^2 =.070	HW	F(1, 54)=1.199, p=.278, η^2 =.022
Novel bandit	short	36.87 (9.49)	39.02 (10.94)	40.15 (12.43)	D	F(2, 54)=1.498, p=.233, η^2 =.053	DH	F(2, 54)=.542, p=.584, η^2 =.020
	long	46.82 (12.1)	43.62 (16.27)	48.55 (16.59)	H	F(1, 54)=5.593, p=.022, η^2 =.094	HW	F(1, 54)=13.897, p<.001, η^2 =.205
Consistency	short	64.16 (12.27)	62.70 (12.59)	73.00 (11.33)	D	F(2, 54)=7.154, p=.002, η^2 =.209	DH	F(2, 54)=3.352, p=.042, η^2 =.110
	long	68.11 (10.34)	64.00 (8.93)	70.55 (9.91)	H	F(1, 54)=1.333, p=.253, η^2 =.024	HW	F(1, 54)=.409, p=.525, η^2 =.008

1286 **Appendix 2 Table 4.**

1287 Table of statistics and behavioural measure values of Figure 3. The drug groups differed in low-value bandit picking
1288 frequency (pairwise comparisons: placebo vs propranolol: $t(40)=2.923$, $p=.005$, $d=.654$; amisulpride vs placebo:
1289 $t(38)=-.587$, $p=.559$, $d=.133$; amisulpride vs propranolol: $t(38)=2.171$, $p=.034$, $d=.496$) and choice consistency
1290 (placebo vs propranolol: $t(40)=-3.525$, $p=.01$, $d=.788$; amisulpride vs placebo: $t(38)=1.107$, $p=.272$, $d=.251$;
1291 amisulpride vs propranolol: $t(38)=-2.267$, $p=.026$, $d=.514$). The main effect is either of drug group (D) or of horizon
1292 (H). The interaction is either drug-by-horizon (DH) or horizon-by-WASI (measure of IQ; HW).

	Model	Thompson				UCB				Hybrid			
			+ ϵ	+ η	+ ϵ + η		+ ϵ	+ η	+ ϵ + η		+ ϵ	+ η	+ ϵ + η
Parameters	Horizon independent	Q_0	Q_0	Q_0	Q_0	Q_0	Q_0	Q_0	Q_0	w, Q_0	w, Q_0	w, Q_0	w, Q_0
	Horizon dependent	σ_0	σ_0, ϵ	σ_0, η	σ_0, ϵ, η	γ, β	γ, β, ϵ	γ, β, η	$\gamma, \beta, \epsilon, \eta$	σ_0, γ, β	$\sigma_0, \gamma, \beta, \epsilon$	$\sigma_0, \gamma, \beta, \eta$	$\sigma_0, \gamma, \beta, \epsilon, \eta$
Model selection	Mean held-out data likelihood	50.2 (8.1)	52.7 (7.1)	52.2 (8.7)	55.3 (8.4)	52.9 (8.0)	52.9 (8.0)	53.4 (8.1)	55.1 (8.8)	53.5 (8.1)	53.8 (8.4)	55.0 (8.4)	55.1 (8.5)
	Subjects' for which model fits best (out of 12)	0	3	2	20	0	0	1	20	0	0	7	6
	Subjects' for which model fits best (out of 3 best)	-	-	-	27	-	-	-	22	-	-	-	10

1293 **Appendix 2 Table 5.**

1294 Table of parameters used for each model compared during model selection (Figure 4). Each of the 12 columns
1295 indicate a model. The three 'main models' studied were the Thompson model, the UCB model and a hybrid of both.
1296 Variants were then created by adding the ϵ -greedy parameter, the novelty bonus and a combination of both. All the
1297 parameters besides Q_0 and w were fitted to each horizon separately. Parameters: Q_0 =prior mean (initial estimate of a
1298 bandits mean); σ_0 =prior variance (uncertainty about Q_0); w =contribution of UCB vs Thompson; γ =information
1299 bonus; β =softmax inverse temperature; ϵ = ϵ -greedy parameter (stochasticity); η =novelty bonus. Model selection
1300 measures include the cross-validation held-out data likelihood averaged over subjects, mean (SD), as well as the
1301 subject count for which this model performed better over either 12 models or over the 3 best models.
1302

	Horizon	Mean (sd)			Two-way repeated-measures ANOVA			
		Amisulpride	Placebo	Propranolol	Main effect		Interaction	
ϵ -greedy parameter	short	0.10 (0.10)	0.12 (0.08)	0.07 (0.08)	D	F(2, 54)=6.722, p=.002, η^2 =.199	DH	F(2, 54)=1.305, p=.280, η^2 =.046
	long	0.17 (0.14)	0.14 (0.10)	0.08 (0.06)	H	F(1, 54)=1.968, p=.166, η^2 =.035	HW	F(1, 54)=6.08, p=.017, η^2 =.101
Novelty bonus η	short	2.07 (0.98)	2.26 (1.37)	2.05 (1.16)	D	F(2, 54)=.249, p=.780, η^2 =.009	DH	F(2, 54)=.03, p=.971, η^2 =.001
	long	3.24 (1.19)	3.12 (1.63)	2.95 (1.70)	H	F(1, 54)=1.839, p=.181, η^2 =.033	HW	F(1, 54)=8.416, p=.005, η^2 =.135
Prior variance σ_0	short	1.18 (0.20)	1.12 (0.43)	1.25 (0.34)	D	F(2, 54)=.060, p=.942, η^2 =.002	DH	F(2, 54)=2.162, p=.125, η^2 =.074
	long	1.41 (0.61)	1.42 (0.59)	1.21 (0.44)	H	F(1, 54)=.129, p=.721, η^2 =.002	HW	F(1, 54)=.022, p=.882, η^2 <.001
Prior mean Q_0		3.22 (1.05)	3.20 (1.36)	3.44 (1.05)	D	F(2, 54)=.118, p=.889, η^2 =.004		

1303 **Appendix 2 Table 6.**

1304 Table of statistics and fitted model parameters of Figure 5. The drug groups differed in ϵ -greedy parameter value
1305 (pairwise comparisons: placebo vs propranolol: $t(40)=3.177$, $p=.002$, $d=.71$; amisulpride vs placebo: $t(38)=.251$,
1306 $p=.802$, $d=.057$; amisulpride vs propranolol: $t(38)=2.723$, $p=.009$, $d=.626$). The main effect is either of drug group
1307 (D) or of horizon (H). The interaction is either drug-by-horizon (DH) or horizon-by-WASI (measure of IQ; HW).

	Horizon	Low exploration	High exploration	Additional parameters
Value-free random exploration	short	$\epsilon = 0.1$	$\epsilon = 0.2$	$\eta = 0$
	long	$\epsilon = 0.3$	$\epsilon = 0.4$	$\eta = 2$
Novelty exploration	short	$\eta = 0$	$\eta = 1$	$\epsilon = 0$
	long	$\eta = 2$	$\eta = 3$	$\epsilon = 0.2$
Thompson-sampling exploration	short	$\sigma_0 = 0.8$	$\sigma_0 = 1.2$	$\eta = 0, \epsilon = 0$
	long	$\sigma_0 = 1.6$	$\sigma_0 = 2$	$\eta = 2, \epsilon = 0.2$
UCB exploration	short	$\gamma = 0.1$	$\gamma = 0.3$	$\beta = 5, \epsilon = 0$
	long	$\gamma = 0.7$	$\gamma = 1.5$	$\beta = 1.5, \epsilon = 0.2$

1308 **Appendix 2 Table 7**

1309 Parameter values used for simulations on Figure 1- Figure supplement 3-5. Parameter values for high and low
1310 exploration were selected empirically from pilot and task data. Value-free random exploration and novelty
1311 exploration were simulated with an argmax decision function, which always selects the value with the highest
1312 expected value. For simulating the long (versus short) horizon condition, we assumed that not only the key value but
1313 also the other exploration strategies increased, as found in our experimental data. For each simulation $Q_0 = 5$ and
1314 unless otherwise stated, $\sigma_0 = 1.5$.

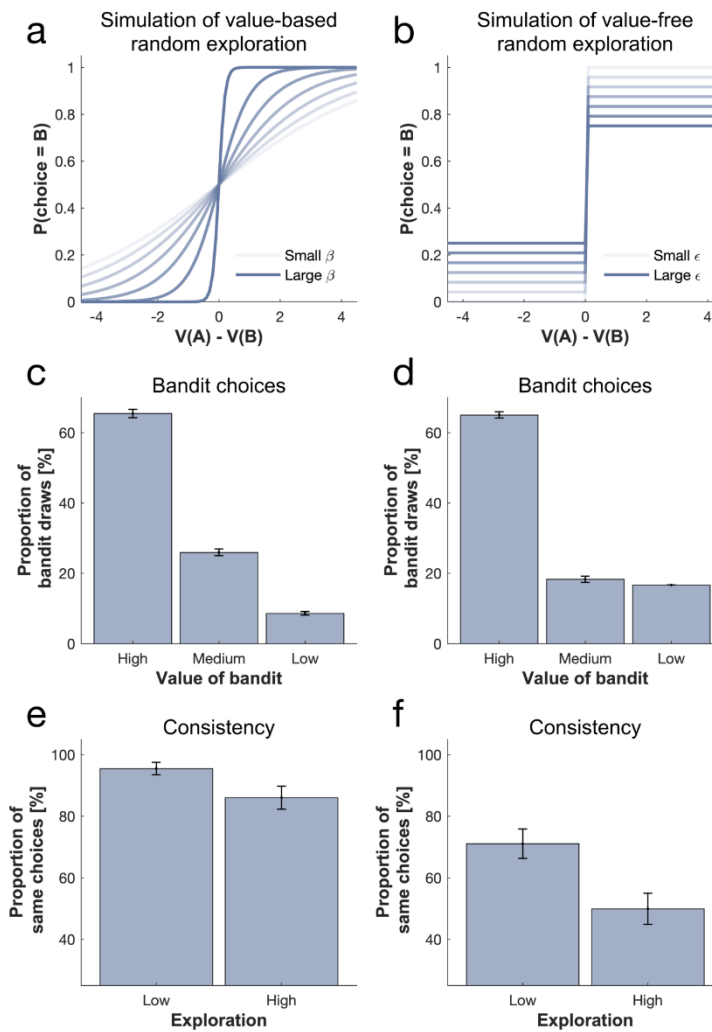


1315

1316 **Figure 1 - Figure supplement 1**

1317 Visualisation of the 9 different sizes that the apples could take. The associated rewards went from 2
1318 (small apple on the left) to 10 (big apple on the right).

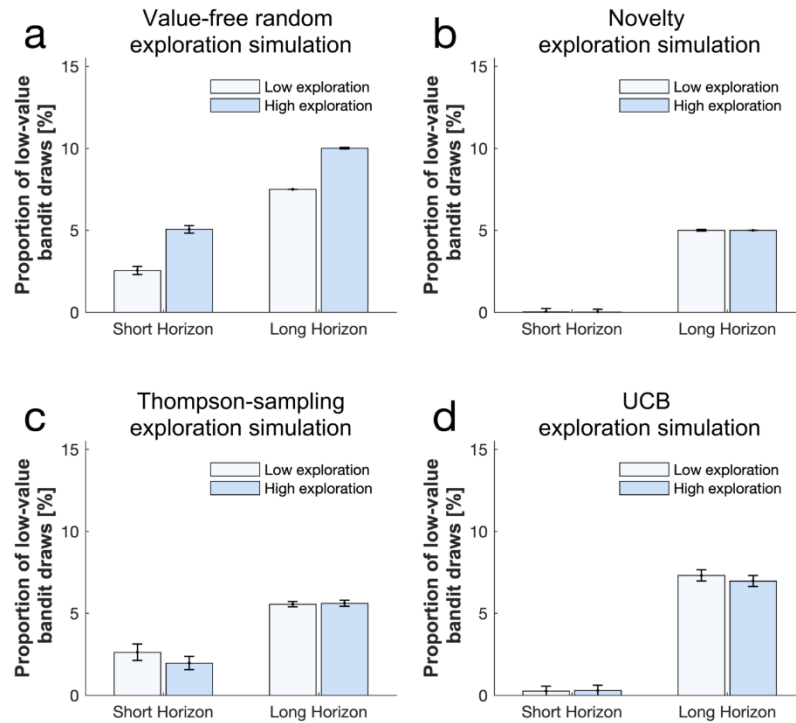
1319



1320

1321 **Figure 1 - Figure supplement 2**

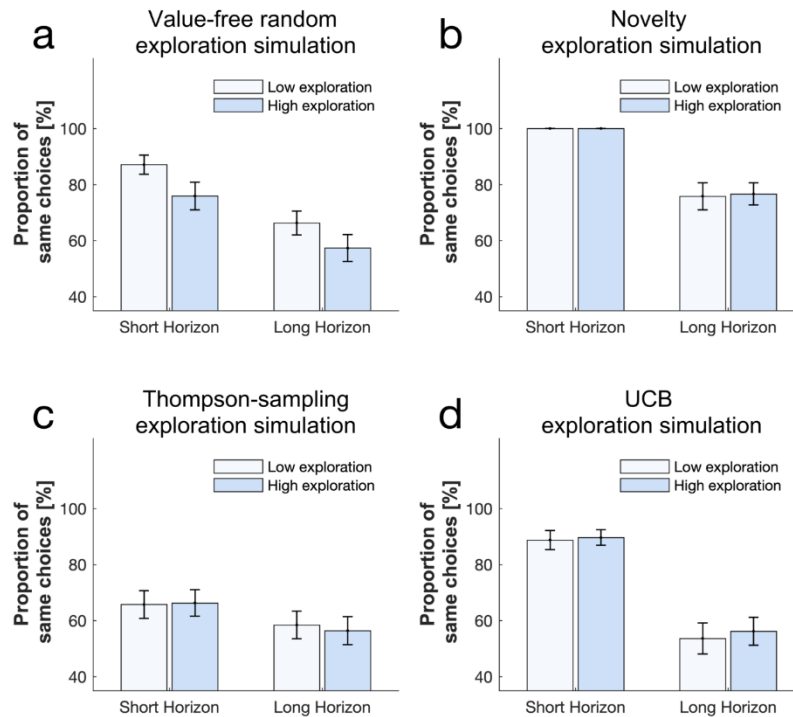
1322 Comparison of value-based (softmax) and value-free (ϵ -greedy) random exploration. (a) Changing the
 1323 softmax inverse temperature affects the slope of the sigmoid while changing the ϵ -greedy parameter (b)
 1324 affects the compression of the sigmoid. Conceptually, in a softmax exploration mode, as each bandits'
 1325 expected value is taken into account, (c) the 2nd best bandit (medium-value bandit) is favoured over one
 1326 with a lower value (low-value bandit) when injecting noise. In contrast, in an ϵ -greedy exploration mode,
 1327 (d) bandits are explored equally often irrespective of their expected value. Both simulations were
 1328 performed on trials without novel bandit. When simulating on all trials we observe that this also has a
 1329 consequence for choice consistency. (e) Choices are more consistent in a low (versus high) softmax
 1330 exploration mode (i.e. high and low values of β respectively), and similarly (f) choices are more
 1331 consistent in a low (versus high) ϵ -greedy exploration mode (i.e. low and high values of ϵ respectively).
 1332 When comparing the overall consistency of the two random exploration strategies, consistency is higher
 1333 in the value-based mode, reflecting a higher probability of (consistently) exploring the 2nd best option,
 1334 compared to an equal probability of exploring any non-optimal option (inconsistently) in the value-free
 1335 mode.



1336

1337 **Figure 1 - Figure supplement 3**

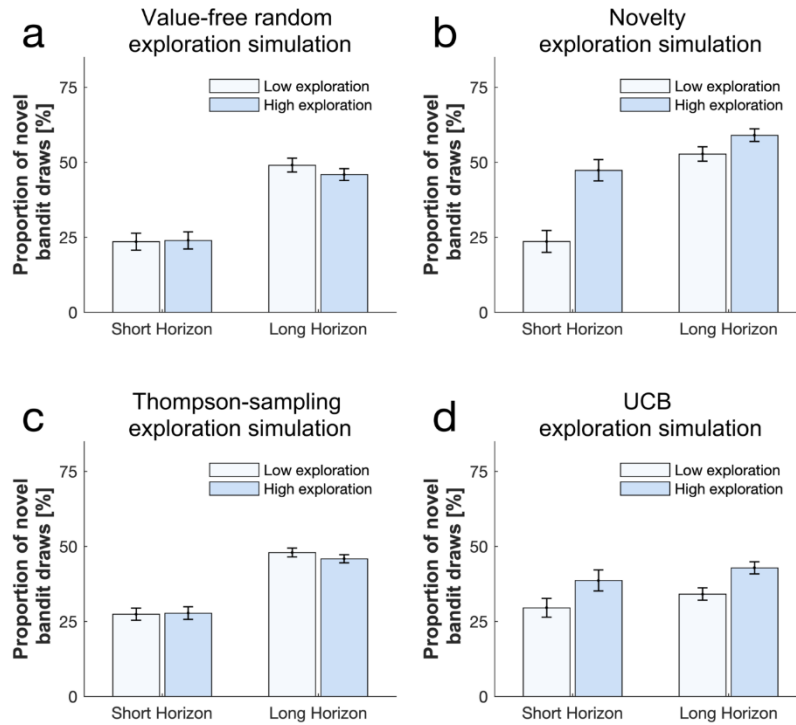
1338 Simulation illustrations of high and low exploration on the frequency of picking the low-value bandit using different
 1339 exploration strategies shows that (a) a high (versus low) value-free random exploration increases the selection of the
 1340 low-value bandit, whereas neither (b) a high (versus low) novelty exploration, (c) a high (versus low) Thompson-
 1341 sampling exploration nor (d) a high (versus low) UCB exploration affected this frequency. To illustrate the long
 1342 (versus short) horizon condition, we accommodated the fact that not only key values but also other exploration
 1343 strategies were enhanced by increasing multiple exploration strategies, as found in our experimental data (cf.
 1344 Appendix 2 Table 7 for parameter values). Please note that the difference between low and high exploration is
 1345 critical here, rather than a comparison of the absolute height of the bars between strategies (which is influences in
 1346 the models by multiple different exploration strategies). For simulations fitting participants' data, please see Figure 5
 1347 - Figure supplement 1 and 3.



1348

1349 **Figure 1 - Figure supplement 4**

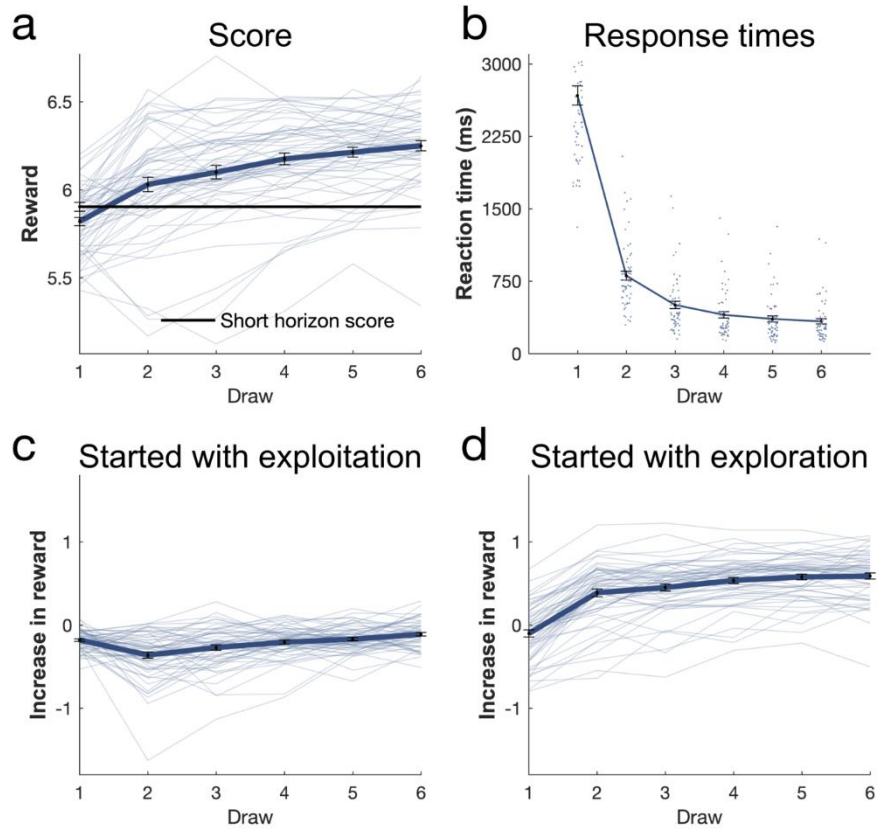
1350 Simulation illustrations of high and low exploration choice consistency using different exploration strategies shows
 1351 that (a) a high (versus low) value-free random exploration decreases the proportion of same choices, whereas neither
 1352 (b) a high (versus low) novelty exploration, (c) a high (versus low) Thompson-sampling exploration nor (d) a high
 1353 (versus low) UCB exploration affected this measure. To illustrate the long (versus short) horizon condition,
 1354 accommodated the fact that not only the key value but also other exploration strategies were enhanced by increasing
 1355 multiple exploration strategies, as found in our experimental data (cf. Appendix 2 Table 7 for parameter values).
 1356 Please note that the difference between low and high exploration is critical here, rather than a comparison of the
 1357 absolute height of the bars between strategies (which is influenced in the models by multiple different exploration
 1358 strategies). For simulations fitting participants' data, please see Figure 5 - Figure supplement 1 and 3.



1359

1360 **Figure 1 - Figure supplement 5**

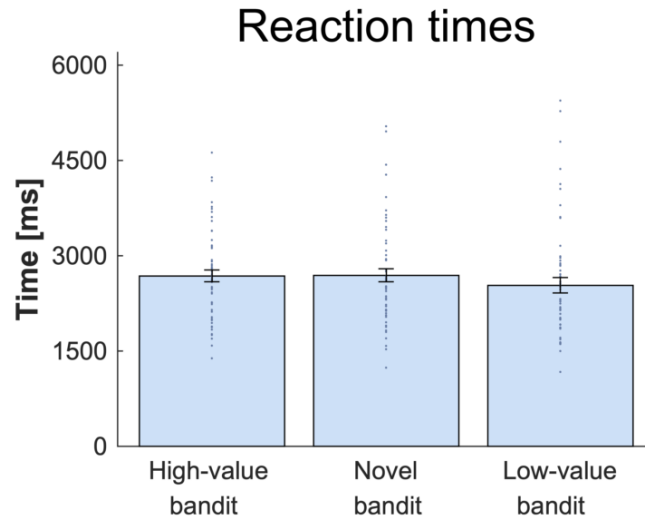
1361 Simulation illustrations of high and low exploration on the frequency of picking the novel bandit using different
 1362 exploration strategies shows that (a) a high (versus low) value-free random exploration has little effect on the
 1363 selection of the novel bandit, whereas (b) a high (versus low) novelty exploration increases this frequency. (c) A
 1364 high (versus low) Thompson-sampling exploration had little effect and (d) a high (versus low) UCB exploration
 1365 affected this frequency but to a lower extent than novelty exploration. To illustrate the long (versus short) horizon
 1366 condition, we accommodated the fact that not only the key value but also other exploration strategies were enhanced
 1367 by increasing multiple exploration strategies, as found in our experimental data (cf. Appendix 2 Table 7 for
 1368 parameter values). Please note that the difference between low and high exploration is critical here, rather than a
 1369 comparison of the absolute height of the bars between strategies (which is influenced in the models by multiple
 1370 different exploration strategies). For simulations fitting participants' data, please see Figure 5- Figure supplement 1
 1371 and 3.



1372

1373 **Figure 2 - Figure supplement 1**

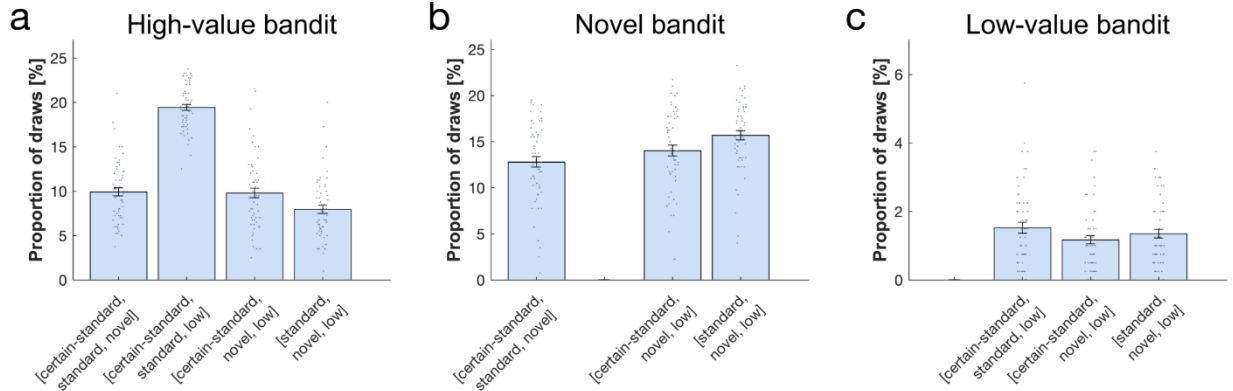
1374 Further analysis of long horizon draws. (a) The first draw in the long horizon led to a lower reward than the short
 1375 horizon, indicating more exploration, while the subsequent draws led to a higher reward indicating that this
 1376 additional information helped making better decisions in the long run. (b) The first draws' response time was the
 1377 highest and then decreased for each draw. Long horizon trials in which subjects started with (c) an exploitation draw
 1378 (choose the bandit with the highest expected value) led to little increase in reward (y-axis: difference between
 1379 obtained reward and highest reward of initial samples; linear regression slope coefficient: mean=0.118, sd=0.038),
 1380 whereas trials in which they started with (d) an exploration draw led to an large increase in reward (linear regression
 1381 slope coefficient: mean=0.028, sd=0.041). This larger increase in reward when starting by exploring (slope is
 1382 higher: $t(58)=-12.161$, $p<.001$, $d=-1.583$) indicates that the information that was gained through exploration led to
 1383 higher long-term outcomes. Data are shown as mean \pm SEM and each dot represent one subject.



1384

1385 **Figure 3 - Figure supplement 1**

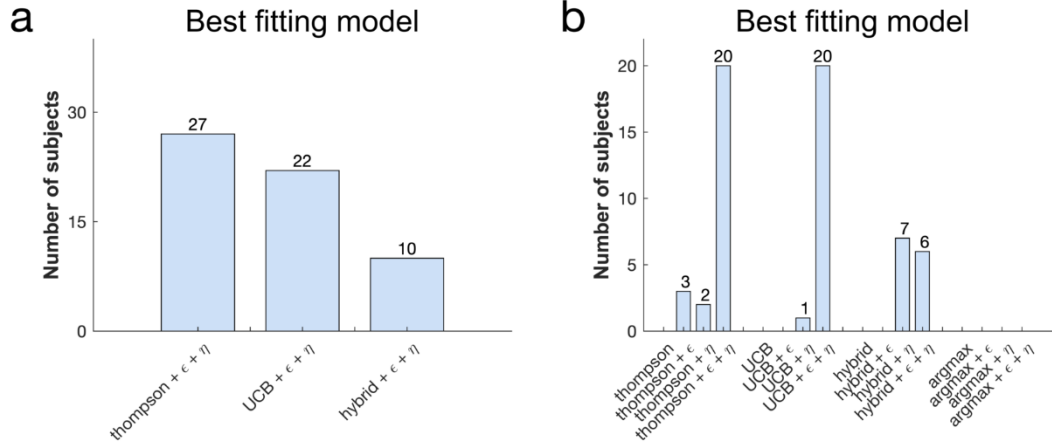
1386 Response time analysis per bandit. There was no difference in RT depending which bandit was chosen. For details
1387 and statistics cf. Appendix 1.



1388
 1389
 1390
 1391
 1392
 1393
 1394

Figure 3 - Figure supplement 2

Proportion of draws per bandit combination (x-axis). (a) The high-value bandit was picked more when there was no novel bandit, and less when the high-value bandit was less certain. (b) The novel bandit was picked the most when the high-value bandit was less certain, then when the high-value bandit was more certain, and it was picked the least when both certain and certain standard bandits were present. (c) The low-value bandit was picked less when the high-value bandit was more certain. For statistics see Appendix 1.

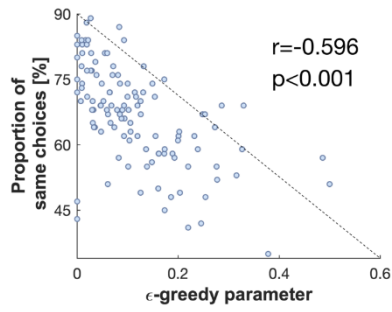
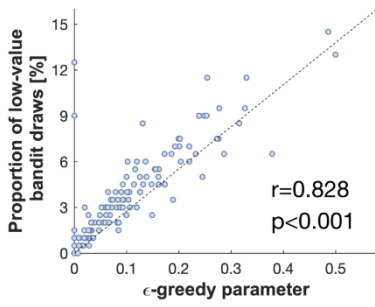


1395

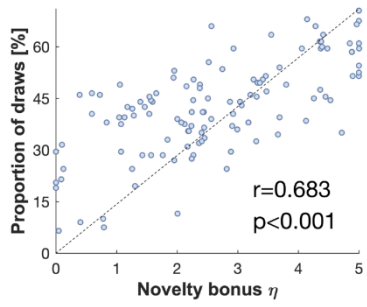
1396 **Figure 4 - Figure supplement 1**

1397 Model comparison: further evaluations. (a) The winning model at the group level (the Thompson model with both ϵ
 1398 and η) was also the one that accounted best for the largest number of subjects. (b) The Thompson+ ϵ + η model and
 1399 the UCB+ ϵ + η are equally first in subject count when comparing all models, the Thompson+ ϵ + η model is therefore
 1400 still the winning model as it has the highest average likelihood of held-out data.

a Value-free random exploration



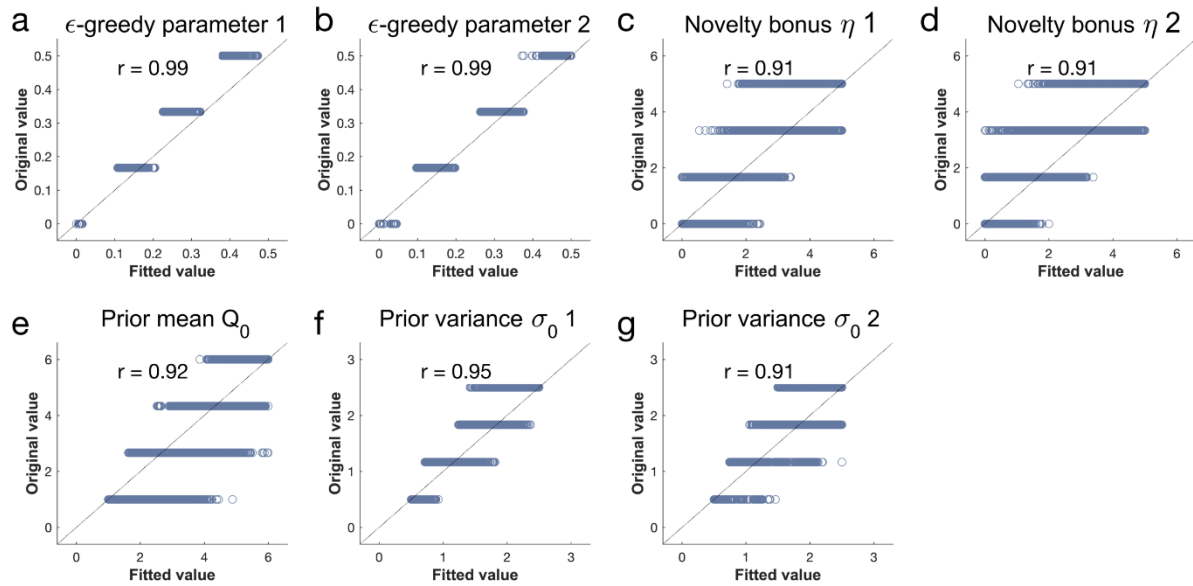
b Novelty exploration



1401

1402 **Figure 4 - Figure supplement 2**

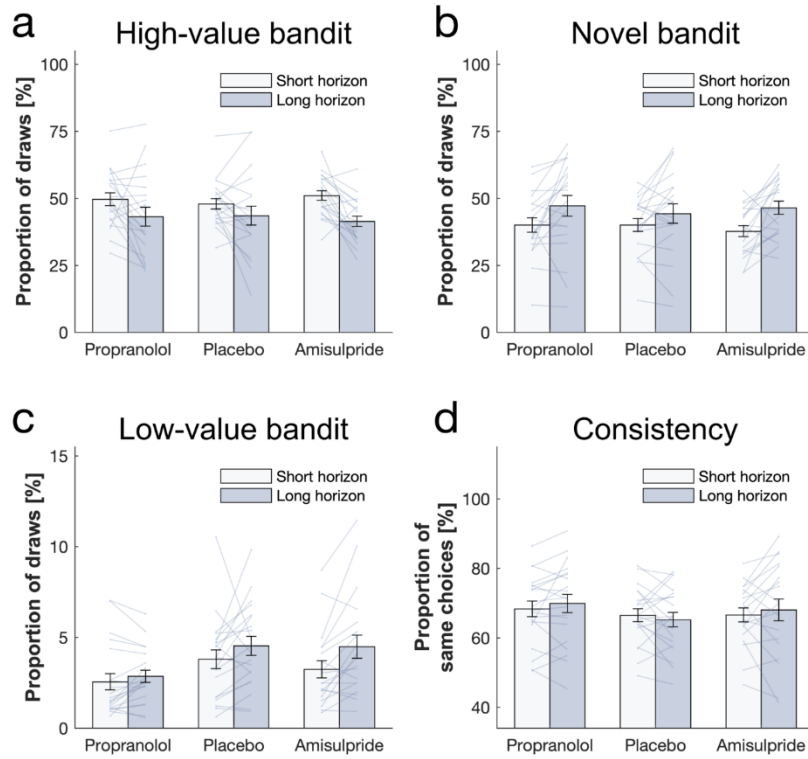
1403 Correlations between model parameters and behaviour. The behavioural indicators of (a) value-free random
1404 exploration (left panel: draws from the low-value bandit; right panel: consistency) correlated with the ϵ -greedy
1405 parameter values, and of (b) novelty exploration (draws from the novel bandit) correlated with the novelty bonus η .



1406

1407 **Figure 4 - Figure supplement 3**

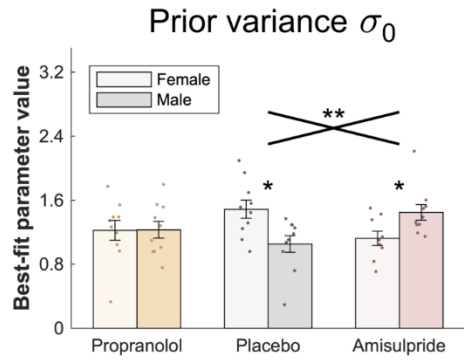
1408 Parameter recovery analysis details. For each of the 7 parameters of the winning model, we took 4 values, equally
 1409 spread within the parameter range. We simulated behaviour using every combination ($4^7 = 16384$), fitted the
 1410 model and analysed how well the generative parameters (original values) correlated with the recovered ones (fitted
 1411 parameters). Pearson correlation coefficient = r. Each dot represents one simulation.



1412

1413 **Figure 5 - Figure supplement 1**

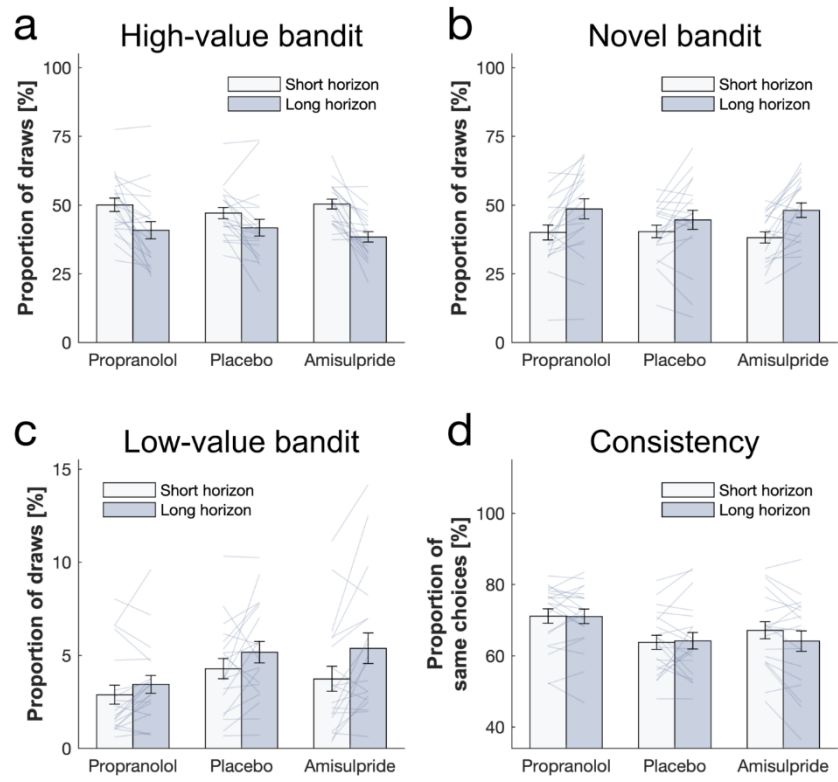
1414 Simulated behaviour for Thompson+ ϵ + η model. We used each subjects' fitted parameters to simulate behaviour
 1415 ($N_{trials}=4000$). Data are shown as mean \pm SEM and each dot/line represent one agent.



1416

1417 **Figure 5 – Figure supplement 2**

1418 Gender effect on prior variance parameter. Mean values (across horizon conditions) of σ_0 were larger for female
 1419 subjects, whereas in the amisulpride group, they were larger for male subjects. Data are shown as mean \pm SEM and
 1420 each dot represent one subject.
 1421



1422
 1423
 1424
 1425

Figure 5 - Figure supplement 3

Simulated behaviour for UCB+ ϵ + η model. We used each subjects' fitted parameters to simulate behaviour ($N_{trials}=4000$). Data are shown as mean \pm SEM and each dot/line represent one agent.