# Noradrenaline modulates tabula-rasa exploration — **Source link** ↗

Magda Dubois, Johanna Habicht, Jochen Michely, Rani Moran ...+2 more authors

Related papers:

- Dopamine regulates the exploration-exploitation trade-off in rats

1 **Human complex exploration strategies are extended via noradrenaline-modulated**
2 **heuristics**

3 Dubois M[1,2], Habicht J[1,2], Michely J[1,2], Moran R[1,2], Dolan RJ[1,2] & Hauser TU[1,2]

4 [1]Max Planck UCL Centre for Computational Psychiatry and Ageing Research, London WC1B
5 5EH, United Kingdom.
6 [2]Wellcome Centre for Human Neuroimaging, University College London, London WC1N 3BG,
7 United Kingdom.

8 **Corresponding author**

9 Tobias U. Hauser
10 Max Planck UCL Centre for Computational Psychiatry and Ageing Research
11 University College London
12 10-12 Russell Square
13 London WC1B 5EH
14 United Kingdom
15 Phone: +44 / 207 679 5264
16 Email: t.hauser@ucl.ac.uk

17 Number of pages: 67
18 Number of Figures: 5
19 Number of Tables: 0
20 Abstract: 123 words
21 Introduction: 1031 words
22 Discussion: 2393 words

23 **Data and materials availability:** Data and code will be provided upon acceptance.

**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.

**Introduction**

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

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

59   exploration is essentially a more sophisticated, uncertainty-driven, version of a softmax. By

60   accounting for stochasticity when comparing choice options' expected values, in effect choosing

61   based on both uncertainty and value, these exploration strategies increase the likelihood of

62   choosing 'good' options that are only slightly less valuable than the best (e.g. the Toblerone ice-

63   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 options.

77   Such 'value-free' random exploration, as opposed to the two previously considered 'value-based'

78   random explorations (for simulations comparing their effects cf. Figure 1 – Figure supplement 2)

79   that add stochasticity during choice value computation, forgoes any costly computation (i.e. value

80   mean and uncertainty), known as an $\epsilon$-greedy algorithmic strategy in reinforcement learning (*20*).

81  Computational efficiency, however, comes at the cost of sub-optimality due to occasional selection

82  of options of low expected value (e.g. the repulsive spinach ice cream).

83  Despite its sub-optimality, value-free random exploration has neurobiological plausibility.

84  Of relevance in this context is a view that exploration strategies depend on dissociable neural

85  mechanisms (*21*). Influences from noradrenaline and dopamine are plausible candidates in this

86  regard based on prior evidence (*9, 22*). Amongst other roles (such as memory (*23*), or energisation

87  of behaviour (*24, 25*)), the neuromodulator noradrenaline has been ascribed a function of indexing

88  uncertainty (*26–28*) or as acting as a 'reset button' that interrupts ongoing information processing

89  (*29–31*). Prior experimental work in rats shows boosting noradrenaline leads to more tabula-rasa-

90  like random behaviour (*32*), while pharmacological manipulations in monkeys indicates reducing

91  noradrenergic activity increases choice consistency (*33*).

92  In human pharmacological studies, interpreting the specific function of noradrenaline on

93  exploration strategies is problematic as many drugs, such as atomoxetine (e.g. (*34*)), impact

94  multiple neurotransmitter systems. Here, to avoid this issue, we chose the highly specific β-

95  adrenoceptor antagonist propranolol, which has only minimal impact on other neurotransmitter

96  systems (*35–37*). Using this neuromodulator, we examine whether signatures of value-free random

97  exploration are impacted by administration of propranolol.

98  An alternative computationally efficient exploration heuristic to random exploration is to

99  simply choose an option not encountered previously, which we term novelty exploration. Humans

100 often show novelty seeking (*38–41*), and this strategy can be used in exploration as implemented

101 by a low-cost version of the UCB algorithm. Here a novelty bonus (*42*) is added if a choice option

102 has not been seen previously (i.e. it does not have to rely on precise uncertainty estimates). The

103 neuromodulator dopamine is implicated not only in exploration in general (*43*), but also in

104 signalling such types of novelty bonuses, where evidence indicates a role in processing and

105 exploring novel and salient states (*39*, *44–47*). Although pharmacological dopaminergic studies in

106 humans have demonstrated effects on exploration as a whole (*48*), they have not identified specific

107 exploration strategies. Here, we used the highly specific D2/D3 antagonist, amisulpride, to

108 disentangle the specific role of dopamine and noradrenaline on different exploration strategies.

109       Thus, in the current study, we examine the contributions of value-free random exploration

110 and novelty exploration in human choice behaviour. We developed a novel exploration task

111 combined with computational modeling to probe the contributions of noradrenaline and dopamine.

112 Under double-blind, placebo-controlled, conditions we tested the impact of two antagonists with

113 a high affinity and specificity for either dopamine (amisulpride) or noradrenaline (propranolol).

114 Our results provide evidence that both exploration heuristics supplement computationally more

115 demanding exploration strategies, and that value-free random exploration is particularly sensitive
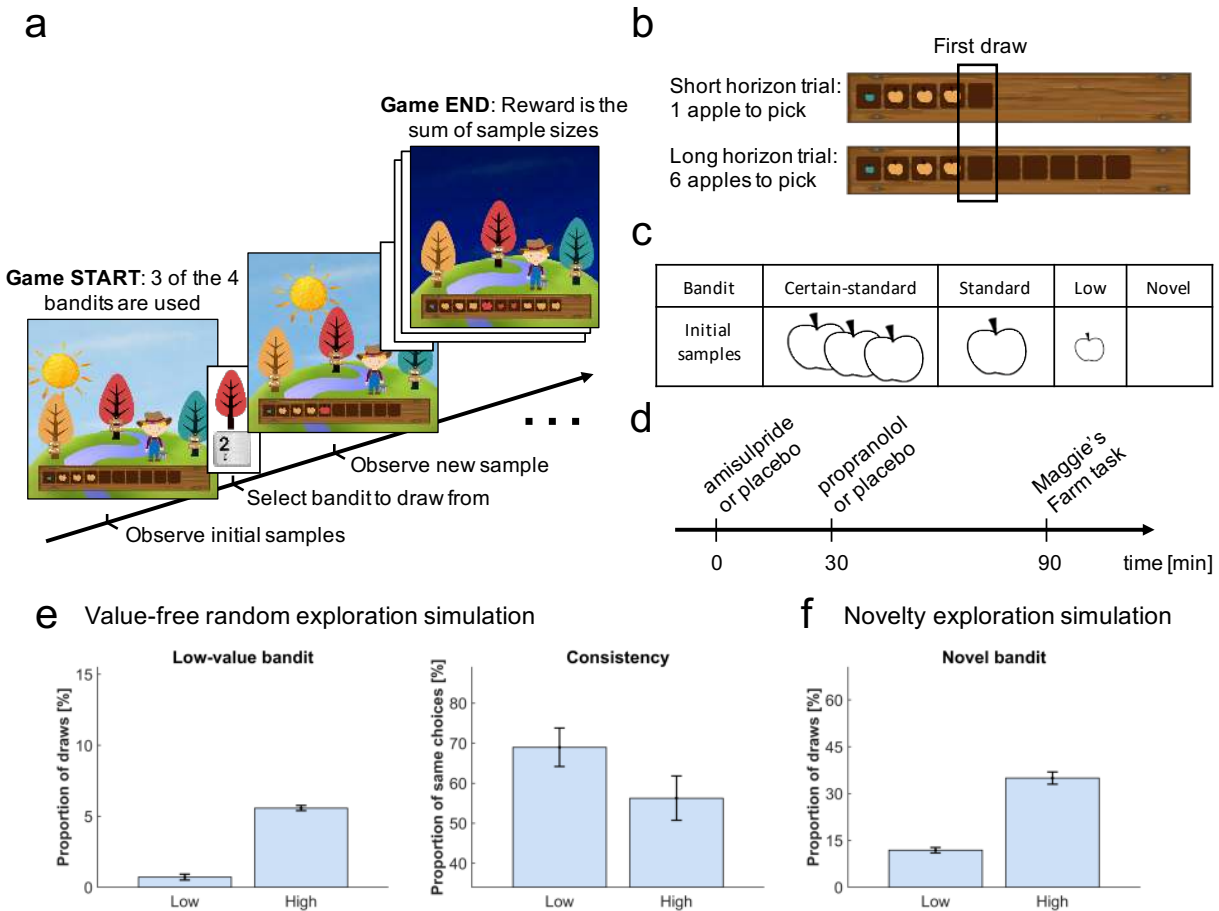
116 to noradrenergic modulation.

117 **Results**

118 *Probing the contributions of heuristic exploration strategies*

119 We developed a novel multi-round three-armed bandit task (Figure 1; bandits depicted as

120 trees), enabling us to assess the contributions of value-free random exploration and novelty

121 exploration in addition to Thompson sampling and UCB (combined with a softmax). In particular,

122 we exploited the fact that both heuristic strategies make specific predictions about choice patterns.

123 The novelty exploration assigns a 'novelty bonus' only to bandits for which subjects have no prior

124 information, but not to other bandits. This can be seen as a low-resolution version of UCB, which

125 assigns a bonus to all choice options proportionally to how informative they are, in effect a graded

126 bonus which scales to each bandits' uncertainty. Thus, to capture this heuristic, we manipulated

127 the amount of prior information with bandits carrying only little information (i.e. 1 vs 3 initial

128 samples) or no information (0 initial samples). A high novelty exploration predicts a higher

129 frequency of selecting the novel option (Figure 1f). This is in contrast to high exploration using

130 other strategies which does not predict such a strong effect on the novel option (cf. Figure 1 -

131 Figure supplement 5).

132 Value-free random exploration, captured here by $\epsilon$-greedy, predicts that all prior

133 information is discarded entirely and that there is equal probability attached to all choice options.

134 This strategy is distinct from other exploration strategies as it is likely to choose bandits known to

135 be substantially worse than the other bandits. Thus, a high value-free random exploration predicts

136 a higher frequency of selecting the low-value option (Figure 1e), whereas high exploration using

137 other strategies does not predict such effect (cf. Figure 1 - Figure supplement 3). A second

138 prediction is that choice consistency, across repeated trials, is substantially affected by value-free

139 random exploration. Given that value-free random exploration splits its choice probability equally

140 (i.e. 33.3% of choosing any bandit out of the three displayed), an increase in such exploration

141 predicts a lower likelihood of choosing the same bandit again, even under identical choice options

142 (Figure 1e). This contrasts to other strategies that make consistent exploration predictions (e.g.

143 UCB would consistently explore the choice option that carries a high information bonus; Figure 1

144 - Figure supplement 4).

145      We generated bandits from four different generative processes (Figure 1c) with distinct

146 sample means (but a fixed sampling variance) and number of initial samples (i.e. samples shown

147 at the beginning of a trial for this specific bandit). Subjects were exposed to these bandits before

148 making their first draw. The 'certain-standard bandit' and the (less certain) 'standard bandit' were

149 bandits with comparable means but varying levels of uncertainty, providing either three or one

150 initial samples (depicted as apples; similar to the horizon task (*7*)). The 'low-value bandit' was a

151 bandit with one initial sample from a substantially lower generative mean, thus appealing to a

152 value-free random exploration strategy alone. The last bandit, with a mean comparable with that

153 of the standard bandits, was a 'novel bandit' for which no initial sample was shown, primarily

154 appealing to a novelty exploration strategy (cf. Materials and Methods for a full description of

155 bandit generative processes). To assess choice consistency, all trials were repeated once. In the

156 pilot experiments (data not shown), we noted some exploration strategies tended to overshadow

157 other strategies. To effectively assess all exploration strategies, we opted to present only three of

158 the four different bandit types on each trial, as different bandit triples allow different explorations

159 to manifest. Lastly, to assess whether subjects' behaviour captured exploration, we manipulated

160 the degree to which subjects could interact with the same bandits. Similar to previous studies (*7*),

161 subjects could perform either one draw, encouraging exploitation (short horizon condition) or six

162 draws encouraging more substantial explorative behaviour (long horizon condition) (*7, 34*).

8

**Figure 1.** Study design. In the Maggie's farm task, subjects had to choose from three bandits (depicted as trees) to maximise an outcome (sum of reward). The rewards (apple size) of each bandit followed a normal distribution with a fixed sampling variance. (a) At the beginning of each trial, subjects were provided with some initial samples on the wooden crate at the bottom of the screen and had to select which bandit they wanted to sample from next. (b) Depending the condition, they could either perform one draw (short horizon) or six draws (long horizon). The empty spaces on the wooden crate (and the suns' position) indicated how many draws they had left. The first draw in both conditions was the main focus of the analysis. (c) In each trial, three bandits were displayed, selected from four possible bandits, with different generative processes that varied in terms of their sample mean and number of initial samples (i.e. samples shown at the beginning of a trial). The 'certain-standard bandit' and the 'standard bandit' had comparable means but different levels of uncertainty about their expected mean: they provided three and one initial sample respectively; the 'low-value bandit' had a low mean and displayed one initial sample; the 'novel bandit' did not show any initial sample and its mean was comparable with that of the standard bandits. (d) Prior to the task, subjects were administered different drugs: 400mg amisulpride that blocks dopaminergic D2/D3 receptors, 40mg propranolol to block noradrenergic β-receptors, and inert substances for the placebo group. Different administration times were chosen to comply with the different drug pharmacokinetics (placebo matching the other groups' administration schedule). (e) Simulating value-free random behaviour with a low vs high model parameter ($\epsilon$) in this task shows that in a high regime, agents choose the low-value bandit more often (left panel; mean ± SD) and are less consistent in their choices when facing identical choice

9

185 options (right panel). (f) Novelty exploration exclusively promotes choosing choice options for
186 which subjects have no prior information, captured by the 'novel bandit' in our task. For details
187 about simulations cf. Materials and Methods. For details about the task display cf. Figure 1 –
188 Figure supplement 1. For simulations of different exploration strategies and their impact of
189 different bandits cf. Figure 1 – Figure supplement 2-5.

190

191 *Testing the role of catecholamines noradrenaline and dopamine*

192 In a double-blind, placebo-controlled, between-subjects, study design we assigned subjects

193 (N=60) randomly to one of three experimental groups: amisulpride, propranolol or placebo. The

194 first group received 40mg of the $\beta$-adrenoceptor antagonist propranolol to alter noradrenaline

195 function, while the second group was administered 400mg of the D2/D3 antagonist amisulpride

196 that alters dopamine function. Because of different pharmacokinetic properties, these drugs were

197 administered at different times (Figure 1d) and compared to a placebo group that received a

198 placebo at both drug times to match the corresponding antagonists' time. One subject (amisulpride

199 group) was excluded from the analysis due to a lack of engagement with the task. Reported

200 findings were corrected for IQ and mood, as drug groups differed marginally in those measures

201 (cf. Appendix 2 Table 1), by adding WASI (*49*) and PANAS (*50*) negative scores as covariates in

202 each ANOVA. Similar results were obtained in an analysis that corrected for physiological effects

203 as from the analysis without covariates (cf. Appendix 1).

204 *Increased exploration when information can subsequently be exploited*

205 Our task embodied two decision-horizon conditions, a short and a long. To assess whether

206 subjects explored more in a long horizon condition, in which additional information can inform

207 later choices, we examined which bandit subjects chose in their first draw (in accordance with the

208 horizon task (*7*)), irrespective of their drug group. A marker of exploration here is evident if

209 subjects chose bandits with lower expected values, computed as the mean value of their initial

210    samples shown (trials where the novel bandit was chosen were excluded). As expected, subjects

211    chose bandits with a lower expected value in the long compared to the short horizon (repeated-

212    measures ANOVA for the expected value: $F(1, 56)=19.457$, $p<.001$, $\eta2=.258$; Figure 2a). To

213    confirm that this was a consequence of increased exploration, we analysed the proportion of how

214    often the high-value option was chosen (i.e. the bandit with the highest expected reward based on

215    its initial samples) and we found that subjects (especially those with higher IQ) sampled from it

216    more in the short compared to the long horizon, (WASI-by-horizon interaction: $F(1,54)=13.304$,

217    $p=.001$, $\eta2=.198$; horizon main effect: $F(1, 54)=3.909$, $p=.053$, $\eta2=.068$; Figure 3a), confirming a

218    reduction in exploitation when this information could be subsequently used. Interestingly, this

219    frequency seemed to be marginally higher in the amisulpride group, suggesting an overall higher

220    tendency to exploitation following dopamine blockade (cf. Appendix 1). This horizon-specific

221    behaviour resulted in a lower reward on the $1^{st}$ sample in the long compared to the short horizon

222    $(F(1, 56)=23.922$, $p<.001$, $\eta2=.299$; Figure 2c). When we tested whether subjects were more likely

223    to choose options they knew less about (computed as the mean number of initial samples shown),

224    we found that subjects chose less known (i.e. more informative) bandits more often in the long

225    horizon compared to the short horizon $(F(1, 56)=58.78$, $p<.001$, $\eta2=.512$; Figure 2b).

226        Next, to evaluate whether subjects used the additional information beneficially in the long

227    horizon condition, we compared the average reward (across six draws) obtained in the long

228    compared to short horizon (one draw). We found that the average reward was higher in the long

229    horizon $(F(1, 56)=103.759$, $p<.001$, $\eta2=.649$; Figure 2c), indicating that subjects tended to choose

230    less optimal bandits at first but subsequently learnt to appropriately exploit the harvested

231    information to guide choices of better bandits in the long run. Additionally, when looking

232    specifically at the long horizon condition, we found that subjects earned more when their first draw

233    was explorative versus exploitative (Figure 2 - Figure supplement 1c-d; cf. Appendix 2 for details).

234



235
236    **Figure 2.** Benefits of exploration. To investigate the effect of information on performance we
237    collapsed subjects over all three treatment groups. (a) The expected value (average of its initial
238    samples) of the first chosen bandit as a function of horizon. Subjects chose bandits with a lower
239    expected value (i.e. they explored more) in the long horizon compared to the short horizon. (b)
240    The mean number of samples for the first chosen bandit as a function of horizon. Subjects chose
241    less known (i.e. more informative) bandits more in the long compared to the short horizon. (c) The
242    first draw in the long horizon led to a lower reward than the first draw in the short horizon,
243    indicating that subjects sacrificed larger initial outcomes for the benefit of more information. This
244    additional information helped making better decisions in the long run, leading to a higher earning
245    over all draws in the long horizon. For values and statistics cf. Appendix 2 Table 3. For response
246    times and details about all long horizons' samples cf. Figure 2 – Figure supplement 1. *** =p<.001.
247    Data are shown as mean ± SEM and each dot/line represent a subject.

248

249        *Subjects demonstrate value-free random behaviour*

250        Value-free random exploration (analogue to $\epsilon$-greedy) predicts that $\epsilon$ % of the time each

251    option will have an equal probability of being chosen. In such a regime (compared to more

252    complex strategies that would favour options with a higher expected value with a similar

253    uncertainty), the probability of choosing bandits with a low expected value (here the low-value

254    bandit; Fig. 1e) will be higher (cf. Figure 1 – Figure supplement 3). We investigated whether the

255    frequency of picking the low-value bandit was increased in the long horizon condition across all

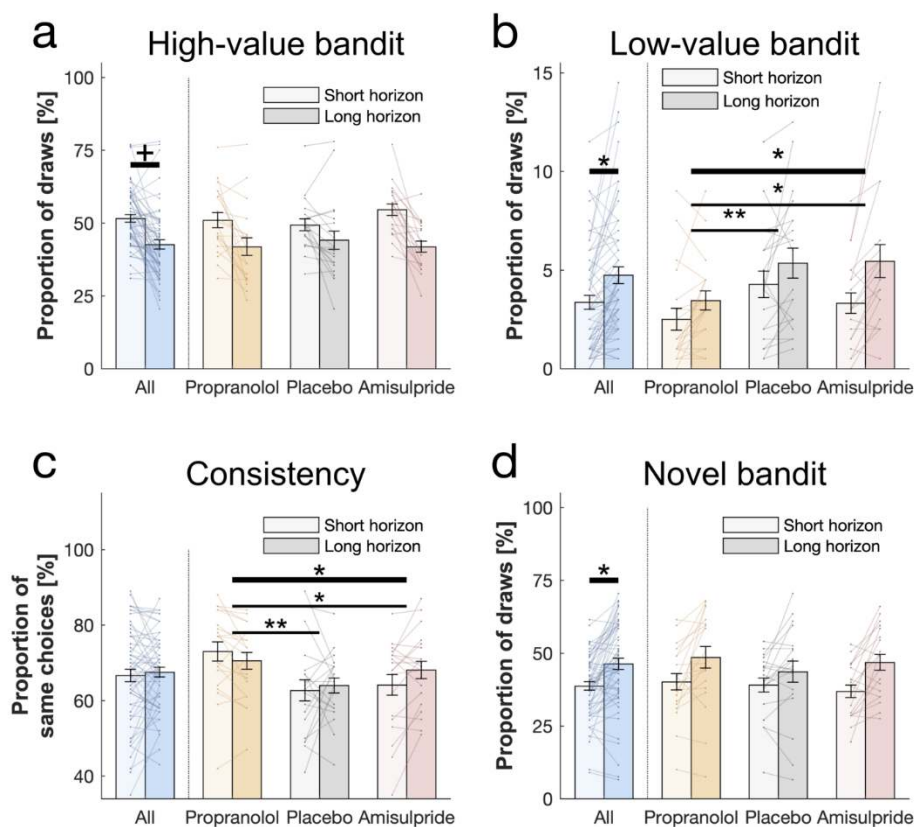256    subjects (i.e. when exploration is useful), and we found a significant main effect of horizon (F(1,

257 54)=4.069, p=.049, η2=.07; Figure 3b). This demonstrates that value-free random exploration is

258 utilised more when exploration is beneficial.

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

260 When we tested whether value-free random exploration was sensitive to neuromodulatory

261 influences, we found a difference in how often drug groups sampled from the low-value option

262 (drug main effect: F(2, 54)=7.003, p=.002, η2=.206; drug-by-horizon interaction: F(2, 54)=2.154,

263 p=.126, η2=.074; Figure 3b). This was driven by the propranolol group choosing the low-value

264 option significantly less often than the other two groups (placebo vs propranolol: t(40)=2.923,

265 p=.005, d=.654; amisulpride vs propranolol: t(38)=2.171, p=.034, d=.496) with no difference

266 between amisulpride and placebo: (t(38)=-0.587, p=.559, d=.133). These findings demonstrate that

267 a key feature of value-free random exploration, the frequency of choosing low-value bandits, is

268 sensitive to influences from noradrenaline.

269 To further examine drug effects on value-free random exploration, we assessed a second

270 prediction, namely choice consistency. Because value-free random exploration ignores all prior

271 information and chooses randomly, it should result in a decreased choice consistency when

272 presented identical choice options (cf. Figure 1 – Figure supplement 2 & 4, compared to more

273 complex strategies which are always biased towards the rewarding or the information providing

274 bandit for example). To this end, each trial was duplicated in our task, allowing us to compute the

275 consistency as the percentage of time subjects sampled from an identical bandit when facing the

276 exact same choice options. In line with the above analysis, we found a difference in consistency

277 by which drug groups sampled from different option (drug main effect: F(2, 54)=7.154, p=.002,

278 η2=.209; horizon main effect: F(1, 54)=1.333, p=.253, η2=.024; drug-by-horizon interaction: F(2,

279 54)=3.352, p=.042, η2=.11; Figure 3c), driven by the fact that the propranolol group chose

280     significantly more consistently than the other two groups (pairwise comparisons: placebo vs

281     propranolol: $t(40)=-3.525$, $p=.001$, $d=.788$; amisulpride vs placebo: $t(38)=1.107$, $p=.272$, $d=.251$;

282     amisulpride vs propranolol: $t(38)=-2.267$, $p=.026$, $d=.514$). Please see Appendix 1 for further

283     discussion and analysis of the drug-by-horizon interaction. Taken together, these results indicate

284     that value-free random exploration depends critically on noradrenaline functioning, such that an

285     attenuation of noradrenaline leads to a reduction in value-free random exploration.



286

**Figure 3.** Behavioural horizon and drug effects. Choice patterns in the first draw for each horizon and drug group (propranolol, placebo and amisulpride). (a) Subjects sampled from the high-value bandit (i.e. bandit with the highest average reward of initial samples) more in the short horizon compared to the long horizon indicating reduced exploitation. (b) Subjects sampled from the low-value bandit more in the long horizon compared to the short horizon indicating value-free random exploration, but subjects in the propranolol group sampled less from it overall, and (c) were more consistent in their choices overall, indicating that noradrenaline blockade reduces value-free random exploration. (d) Subjects sampled from the novel bandit more in the long horizon compared to the short horizon indicating novelty exploration. Please note that some horizon effects were modulated by subjects' intellectual abilities when additionally controlling for them (cf.

297      Appendix 2 Table 4). Horizontal bars represent rm-ANOVA (thick) and pairwise comparisons
298      (thin). † =p<.07, * =p<.05, ** =p<.01. Data are shown as mean ± SEM and each line represent one
299      subject. For values and statistics cf. Appendix 2 Table 4. For response times and frequencies
300      specific to the displayed bandits cf. Figure 3 – Figure supplement 1-2.

301

302      *Novelty exploration is unaffected by catecholaminergic drugs*

303      Next, we examined whether subjects show evidence for novelty exploration by choosing the

304      novel bandit for which there was no prior information (i.e. no initial samples), as predicted by

305      model simulations (Figure 1f). We found a significant main effect of horizon ($F_{(1, 54)}$=5.593,

306      p=.022, $\eta2$=.094; WASI-by-horizon interaction: $F_{(1, 54)}$ =13.897, p<.001, $\eta2$=.205; Figure 3d)

307      indicating that subjects explored the novel bandit significantly more often in the long horizon

308      condition, and this was particularly strong for subjects with a higher IQ. We next assessed whether

309      novelty exploration was sensitive to our drug manipulation, but found no drug effects on the novel

310      bandit ($F_{(2, 54)}$=1.498, p=.233, $\eta2$=.053; drug-by-horizon interaction: $F_{(2, 54)}$=.542, p=.584,

311      $\eta2$=.02; Figure 3d). Thus, there was no evidence that an attenuation of dopamine or noradrenaline

312      function impact novelty exploration in this task.
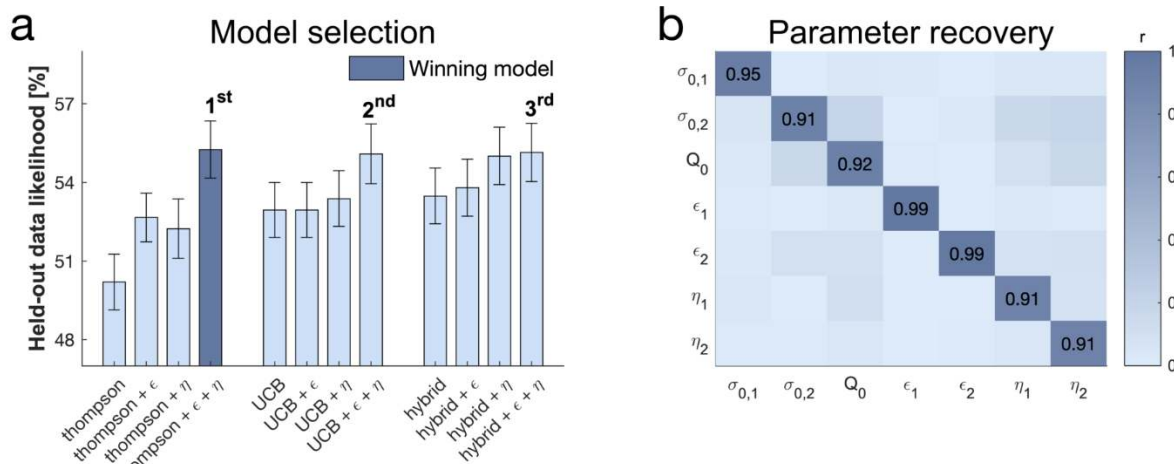
313      *Subjects combine computationally demanding strategies and exploration heuristics*

314      To examine the contributions of different exploration strategies to choice behaviour, we

315      fitted a set of computational models to subjects' behaviour, building on models developed in

316      previous studies (*1*). In particular, we compared models incorporating UCB, Thompson sampling,

317      an $\epsilon$-greedy algorithm and the novelty bonus (cf. Materials and Methods). Essentially, each model

318      makes different exploration predictions. In the Thompson model, Thompson sampling (*8, 51*) leads

319      to an uncertainty-driven value-based random exploration, where both expected value and

320      uncertainty contribute to choice. In this model higher uncertainty leads to more exploration such

321      that instead of selecting a bandit with the highest mean, bandits are chosen relative to how often a

322    random sample would yield the highest outcome, thus accounting for uncertainty (*2*). The UCB

323    model (*3, 4*), capturing directed exploration, predicts that each bandit is chosen according to a

324    mixture of expected value and an additional expected information gain (*2*). This is realised by

325    adding a bonus to the expected value of each option, proportional to how informative it would be

326    to select this option (i.e. the higher the uncertainty in the options' value, the higher the information

327    gain). This computation is then passed through a softmax decision model, capturing value-based

328    random exploration. Novelty exploration is a simplified version of the information bonus in the

329    UCB algorithm, which only applies to entirely novel options. It defines the intrinsic value of

330    selecting a bandit about which nothing is known, and thus saves demanding computations of

331    uncertainty for each bandit. Lastly, the value-free random $\epsilon$-greedy algorithm selects any bandit $\epsilon$

332    % of the time, irrespective of the prior information of this bandit. For additional models cf.

333    Appendix 1.

334         We used cross-validation for model selection (Figure 4a) by comparing the likelihood of

335    held-out data across different models, an approach that adequately arbitrates between model

336    accuracy and complexity. The winning model encompasses uncertainty-driven value-based

337    random exploration (Thompson sampling) with value-free random exploration ($\epsilon$-greedy

338    parameter) and novelty exploration (novelty bonus parameter $\eta$). The winning model predicted

339    held-out data with a 55.25% accuracy (SD=8.36%; chance level =33.33%). Similarly to previous

340    studies (*1*), the hybrid model combining UCB and Thompson sampling explained the data better

341    than each of those processes alone, but this was no longer the case when accounting for novelty

342    and value-free random exploration (Figure 4a). The winning model further revealed that all

343    parameter estimates could be accurately recovered (Figure 4b; Figure 4 – Figure supplement 3).

344    Interestingly, although the 2[nd] and 3[rd] place models made different prediction about the complex

345     exploration strategy, using a directed exploration with value-based random exploration (UCB) or

346     a combination of complex strategies (hybrid) respectively, they share the characteristic of

347     benefitting from value-free random and novelty exploration. This highlights that subjects used a

348     mixture of computationally demanding and heuristic exploration strategies.



349

350     **Figure 4.** Subjects use a mixture of exploration strategies. (a) A 10-fold cross-validation of the
351     likelihood of held-out data was used for model selection (chance level =33.3%; for model selection
352     at the individual level cf. Figure 4 – Figure supplement 1). The Thompson model with both the $\epsilon$-
353     greedy parameter and the novelty bonus $\eta$ best predicted held-out data (b) Model simulation with
354     $4^7$ simulations predicted good recoverability of model parameters (for correlations between
355     behaviour and model parameters cf. Figure 4 – Figure supplement 2); $\sigma_0$ is the prior variance and
356     $Q_0$ is the prior mean (for parameter recovery correlation plots cf. Figure 4 – Figure supplement 3).
357     1 stands for short horizon-, and 2 for long horizon-specific parameters. For values and parameter
358     details cf. Appendix 2 Table 5.

359

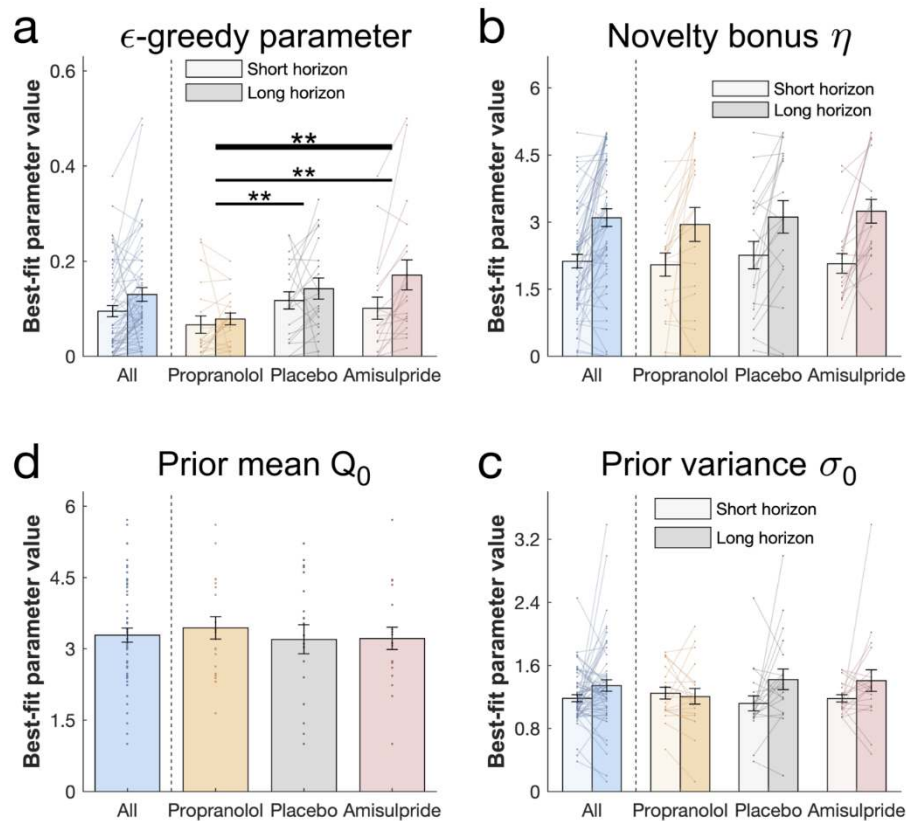360     *Noradrenaline controls value-free random exploration*

361        To more formally compare the impact of catecholaminergic drugs on different exploration

362     strategies, we assessed the free parameters of the winning model between drug groups (Figure 5,

363     cf. Appendix 2 Table 6 for exact values). First, we examined the $\epsilon$-greedy parameter that captures

364     the contribution of value-free random exploration to choice behaviour. We assessed how this

365     value-free random exploration differed between drug groups. A significant drug main effect (drug

366     main effect: $F_{(2, 54)}=6.722$, $p=.002$, $\eta2=.199$; drug-by-horizon interaction: $F_{(2, 54)}=1.305$, $p=.28$,

17

367  η2=.046; Figure 5a) demonstrates that the drug groups differ in how strongly they deploy this

368  exploration strategy. Post-hoc analysis revealed that subjects with reduced noradrenaline

369  functioning had the lowest values of $\epsilon$ (pairwise comparisons: placebo vs propranolol:

370  t(40)=3.177, p=.002, d=.71; amisulpride vs propranolol: t(38)=2.723, p=.009, d=.626) with no

371  significant difference between amisulpride vs placebo: (t(38)=.251, p=.802, d=.057). Critically,

372  the effect on $\epsilon$ was also significant when the complex exploration strategy was a directed

373  exploration with value-based random exploration (2nd place model) and, marginally significant,

374  when it was a combination of the above (3rd place model; cf. Appendix 1).

375      The $\epsilon$-greedy parameter was also closely linked to the above behavioural metrics (correlation

376  between the $\epsilon$-greedy parameter with draws from the low-value bandit: $R_{Pearson}$=.828, p<.001;

377  and with choice consistency: $R_{Pearson}$=-.596, p<.001; Figure 4 – Figure supplement 2), and

378  showed a similar horizon effect (horizon main effect: F(1, 54)=1.968, p=.166, η2=.035; WASI-

379  by-horizon interaction: F(1, 54)=6.08, p=.017, η2=.101; Figure 5a). Our findings thus accord with

380  the model-free analyses and demonstrate that noradrenaline blockade reduces value-free random

381  exploration.

382

383



384

**Figure 5.** Drug effects on model parameters. The winning model's parameters were fitted to each subject's first draw (for model simulations cf. Figure 5 – Figure supplement 1). (a) Subjects had higher values of $\epsilon$ (value-free random exploration) in the long compared to the short horizon. Notably, subjects in the propranolol group had lower values of $\epsilon$ overall, indicating that attenuation of noradrenaline functioning reduces value-free random exploration. Subjects from all groups (b) assigned a similar value to novelty, captured by the novelty bonus η, which was higher (more novelty exploration) in the long compared to the short horizon. (c) The groups had similar beliefs $Q_0$ about a bandits' mean before seeing any initial samples and (d) were similarly uncertain $\sigma_0$ about it (for gender effects cf. Figure 5 – Figure supplement 2). Please note that some horizon effects were modulated by subjects' intellectual abilities when additionally controlling for them (cf. Appendix 2 Table 6). ** =p<.01. Data are shown as mean ± SEM and each dot/line represent one subject. For parameter values and statistics cf. Appendix 2 Table 6.

397

398          *No drug effects on other parameters*

399          The novelty bonus $\eta$ captures the intrinsic reward of selecting a novel option. In line with the

400   model-free behavioural findings, there was no difference between drug groups in terms of this

19

401    effect (F(2, 54)=.249, p=.78, $\eta^2$=.009; drug-by-horizon interaction: F(2, 54)=.03, p=.971,

402    $\eta^2$=.001). There was also a close alignment between model-based and model-agnostic analyses

403    (correlation between the novelty bonus $\eta$ with draws from the novel bandit: $R_{Pearson}$=.683,

404    p<.001; Figure 4 – Figure supplement 2), and we found a similarly increased novelty bonus effect

405    in the long horizon in subjects with a higher IQ (WASI-by-horizon interaction: F(1, 54) =8.416,

406    p=.005, $\eta^2$=.135; horizon main effect: F(1, 54)=1.839, p=.181, $\eta^2$=.033; Figure 5b).

407        When analysing the additional model parameter, we found that subjects had similar prior

408    beliefs about bandits, given by the initial estimate of a bandit's mean (prior mean $Q_0$: F(2,

409    54)=.118, p=.889, $\eta^2$=.004; Figure 5c) and their uncertainty about it (prior variance $\sigma_0$: horizon

410    main effect: F(1, 54)=.129, p=.721, $\eta^2$=.002; drug main effect: F(2, 54)=.06, p=.942, $\eta^2$=.002;

411    drug-by-horizon interaction: F(2, 54)=2.162, p=.125, $\eta^2$=.074; WASI-by-horizon interaction: F(1,

412    54)=.022, p=.882, $\eta^2$<.001; Figure 5d). Interestingly, our dopamine manipulation seemed to affect

413    this uncertainty in a gender-specific manner, with female subjects having larger values of $\sigma_0$

414    compared to males in the placebo group, and with the opposite being true in the amisulpride group

415    (cf. Appendix 1). Taken together, these findings show that value-free random exploration was most

416    sensitive to our drug manipulations.

**Discussion**

417

418    Solving the exploration-exploitation problem is non trivial, and one suggestion is that

419    humans solve it using computationally demanding exploration strategies (1, 2), taking account of

420    the uncertainty (variance) as well as the expected reward (mean) of each choice. Although tracking

421    the distribution of summary statistics (e.g. mean and variance) is less resource costly than keeping

422    track of full distributions (*52*), it nevertheless carries considerable costs when one has to keep track

423    of multiple options, as in exploration. Indeed, in a three-bandit task such as that considered here,

424    this results in a necessity to compute 6 key-statistics, drastically limiting computational resources

425    when selecting among choice options (10). Real-life decisions often comprise an unlimited range

426    of options, which results in a tracking of a multitude of key-statistics, potentially mandating a

427    deployment of alternative more efficient strategies. Here, we demonstrate that two additional, less

428    resource-hungry heuristics are at play during human decision-making, value-free random

429    exploration and novelty exploration.

430    By assigning intrinsic value (novelty bonus (*42*)) to an option not encountered before (*53*),

431    a novelty bonus can be seen as an efficient simplification of demanding algorithms, such as UCB

432    (*3*, *4*). It is interesting to note that our winning model did not include UCB, but instead novelty

433    exploration. This indicates humans might use such a novelty shortcut to explore unseen, or rarely

434    visited, states to conserve computational costs when such a strategy is possible. A second

435    exploration heuristic that also requires minimal computational resources, value-free random

436    exploration, also plays a role in our task. Even though less optimal, its simplicity and neural

437    plausibility renders it a viable strategy. We show through converging behavioural and modelling

438    measures that both value-free random and novelty exploration were deployed in a goal-directed

439    manner, coupled with increased levels of exploration when this was strategically useful.

21

440  Importantly, these heuristics were observed in all best models (1st, 2nd and 3rd position) even though

441  each incorporated different exploration strategies. This suggests that the complex models made

442  similar predictions in our task, and demonstrates that value-free random exploration is at play even

443  when accounting for other value-based forms of random exploration (*1*, *7*), whether fixed or

444  uncertainty-driven.

445  Exploration was captured in a similar manner to previous studies (*7*), by comparing in the

446  same setting (i.e. same prior information) the first choice in a long decision horizon, where reward

447  can be increased in the long term through information gain, and in a short decision horizon where

448  information cannot subsequently be put to use. This means that by changing the opportunity to

449  benefit from the information gained for the first sample, the long horizon invites extended

450  exploration (7), what we find also in our study. This experimental manipulation is a well-

451  established means for altering exploration and has been used extensively in previous studies (*7*,

452  *21*, *34*, *54*). Nevertheless, there remains a possibility that a longer horizon may also affect the

453  psychological nature of the task. In our task, reward outcomes were presented immediately after

454  every draw, rendering it unlikely that perception of reward delays (i.e. delay discounting) is

455  impacted. Moreover, a monetary bonus was given only at the end of the task, and thus did not

456  impact a horizon manipulation. We also consider our manipulation was unlikely to change effort

457  in each horizon, because the reward (i.e. size of the apple) remains the same at every draw,

458  resulting in an equivalent reward-effort ratio (*55–58*). However, this issue can be addressed in

459  further studies, for example, by equating the amount of button presses across both conditions.

460  Value-free random exploration might reflect other influences, such as attentional lapses or

461  impulsive motor responses. We consider these as unlikely to a significant factor at play here.

462  Indeed, there are two key features that would signify such effects. Firstly, these influences would

463     be independent of task condition. Secondly, they would be expected to lead to shorter, or more

464     variable, response latencies. In our data, we observe an increase in value-free exploration in the

465     long horizon condition in both behavioural measures and model parameters, speaking against an

466     explanation based upon simple mistakes. Moreover, we did not observe a difference in response

467     latency for choices that were related to value-free random exploration (cf. Appendix 1), further

468     arguing against mistakes. Lastly, the sensitivity of value-free random exploration to propranolol

469     supports this being a separate process, and previous studies using the same drug did not find an

470     effect on task mistakes (e.g. on accuracy (*59*); (*33*, *58–60*)). However, future studies could explore

471     these exploration strategies in more detail including by reference to subjects' own self-reports.

472        It is still unclear how exploration strategies are implemented neurobiologically.

473     Noradrenaline inputs, arising from the locus coeruleus (*63*) (LC) are thought to modulate

474     exploration (*2*, *64*, *65*), though empirical data on its precise mechanisms and means of action

475     remains limited. In this study, we found that noradrenaline impacted value-free random

476     exploration, in contrast to novelty exploration and complex exploration. This might suggest that

477     noradrenaline influences ongoing valuation or choice processes that discards prior information.

478     Importantly, this effect was observed whether the complex exploration was an uncertainty-driven

479     value-based random exploration (winning model), a directed exploration with value-based random

480     exploration (2nd place model) or a combination of the above (3rd place model; cf. Appendix 1).

481     This is consistent with findings in rodents where enhanced anterior cingulate noradrenaline release

482     leads to more random behaviour (*32*). It is also consistent with pharmacological findings in

483     monkeys that show enhanced choice consistency after reducing LC noradrenaline firing rates (*33*).

484     It would be interesting for future studies to determine, in more detail, whether value-free random

485    exploration is corrupting a value computation itself, or whether it exclusively biases the choice

486    process.

487        We note that pupil diameter has been used as an indirect marker of noradrenaline activity

488    (*66*), although the link between the two it not always straightforward (*36*). Because the effect of

489    pharmacologically induced changes of noradrenaline levels on pupil size remains poorly

490    understood (*36*, *67*), including the fact that previous studies found no effect of propranolol on pupil

491    diameter (*36*, *68*), we opted against using pupillometry in this study. However, our current findings

492    align with previous human studies that show an association between this indirect marker and

493    exploration, but that study did not dissociate between the different potential exploration strategies

494    that subjects could deploy (*69*). Future studies might usefully include indirect measures of

495    noradrenaline activity, for example pupillometry, to examine a potential link between natural

496    variations in noradrenaline levels and a propensity towards value-free random exploration.

497        The LC has two known modes of synaptic signalling (*63*), tonic and phasic, thought to have

498    complementary roles (*31*). Phasic noradrenaline is thought to act as a reset button (*31*), rendering

499    an agent agnostic to all previously accumulated information, a de facto signature of value-free

500    random exploration. Tonic noradrenaline has been associated, although not consistently (*70*), with

501    increased exploration (*64*, *71*), decision noise in rats (*72*) and more specifically with random as

502    opposed to directed exploration strategies (*34*). This later study unexpectedly found that boosting

503    noradrenaline decreased (rather than increased) random exploration, which the authors speculated

504    was due to an interplay with phasic signalling. Importantly, the drug used in that study also affects

505    dopamine function making it difficult to assign a precise interpretation to the finding. A

506    consideration of this study influenced our decision to opt for drugs with high specificity for either

507    dopamine or noradrenaline (*59*), enabling us to reveal highly specific effects on value-free random

508    exploration. Although the contributions of tonic and phasic noradrenaline signalling cannot be

509    disentangled in our study, our findings align with theoretical accounts and non-primate animal

510    findings, indicating that phasic noradrenaline promotes value-free random exploration.

511         Aside from this 'reset signal' role, noradrenaline has been assigned other roles, including

512    a role in memory function (*23*, *73*, *74*). To minimise a possible memory-related impact, we

513    designed the task such that all necessary information was visible on the screen at all times. This

514    means subjects did not have to memorise values for a given trial, rendering the task less susceptible

515    to forgetting or other memory effects. Another role for noradrenaline relates to volatility and

516    uncertainty estimation (*26–28*), as well as the energisation of behaviour (*24*, *25*). Non-human

517    primates studies demonstrate a higher LC activation for high effort choices, suggesting that

518    noradrenaline release facilitates energy mobilisation (*24*). Theoretical models also suggest that the

519    LC is involved in the control of effort exertion. Thus, it is thought to contribute to trading off

520    between effortful actions leading to large rewards and "effortless" actions leading to small rewards

521    by modulating "raw" reward values as a function of the required effort (*25*). Our task can be

522    interpreted as encapsulating such a trade-off: complex exploration strategies are effortful but

523    optimal in terms of reward gain, while value-free random exploration requires little effort while

524    occasionally leading to low reward. Applying this model, a noradrenaline boost could optimise

525    cognitive effort allocation for high reward gain (*25*), thereby facilitating complex exploration

526    strategies compared to value-free random exploration. In such a framework, blocking

527    noradrenaline release should decrease usage of complex exploration strategies, leading to an

528    increase of value-free random exploration which is the opposite of what we observed in our data.

529    Another interpretation of an effort-facilitation model of noradrenaline is that a boost would help

530    overcoming cost, i.e. the lack of immediate reward when selecting the low-value bandit, essentially

531     providing a significant increase to the value of information gain. In line with our results, a decrease

532     would interrupt this boost in valuation, removing an incentive to choose the low-value option.

533     However, this theory is currently limited by the absence of empirical evidence for noradrenaline

534     boosting valuation.

535     Noradrenaline blockade by propranolol has been shown previously to enhance

536     metacognition (*75*), decrease information gathering (*59*), and attenuate arousal-induced boosts in

537     incidental memory (*36*). All of these findings, including a decrease in value-free random

538     exploration found here, suggests propranolol may influence how neural noise affects information

539     processing. In particular, the results indicate that under propranolol behaviour is more

540     deterministic and less influenced by 'task-irrelevant' distractions. This aligns with theoretical

541     ideas, as well as recent optogenetic evidence (*32*), that propose noradrenaline infuses noise in a

542     temporally targeted way (*31*). It also accords with studies implicating noradrenaline in attention

543     shifts (for a review cf. (*76*)). Other theories of noradrenaline/catecholamine function can link to

544     determinism (*64, 65*), although the hypothesized direction of effect is different (i.e. noradrenaline

545     increases determinism). This idea can be extended also to tasks where propranolol has been shown

546     to attenuate a discrimination between different levels of loss (with no effect on the value-based

547     exploration parameter, referred to in these studies as consistency) (*62*) and a reduction in loss

548     aversion (*60*). This hints at additional roles for noradrenaline on prior information and task-

549     distractibility during exploration in loss frame environments. Future studies investigating

550     exploration in loss contexts might provide important additional information on these questions.

551     It is important to mention here that β-adrenergic receptors, the primary target of

552     propranolol, have been shown (unlike $\alpha$-adrenergic receptors) to increase synaptic inhibition

553     within rat cortex (*77*), specifically through inhibitory GABA-mediated transmission (*78*).

554 Additionally *β*-adrenergic receptors are more concentrated in the intermediate layers in the

555 prefrontal area (*79*), within which inhibition is favoured (*80*). Thus inhibitory mechanisms might

556 account for noradrenaline-related task-distractibility and randomness, or the role of β-adrenergic

557 receptors in executive function impairments (*81*). This raises the question of whether blocking β-

558 adrenergic receptors might lead to an accumulation of synaptic noradrenaline, and therefore act

559 via α-adrenergic receptors. To the best of our knowledge, evidence for such an effect is limited. A

560 second question is whether the observed effects are a pure consequence of propranolol's impact

561 on the brain, or whether they reflect peripheral effects of propranolol. When we examined

562 peripheral markers (i.e. heart rate) and behaviour we found no evidence for an effect on any of our

563 findings, rendering such influences unlikely. However, future studies using drugs that exclusively

564 targets peripheral, but not central, noradrenaline receptors (e.g. (*82*)) are needed to answer this

565 question conclusively.

566  Dopamine has been ascribed multiple functions besides reward learning (*83*), such as

567 novelty seeking (*46*, *84*, *85*) or exploration in general (*43*). In fact, studies have demonstrated that

568 there are different types of dopaminergic neurons in the ventral tegmental area, and that some

569 contribute to non-reward signals, such as saliency and novelty (*44*). This suggests a role in novelty

570 exploration. Moreover, dopamine has been suggested as important in an exploration-exploitation

571 arbitration (*21*, *86*, *87*), although its precise role remains unclear, given reported effects on random

572 exploration (*88*), on directed exploration (*45*, *89*), or no effects at all (*90*). A recent study found

573 no effect following dopamine blockade using haloperidol (*87*), which interestingly also affects

574 noradrenaline function (e.g. (*91*, *92*)). Our results did not demonstrate any main effect of dopamine

575 manipulation on exploration strategies, even though blocking dopamine was associated with a

576 trend level increase in exploitation (cf. Appendix 1). We believe it unlikely this reflects an

577 ineffective drug dose as previous studies have found neurocognitive effects with the same dose

578 (*36, 59, 93, 94*).

579       One possible reason for an absence of significant findings is that our dopaminergic

580 blockade targets D2/D3 receptors rather than D1 receptors, a limitation due a lack of available

581 specific D1 receptor blockers for use in humans. An expectation of greater D1 involvement arises

582 out of theoretical models (*95*) and a prefrontal hypothesis of exploration (*89*). Interestingly, we

583 observed a weak gender-specific differential drug effect on subjects' uncertainty about an expected

584 reward, with women being more uncertain than men in the placebo setting, but more certain in the

585 dopamine blockade setting (cf. Appendix 1). This might be meaningful as other studies using the

586 same drug have also found behavioural gender-specific drug effects (*96*). Upcoming, novel drugs

587 (*97*) might be able help unravel a D1 contribution to different forms of exploration. Additionally,

588 future studies could use approved D2/D3 agonists (e.g. ropinirole) in a similar design to probe

589 further whether enhancing dopamine leads to a general increase in exploration.

590       In conclusion, humans supplement computationally expensive exploration strategies with

591 less resource demanding exploration heuristics, and as shown here the latter include value-free

592 random and novelty exploration. Our finding that noradrenaline specifically influences value-free

593 random exploration demonstrates that distinct exploration strategies may be under specific

594 neuromodulator influence. Our current findings may also be relevant to enabling a richer

595 understanding of disorders of exploration, such as attention-deficit/hyperactivity disorder (*22, 98*)

596 including how aberrant catecholamine function might contribute to its core behavioural

597 impairments.

**Materials and Methods**

*Subjects*

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

*Pharmacological manipulation*

29

621     To reduce noradrenaline functioning, we administered 40mg of the non-selective β-

622     adrenoceptor antagonist propranolol 60 minutes before the task (Fig 1D). To reduce dopamine

623     functioning, we administered 400mg of the selective D2/D3 antagonist amisulpride 90 minutes

624     before the task. Because of different pharmacokinetic properties, drugs were administered at

625     different times. Each drug group received the drug on its corresponding time point and a placebo

626     at the other time point. The placebo group received placebo at both time points, in line with our

627     previous studies (*36*, *59*, *75*).

628     *Experimental paradigm*

629     To quantify different exploration strategies, we developed a multi-armed bandit task

630     implemented using Cogent (http://www.vislab.ucl.ac.uk/cogent.php) for MATLAB (R2018a).

631     Subjects had to choose between bandits (i.e. trees) that produced samples (i.e. apples) with varying

632     reward (i.e. size) in two different horizon conditions (Figure 1a-b). Bandits were displayed during

633     the entire duration of a trial and there was no time limit for sampling from (choosing) the bandits.

634     The sizes of apples they collected were summed and converted to an amount of juice (feedback),

635     which was displayed during 2000 ms at the end of each trial. Subjects were instructed to endeavour

636     to make the most juice and that they would receive a cash bonus proportional to their performance.

637     Overall subjects received £10 per hour and a mean bonus of £1.12 (std: £0.06).

638     Similar to the horizon task (*7*), to induce different extents of exploration, we manipulated

639     the horizon (i.e. number of apples to be picked: 1 in the short horizon, 6 in the long horizon)

640     between trials. This horizon-manipulation, which has been extensively used to modulate

641     exploratory behaviour (*21*, *34*, *54*, *99*), promotes exploration in the long horizon condition as there

642     are more opportunities to gather reward.

643     Within a single trial, each bandit had a different mean reward $\mu$ (i.e. apple size) and

644     associated uncertainty as captured by the number of initial samples (i.e. number of apples shown

645     at the beginning of the trial). Each bandit (i.e. tree) $i$ was from one of four generative processes

646     (Figure 1c) characterised by different means $\mu_i$ and number of initial samples. The rewards (apple

647     sizes) for each bandit were sampled from a normal distribution with mean $\mu_i$, specific to the bandit,

648     and with a fixed variance, $S^2=0.8$. The rewards were those sampled values rounded to the closest

649     integer. Each distribution was truncated to [2, 10], meaning that rewards with values above or

650     below this interval were excluded, resulting in a total of 9 possible rewards (i.e. 9 different apple

651     sizes; cf. Figure 1 - Figure supplement 1 for a representation). The 'certain standard bandit'

652     provided three initial samples and on every trial its mean $\mu_{cs}$ was sampled from a normal

653     distribution: $\mu_{cs} \sim N(5.5, 1.4)$. The 'standard bandit' provided one initial sample and to make sure

654     that its mean $\mu_s$ was comparable to $\mu_{cs}$, the trials were split equally between the four following:

655     $\{\mu_s = \mu_{cs} + 1; \mu_s = \mu_{cs} - 1; \mu_s = \mu_{cs} + 2; \mu_s = \mu_{cs} - 2\}$. The 'novel bandit' provided no

656     initial samples and its mean $\mu_n$ was comparable to both $\mu_{cs}$ and $\mu_s$ by splitting the trials equally

657     between the eight following:$\{\mu_n = \mu_{cs} + 1; \mu_n = \mu_{cs} - 1; \mu_n = \mu_{cs} + 2; \mu_n = \mu_{cs} - 2; \mu_n =$

658     $\mu_s + 1; \mu_n = \mu_s - 1; \mu_n = \mu_s + 2; \mu_n = \mu_s - 2\}$. The 'low bandit' provided one initial sample

659     which was smaller than all the other bandits' means on that trial: $\mu_l = min(\mu_{cs}, \mu_s, \mu_n) - 1$. We

660     ensured that the initial sample from the low-value bandit was the smallest by resampling from each

661     bandit in the trials were that was not the case. To make sure that our task captures heuristic

662     exploration strategies, we simulated behaviour (cf. Figure 1). Additionally, in each trial, to avoid

663     that some exploration strategies overshadow other ones, only three of the four different groups

664     were available to choose from. Based on the mean of the initial samples, we identified the high-

665    value option (i.e. the bandit with the highest expected reward) in trials where both the certain-

666    standard and the standard bandit were present.

667         There were 25 trials of each of the four three-bandit combination making it a total of 100

668    different trials. They were then duplicated to measure choice consistency, defined as the frequency

669    of making the same choice on identical trials (in contrast to a previous propranolol study where

670    consistency was defined in terms of a value-based exploration parameter (*60*)). Each subject

671    played these 200 trials both in a short and in a long horizon setting, resulting in a total of 400 trials.

672    The trials were randomly assigned to one of four blocks and subjects were given a short break at

673    the end of each of them. To prevent learning, the bandits' positions (left, middle or right) as well

674    as their colour (8 sets of 3 different colours) where shuffled between trials. To ensure subjects

675    distinguished different apple sizes and understood that apples from the same tree were always of

676    similar size (generated following a normal distribution), they needed to undergo training prior to

677    the main experiment. In training, based on three displayed apples of similar size, they were tasked

678    to guess between two options, namely which apple was most likely to come from the same tree

679    and then received feedback about their choice.

680        *Statistical analyses*

681         All statistical analyses were performed using the R Statistical Software (*100*). For

682    computing ANOVA tests and pairwise comparisons the 'rstatix' package was used, and for

683    computing effect sizes the 'lsr' package (*101*) was used. To ensure consistent performance across

684    all subjects, we excluded one outlier subject (belonging to the amisulpride group) from our analysis

685    due to not engaging in the task and performing at chance level (defined as randomly sampling

686    from one out of three bandits, i.e. 33%). Each bandits' selection frequency for a horizon condition

687    was computed over all 200 trials and not only over the trials where this specific bandit was present

688    (i.e. 3/4 of 200 = 150 trials). In all the analysis comparing horizon conditions, except when looking

689    at score values (Figure 2c), only the 1st draw of the long horizon was used. We compared

690    behavioural measures and model parameters using (paired-samples) t-tests and repeated-measures

691    (rm-) ANOVAs with a between-subject factor of drug group (propranolol group, amisulpride

692    group, placebo group) and a within-subject factor horizon (long, short). Information seeking,

693    expected values and scores were analysed using rm-ANOVAS with a within-subject factor

694    horizon. Measures that were horizon-independent (e.g. prior mean), were analysed using one-way

695    ANOVAs with a between-subject factor drug group. As drug groups differed in negative affect

696    (cf. Appendix 2 Table 1), which, through its relationship to anxiety (*102*) is thought to affect

697    cognition (*103*) and potentially exploration (*104*). We corrected for negative affect (PANAS) and

698    IQ (WASI) in each analysis by adding those two measures as covariates in each ANOVA

699    mentioned above (cf. Appendix 1 for analysis without covariates and analysis with physiological

700    effect as an additional covariates). We report effect sizes using partial eta squared (η2) for

701    ANOVAs and Cohen's d (d) for t-tests (*105*).

702    *Computational modelling*

703        We adapted a set of Bayesian generative models from previous studies (*1*), where each

704    model assumed that different characteristics account for subjects' behaviour. The binary indicators

705    $(c_{tr}, c_n)$ indicate which components (value-free random and novelty exploration respectively)

706    were included in the different models. The value of each bandit is represented as a distribution

707    $N(Q, S)$ with $S = 0.8$, the sampling variance fixed to its generative value. Subjects have prior

708    beliefs about bandits' values which we assume to be Gaussian with mean $Q_0$ and uncertainty $\sigma_0$.

709    The subjects' initial estimate of a bandit's mean ($Q_0$; prior mean) and its uncertainty about it ($\sigma_0$;

710    prior variance) are free parameters.

33

711    These beliefs are updated according to Bayes rule (detailed below) for each initial sample (note

712    that there are no updates for the novel bandit).

713          *Mean and variance update rules*

714          At each time point $t$, in which a sample $m$, of one of the bandits is presented, the expected

715    mean $Q$ and precision $\tau = \frac{1}{\sigma^2}$ of the corresponding bandit $i$ are updated as follows:

716
$$Q_{i,t+1} = \frac{\tau_{i,t} * Q_{i,t} + \tau_{samp} * m}{\tau_{i,t} + \tau_{samp}}$$

717
$$\tau_{t+1}^i = \tau_{samp} + \tau_t^i$$

718    where $\tau_{samp} = \frac{1}{S^2}$ is the sampling precision, with the sampling variance $S = 0.8$ fixed. Those

719    update rules are equivalent to using a Kalman filter (*106*) in stationary bandits.

720          We examined three base models: the UCB model, the Thompson model and the hybrid

721    model. The UCB model encompasses the UCB algorithm (captures directed exploration) and a

722    softmax choice function (captures a value-based random exploration). The Thompson model

723    reflects Thompson sampling (captures an uncertainty-driven value-based random exploration).

724    The hybrid model captures the contribution of the UCB model and the Thompson model,

725    essentially a mixture of the above. We computed three extensions of each model by either adding

726    value-free random exploration $(c_{tr}, c_n) = (1,0)$, novelty exploration $(c_{tr}, c_n) = (0,1)$ or both

727    heuristics $(c_{tr}, c_n) = (1,1)$, leading to a total of 12 models (see the labels on the x-axis in Figure

728    4a; $(c_{tr}, c_n) = (0,0)$ is the model with no extension). For additional models cf. Appendix 1. A

729    coefficient $c_{tr}=1$ indicates that a $\epsilon$-greedy component was added to the decision rule, ensuring that

730    once in a while (every $\epsilon$ % of the time), another option than the predicted one is selected. A

731    coefficient $c_n=1$ indicates that the novelty bonus $\eta$ is added to the computation of the value of

732    novel bandits and the Kronecker delta $\delta$ in front of this bonus ensures that it is only applied to the

733     novel bandit. The models and their free parameters (summarised in Appendix 2 Table 5) are

734     described in detail below.

735     *Choice rules*

736     *UCB model.* In this model, an information bonus $\gamma$ is added to the expected reward of each option,

737     scaling with the option's uncertainty (UCB). The value of each bandit $i$ at timepoint $t$ is:

738
$$V_{i,t} = Q_{i,t} + \gamma\sigma_{i,t} + c_n\eta\delta_{[i=novel]}$$

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

740
$$P(c_t = i) = \frac{e^{\beta V_{i,t}}}{\sum_x e^{\beta V_{i,t}}} * (1 - c_{tr}\epsilon) + c_{tr}\frac{\epsilon}{3}$$

741         where $\beta$ is the inverse temperature of the softmax (lower values producing more

742     stochasticity), and the coefficient $c_{tr}$ adds the value-free random exploration component.

743     *Thompson model.* In this model, based on Thompson sampling, the overall uncertainty can be seen

744     as a more refined version of a decision temperature (*1*). The value of each bandit $i$ is as before:

745
$$V_{i,t} = Q_{i,t} + c_n\eta\delta_{[i=novel]}$$

746         A sample $x_{i,t} \sim N(V_{i,t}, \sigma_{i,t}^2)$ is taken from each bandit. The probability of choosing a bandit

747     $i$ depends on the probability that all pairwise differences between the sample from bandit $i$ and the

748     other bandits $j \neq i$ were greater or equal to 0 (see the probability of maximum utility choice rule

749     (*107*)). In our task, because three bandits were present, two pairwise differences scores (contained

750     in the two-dimensional vector u) were computed for each bandit. The probability of choosing

751     bandit $i$ is:

752
$$P(c_t = i) = P(\forall j: x_{i,t} > x_{j,t}) * (1 - c_{tr}\epsilon) + c_{tr}\frac{\epsilon}{3}$$

753
$$= \int_0^\infty \int_0^\infty \phi(u; M_{i,t}, C_{i,t}) \, du \, * (1 - c_{tr}\epsilon) + c_{tr}\frac{\epsilon}{3}$$

754 where $\phi$ is the multivariate Normal density function with mean vector

756
$$M_{i,t} = A_i \begin{pmatrix} V_{1,t} \\ V_{2,t} \\ V_{3,t} \end{pmatrix}$$

755 and covariance matrix

757
$$C_{i,t} = A_i \begin{pmatrix} \sigma_{1,t} & 0 & 0 \\ 0 & \sigma_{2,t} & 0 \\ 0 & 0 & \sigma_{3,t} \end{pmatrix} A_i^T$$

758 Where the matrix $A_i$ computes the pairwise differences between bandit $i$ and the other bandits. For

759 example, for bandit $i = 1$:

760
$$A_1 = \begin{pmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{pmatrix}$$

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

762 probability of choosing bandit $i$ is:

763
$$P(c_t = i) = \left( w P_{UCB}(c_t = i) + (1 - w) P_{Thompson}(c_t = i) \right) * (1 - c_{tr}\epsilon) + c_{tr}\frac{\epsilon}{3}$$

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

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

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

767 All the parameters besides $Q_0$ and $w$ were free to vary as a function of the horizon (cf.

768 Appendix 2 Table 5) as they capture different exploration forms: directed exploration (information

769 bonus $\gamma$; UCB model), novelty exploration (novelty bonus $\eta$), random exploration (inverse

770     temperature $\beta$; UCB model), uncertainty-directed exploration (prior variance $\sigma_0$; Thompson

771     model) and value-free random exploration ($\epsilon$-greedy parameter). The prior mean $Q_0$ was fitted to

772     both horizons together as we do not expect the belief of how good a bandit is to depend on the

773     horizon. The same was done for $w$ as assume the arbitration between the UCB model and the

774     Thompson model does not depend on horizon.

775     *Parameter estimation.*

776     To fit the parameter values, we used the maximum a posteriori probability (MAP) estimate. The

777     optimisation function used was fmincon in MATLAB. The parameters could vary within the

778     following bounds: $\sigma_0 = [0.01, 6], Q_0 = [1, 10], \epsilon = [0, 0.5], \eta = [0, 5]$. The prior distribution

779     used for the prior mean parameter $Q_0$ was the normal distribution: $Q_0 \sim N(5, 2)$ that approximates

780     the generative distributions. For the $\epsilon$-greedy parameter, the novelty bonus $\eta$ and the prior variance

781     parameter $\sigma_0$, a uniform distribution (of range equal to the specific parameters' bounds) was used,

782     which is equivalent to performing MLE. A summary of the parameter values per group and per

783     horizon can be found in Appendix 2 Table 6.

784     *Model comparison.*

785     We performed a K-fold cross-validation with $K = 10$. We partitioned the data of each subject

786     ($N_{trials}$ =400; 200 in each horizon) into K folds (i.e. subsamples). For model fitting in our model

787     selection, we used maximum likelihood estimation (MLE), where we maximised the likelihood

788     for each subject individually (fmincon was ran with 8 randomly chosen starting point to overcome

789     potential local minima). We fitted the model using K-1 folds and validated the model on the

790     remaining fold. We repeated this process K times, so that each of the K fold is used as a validation

791     set once, and averaged the likelihood over held out trials. We did this for each model and each

792     subject and averaged across subjects. The model with the highest likelihood of held-out data (the

793     winning model) was the Thompson sampling with $(c_{tr}, c_n) = \{1,1\}$. It was also the model which

794     accounted best for the largest number of subjects (Figure 4 – Figure supplement 1).

795     *Parameter recovery.*

796     To make sure that the parameters are interpretable, we performed a parameter recovery analysis.

797     For each parameter, we took 4 values, equally spread, within a reasonable parameter range ($\sigma_0 =$

798     $[0.5, 2.5], Q_0 = [1, 6], \epsilon = [0, 0.5], \eta = [0, 5]$). All parameters but $Q_0$ were free to vary as a

799     function of the horizon. We simulated behaviour with one artificial agent for each $4^7$ combinations

800     using a new trial for each. The model was fitted using MAP estimation (cf. Parameter estimation)

801     and analysed how well the generative parameters (generating parameters in Figure 5) correlated

802     with the recovered ones (fitted parameters in Figure 5) using Pearson correlation (summarised in

803     Figure 5c). In addition to the correlation we examined the spread (Figure 4 – Figure supplement

804     3) of the recovered parameters. Overall the parameters were well recoverable.

805     *Model validation*

806     To validate our model, we used each subjects' fitted parameters to simulate behaviour on our task

807     (4000 trials per agent). The stimulated data (Figure 5 – Figure supplement 1), although not perfect,

808     resembles the real data reasonably well. Additionally, to validate the behavioural indicators of the

809     two different exploration heuristics we stimulated the behaviour of 200 agents using the winning

810     model on one horizon condition (i.e. trials = 200). For the indicators of value-free random

811     exploration, we stimulated behaviour with low ($\epsilon = 0$) and high ($\epsilon = 0.2$) values of the $\epsilon$-greedy

812     parameter. The other parameters were set to the mean parameter fits ($\sigma_0 = 1.312, \eta = 2.625, Q_0 =$

813     3.2). This confirms that higher amounts of value-free random exploration are captured by the

814   proportion of low-value bandit selection (Figure 1f) and the choice consistency (Figure 1e).

815   Similarly, for the indicator of novelty exploration, we simulated behaviour with low ($\eta = 0$) and

816   high ($\eta = 2$) values of the novelty bonus $\eta$ to validate the use of the proportion of the novel-bandit

817   selection (Figure 1g). Again, the remaining parameters were set to the mean parameter fits ($\sigma_0 =$

818   $1.312, \epsilon = 0.1, Q_0 = 3.2$). Parameter values for high and low exploration were selected

819   empirically from pilot and task data. Additionally, we simulated the effects of other exploration

820   strategies in short and long horizon conditions (Figure 1 – Figure supplement 3-5). To simulate a

821   long (versus short) horizon condition we increased the overall exploration by increasing other

822   exploration strategies. Details about parameter values can be found in Appendix 2 Table 7.

## References

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

17.  W. Kool, J. T. McGuire, Z. B. Rosen, M. M. Botvinick, Decision Making and the Avoidance of Cognitive Demand. *J. Exp. Psychol. Gen.* **139**, 665–682 (2010).

18.  R. Cools, The cost of dopamine for dynamic cognitive control. *Curr. Opin. Behav. Sci.* (2015), , doi:10.1016/j.cobeha.2015.05.007.

19.  M. I. Froböse, R. Cools, Chemical neuromodulation of cognitive control avoidance. *Curr. Opin. Behav. Sci.* **22**, 121–127 (2018).

20.  R. S. Sutton, A. G. Barto, Introduction to Reinforcement Learning. *MIT Press Cambridge* (1998), doi:10.1.1.32.7692.

21.  W. K. Zajkowski, M. Kossut, R. C. Wilson, A causal role for right frontopolar cortex in directed, but not random, exploration. *Elife*. **6**, 1–18 (2017).

22.  T. U. Hauser, V. G. Fiore, M. Moutoussis, R. J. Dolan, Computational Psychiatry of ADHD: Neural Gain Impairments across Marrian Levels of Analysis. *Trends Neurosci.* **39**, 63–73 (2016).

23.  S. J. Sara, A. Vankov, A. Hervé, Locus coeruleus-evoked responses in behaving rats: A clue to the role of noradrenaline in memory. *Brain Res. Bull.* **35**, 457–465 (1994).

24.  C. Varazzani, A. San-Galli, S. Gilardeau, S. Bouret, Noradrenaline and dopamine neurons in the reward/effort trade-off: A direct electrophysiological comparison in behaving monkeys. *J. Neurosci.* **35**, 7866–7877 (2015).

25.  M. Silvetti, E. Vassena, E. Abrahamse, T. Verguts, *Dorsal anterior cingulate-brainstem ensemble as a reinforcement meta-learner* (2018), vol. 14.

26.  M. Silvetti, R. Seurinck, M. E. van Bochove, T. Verguts, The influence of the noradrenergic system on optimal control of neural plasticity. *Front. Behav. Neurosci.* **7**, 1–6 (2013).

27.  A. J. Yu, P. Dayan, Uncertainty, neuromodulation, and attention. *Neuron*. **46**, 681–692 (2005).

28.  M. R. Nassar, K. M. Rumsey, R. C. Wilson, K. Parikh, B. Heasly, J. I. Gold, Rational regulation of learning dynamics by pupil-linked arousal systems. *Nat. Neurosci.* (2012), doi:10.1038/nn.3130.

29.  J. David Johnson, Noradrenergic control of cognition: global attenuation and an interrupt function. *Med. Hypotheses* (2003), doi:10.1016/s0306-9877(03)00021-5.

30.  S. Bouret, S. J. Sara, Network reset: A simplified overarching theory of locus coeruleus noradrenaline function. *Trends Neurosci.* (2005), doi:10.1016/j.tins.2005.09.002.

31.  P. Dayan, A. J. Yu, Phasic norepinephrine: A neural interrupt signal for unexpected events. *Netw. Comput. Neural Syst.* **17**, 335–350 (2006).

32.  D. G. R. Tervo, M. Proskurin, M. Manakov, M. Kabra, A. Vollmer, K. Branson, A. Y. Karpova, Behavioral variability through stochastic choice and its gating by anterior cingulate cortex. *Cell*. **159**, 21–32 (2014).

33.  C. I. Jahn, S. Gilardeau, C. Varazzani, B. Blain, J. Sallet, M. E. Walton, S. Bouret, Dual contributions of noradrenaline to behavioural flexibility and motivation, 2687–2702 (2018).

34.  C. M. Warren, R. C. Wilson, N. J. Van Der Wee, E. J. Giltay, M. S. Van Noorden, J. D. Cohen, S. Nieuwenhuis, The effect of atomoxetine on random and directed exploration in humans. *PLoS One* (2017), doi:10.1371/journal.pone.0176034.

35.  P. F. Fraundorfer, R. H. Fertel, D. D. Miller, D. R. Feller, Biochemical and pharmacological characterization of high-affinity trimetoquinol analogs on guinea pig and human beta adrenergic receptor subtypes: Evidence for partial agonism. *J. Pharmacol. Exp. Ther.* **270**, 665–674 (1994).

36.  T. U. Hauser, E. Eldar, N. Purg, M. Moutoussis, R. J. Dolan, Distinct roles of dopamine and noradrenaline in incidental memory. *J. Neurosci.* (2019), doi:10.1523/jneurosci.0401-19.2019.

37.  Ki Database, (available at https://pdsp.unc.edu/databases/pdsp.php).

38.  N. Bunzeck, C. F. Doeller, R. J. Dolan, E. Duzel, Contextual interaction between novelty and reward processing within the mesolimbic system. *Hum. Brain Mapp.* (2012), doi:10.1002/hbm.21288.

39.  B. C. Wittmann, N. D. Daw, B. Seymour, R. J. Dolan, Striatal Activity Underlies Novelty-Based Choice in Humans. *Neuron* (2008), doi:10.1016/j.neuron.2008.04.027.

40.  S. J. Gershman, Y. Niv, Novelty and Inductive Generalization in Human Reinforcement Learning. *Top. Cogn. Sci.* **7**, 391–415 (2015).

41.  H. Stojic, E. Shulz, P. P. Analytis, M. Speekenbrink, It's new, but is it good? How generalization and uncertainty guide the exploration of novel options. *PsyArXiv* (2018).

42.  R. M. Krebs, B. H. Schott, H. Schütze, E. Düzel, The novelty exploration bonus and its attentional modulation. *Neuropsychologia* (2009), doi:10.1016/j.neuropsychologia.2009.01.015.

43.  M. J. Frank, B. B. Doll, J. Oas-Terpstra, F. Moreno, Prefrontal and striatal dopaminergic genes predict individual differences in exploration and exploitation. *Nat. Neurosci.* **12**, 1062–1068 (2009).

44.  E. S. Bromberg-Martin, M. Matsumoto, O. Hikosaka, Dopamine in Motivational Control: Rewarding,

927             Aversive, and Alerting. *Neuron* (2010), , doi:10.1016/j.neuron.2010.11.022.

928   45.   V. D. Costa, V. L. Tran, J. Turchi, B. B. Averbeck, Dopamine modulates novelty seeking behavior during
929         decision making. *Behav. Neurosci.* (2014), doi:10.1037/a0037128.

930   46.   E. Düzel, W. D. Penny, N. Burgess, Brain oscillations and memory. *Curr. Opin. Neurobiol.* (2010), ,
931         doi:10.1016/j.conb.2010.01.004.

932   47.   K. Iigaya, T. U. Hauser, Z. Kurth-Nelson, J. P. O'Doherty, P. Dayan, R. J. Dolan, The value of what's to
933         come: neural mechanisms coupling prediction error and reward anticipation. *bioRxiv* (2019),
934         doi:10.1101/588699.

935   48.   A. S. Kayser, J. M. Mitchell, D. Weinstein, M. J. Frank, Dopamine, locus of control, and the exploration-
936         exploitation tradeoff. *Neuropsychopharmacology* (2015), doi:10.1038/npp.2014.193.

937   49.   D. Wechsler, WASI -II: Wechsler abbreviated scale of intelligence - second edition. *J. Psychoeduc. Assess.*
938         (2013), doi:10.1177/0734282912467756.

939   50.   D. Watson, L. A. Clark, A. Tellegen, Development and Validation of Brief Measures of Positive and
940         Negative Affect: The PANAS Scales. *J. Pers. Soc. Psychol.* (1988), doi:10.1037/0022-3514.54.6.1063.

941   51.   S. Agrawal, N. Goyal, Analysis of thompson sampling for the multi-armed bandit problem. *J. Mach. Learn.*
942         *Res.* **23**, 1–26 (2012).

943   52.   M. D'Acremont, P. Bossaerts, Neurobiological studies of risk assessment: A comparison of expected utility
944         and mean-variance approaches. *Cogn. Affect. Behav. Neurosci.* **8**, 363–374 (2008).

945   53.   N. C. Foley, D. C. Jangraw, C. Peck, J. Gottlieb, Novelty enhances visual salience independently of reward
946         in the parietal lobe. *J. Neurosci.* (2014), doi:10.1523/JNEUROSCI.4171-13.2014.

947   54.   C. M. Wu, E. Schulz, M. Speekenbrink, J. D. Nelson, B. Meder, Generalization guides human exploration in
948         vast decision spaces. *Nat. Hum. Behav.* **2**, 915–924 (2018).

949   55.   V. Skvortsova, S. Palminteri, M. Pessiglione, Learning to minimize efforts versus maximizing rewards:
950         Computational principles and neural correlates. *J. Neurosci.* **34**, 15621–15630 (2014).

951   56.   T. U. Hauser, E. Eldar, R. J. Dolan, Separate mesocortical and mesolimbic pathways encode effort and
952         reward learning signals. *Proc. Natl. Acad. Sci. U. S. A.* **114**, E7395–E7404 (2017).

953   57.   M. E. Walton, S. Bouret, What Is the Relationship between Dopamine and Effort? *Trends Neurosci.* **42**, 79–
954         91 (2019).

955   58.   J. D. Salamone, S. E. Yohn, L. López-Cruz, N. San Miguel, M. Correa, Activational and effort-related
956         aspects of motivation: Neural mechanisms and implications for psychopathology. *Brain*. **139**, 1325–1347
957         (2016).

958   59.   T. U. Hauser, M. Moutoussis, N. Purg, P. Dayan, R. J. Dolan, Beta-Blocker Propranolol Modulates Decision
959         Urgency During Sequential Information Gathering. *J. Neurosci.* (2018), doi:10.1523/jneurosci.0192-
960         18.2018.

961   60.   P. Sokol-Hessner, S. F. Lackovic, R. H. Tobe, C. F. Camerer, B. L. Leventhal, E. A. Phelps, Determinants
962         of Propranolol's Selective Effect on Loss Aversion. *Psychol. Sci.* **26**, 1123–1130 (2015).

963   61.   D. Campbell-Meiklejohn, J. Wakeley, V. Herbert, J. Cook, P. Scollo, M. K. Ray, S. Selvaraj, R. E.
964         Passingham, P. Cowen, R. D. Rogers, Serotonin and dopamine play complementary roles in gambling to
965         recover losses. *Neuropsychopharmacology*. **36**, 402–410 (2011).

966   62.   R. D. Rogers, M. Lancaster, J. Wakeley, Z. Bhagwagar, Effects of beta-adrenoceptor blockade on
967         components of human decision-making. *Psychopharmacology (Berl)*. **172**, 157–164 (2004).

968   63.   J. Rajkowski, P. Kubiak, G. Aston-Jones, Locus coeruleus activity in monkey: Phasic and tonic changes are
969         associated with altered vigilance. *Brain Res. Bull.* **35**, 607–616 (1994).

970   64.   G. Aston-Jones, J. D. Cohen, AN INTEGRATIVE THEORY OF LOCUS COERULEUS-
971         NOREPINEPHRINE FUNCTION: Adaptive Gain and Optimal Performance. *Annu. Rev. Neurosci.* **28**, 403–
972         450 (2005).

973   65.   D. Servan-Schreiber, H. Printz, J. D. Cohen, A network model of catecholamiine effects: Gain, signal-to-
974         noise ratio, and behavior. *Science (80-. )*. (1990), doi:10.1126/science.2392679.

975   66.   S. Joshi, Y. Li, R. M. Kalwani, J. I. Gold, Relationships between Pupil Diameter and Neuronal Activity in
976         the Locus Coeruleus, Colliculi, and Cingulate Cortex. *Neuron* (2016), doi:10.1016/j.neuron.2015.11.028.

977   67.   S. Joshi, J. I. Gold, Pupil Size as a Window on Neural Substrates of Cognition. *Trends Cogn. Sci.* **24**, 466–
978         480 (2020).

979   68.   V. Koudas, A. Nikolaou, E. Hourdaki, S. G. Giakoumaki, P. Roussos, P. Bitsios, Comparison of ketanserin,
980         buspirone and propranolol on arousal, pupil size and autonomic function in healthy volunteers.
981         *Psychopharmacology (Berl)*. (2009), doi:10.1007/s00213-009-1508-5.

982   69.   M. Jepma, S. Nieuwenhuis, Pupil diameter predicts changes in the exploration-exploitation trade-off:

Evidence for the adaptive gain theory. *J. Cogn. Neurosci.* (2011), doi:10.1162/jocn.2010.21548.

70. Jepma, The role of the noradrenergic system in the exploration-exploitation trade-off: a pharmacological study. *Front. Hum. Neurosci.* (2010), doi:10.3389/fnhum.2010.00170.

71. M. Usher, J. D. Cohen, D. Servan-Schreiber, J. Rajkowski, G. Aston-Jones, The role of locus coeruleus in the regulation of cognitive performance. *Science (80-. ).* **283**, 549–554 (1999).

72. G. A. Kane, E. M. Vazey, R. C. Wilson, A. Shenhav, N. D. Daw, G. Aston-Jones, J. D. Cohen, Increased locus coeruleus tonic activity causes disengagement from a patch-foraging task. *Cogn. Affect. Behav. Neurosci.* **17**, 1073–1083 (2017).

73. Z. L. Rossetti, S. Carboni, Noradrenaline and dopamine elevations in the rat prefrontal cortex in spatial working memory. *J. Neurosci.* **25**, 2322–2329 (2005).

74. M. E. Gibbs, D. S. Hutchinson, R. J. Summers, Noradrenaline release in the locus coeruleus modulates memory formation and consolidation; roles for α- and β-adrenergic receptors. *Neuroscience.* **170**, 1209–1222 (2010).

75. T. U. Hauser, M. Allen, N. Purg, M. Moutoussis, G. Rees, R. J. Dolan, Noradrenaline blockade specifically enhances metacognitive performance. *Elife* (2017), doi:10.7554/eLife.24901.

76. I. Trofimova, T. W. Robbins, Temperament and arousal systems: A new synthesis of differential psychology and functional neurochemistry. *Neurosci. Biobehav. Rev.* **64**, 382–402 (2016).

77. B. D. Waterhouse, H. C. Moises, H. H. Yeh, D. J. Woodward, Norepinephrine enhancement of inhibitory synaptic mechanisms in cerebellum and cerebral cortex: Mediation by beta adrenergic receptors. *J. Pharmacol. Exp. Ther.* **221**, 495–506 (1982).

78. B. D. Waterhouse, H. C. Moises, H. H. Yeh, H. M. Geller, D. J. Woodward, Comparison of norepinephrine- and benzodiazepine-induced augmentation of Purkinje cell response to γ-aminobutyric acid (GABA). *J. Pharmacol. Exp. Ther.* **228**, 257–267 (1984).

79. P. S. Goldman-Rakic, M. S. Lidow, D. W. Gallager, Overlap of dopaminergic, adrenergic, and serotoninergic receptors and complementarity of their subtypes in primate prefrontal cortex. *J. Neurosci.* **10**, 2125–2138 (1990).

80. J. S. Isaacson, M. Scanziani, How Inhibition Shapes Cortical Activity Excitation and inhibition walk hand in hand. *Neuron.* **72**, 231–243 (2011).

81. H. Salgado, M. Treviño, M. Atzori, Layer- and area-specific actions of norepinephrine on cortical synaptic transmission. *Brain Res.* **1641**, 163–176 (2016).

82. B. De Martino, B. A. Strange, R. J. Dolan, Noradrenergic neuromodulation of human attention for emotional and neutral stimuli. *Psychopharmacology (Berl).* **197**, 127–136 (2008).

83. W. Schultz, P. Dayan, P. R. Montague, A neural substrate of prediction and reward. *Science (80-. ).* **275**, 1593–1599 (1997).

84. B. C. Wittmann, N. D. Daw, B. Seymour, R. J. Dolan, Striatal Activity Underlies Novelty-Based Choice in Humans. *Neuron.* **58**, 967–973 (2008).

85. V. D. Costa, V. L. Tran, J. Turchi, B. B. Averbeck, Dopamine modulates novelty seeking behavior during decision making. *Behav. Neurosci.* **128**, 556–566 (2014).

86. A. S. Kayser, J. M. Mitchell, D. Weinstein, M. J. Frank, Dopamine, locus of control, and the exploration-exploitation tradeoff. *Neuropsychopharmacology.* **40**, 454–462 (2015).

87. K. Chakroun, D. Mathar, A. Wiehler, F. Ganzer, J. Peters, Dopaminergic modulation of the exploration/exploitation trade-off in human decision-making. *bioRxiv*, 706176 (2019).

88. F. Cinotti, V. Fresno, N. Aklil, E. Coutureau, B. Girard, A. R. Marchand, M. Khamassi, Dopamine blockade impairs the exploration-exploitation trade-off in rats. *Sci. Rep.* **9**, 1–14 (2019).

89. M. J. Frank, B. B. Doll, J. Oas-Terpstra, F. Moreno, Prefrontal and striatal dopaminergic genes predict individual differences in exploration and exploitation. *Nat. Neurosci.* (2009), doi:10.1038/nn.2342.

90. L. K. Krugel, G. Biele, P. N. C. Mohr, S. C. Li, H. R. Heekeren, Genetic variation in dopaminergic neuromodulation influences the ability to rapidly and flexibly adapt decisions. *Proc. Natl. Acad. Sci. U. S. A.* (2009), doi:10.1073/pnas.0905191106.

91. J. Fang, P. H. Yu, Effect of haloperidol and its metabolites on dopamine and noradrenaline uptake in rat brain slices. *Psychopharmacology (Berl).* (1995), doi:10.1007/BF02246078.

92. M. Toru, M. Takashima, Haloperidol in large doses reduces the cataleptic response and increases noradrenaline metabolism in the brain of the rat. *Neuropharmacology* (1985), doi:10.1016/0028-3908(85)90079-6.

93. T. Kahnt, S. C. Weber, H. Haker, T. W. Robbins, P. N. Tobler, Dopamine D2-Receptor Blockade Enhances Decoding of Prefrontal Signals in Humans. *J. Neurosci.* (2015), doi:10.1523/jneurosci.4182-14.2015.

94. T. Kahnt, P. N. Tobler, Dopamine Modulates the Functional Organization of the Orbitofrontal Cortex. *J. Neurosci.* (2017), doi:10.1523/jneurosci.2827-16.2016.

95. M. D. Humphries, M. Khamassi, K. Gurney, Dopaminergic control of the exploration-exploitation trade-off via the basal ganglia. *Front. Neurosci.* (2012), doi:10.3389/fnins.2012.00009.

96. A. Soutschek, C. J. Burke, A. Raja Beharelle, R. Schreiber, S. C. Weber, I. I. Karipidis, J. Ten Velden, B. Weber, H. Haker, T. Kalenscher, P. N. Tobler, The dopaminergic reward system underpins gender differences in social preferences. *Nat. Hum. Behav.* (2017), doi:10.1038/s41562-017-0226-y.

97. A. Soutschek, G. Gvozdanovic, R. Kozak, S. Duvvuri, N. de Martinis, B. Harel, D. L. Gray, E. Fehr, A. Jetter, P. N. Tobler, Dopaminergic D1 Receptor Stimulation Affects Effort and Risk Preferences. *Biol. Psychiatry* (2019), doi:10.1016/j.biopsych.2019.09.002.

98. T. U. Hauser, R. Iannaccone, J. Ball, C. Mathys, D. Brandeis, S. Walitza, S. Brem, Role of the medial prefrontal cortex in impaired decision making in juvenile attention-deficit/hyperactivity disorder. *JAMA Psychiatry*. **71**, 1165–1173 (2014).

99. D. Guo, A. J. Yu, in *Advances in Neural Information Processing Systems* (2018).

100. R. R Development Core Team, *R: A Language and Environment for Statistical Computing* (2011).

101. D. Navarro, *Learning statistics with R: A tutorial for psychology students and other beginners. (Version 0.5)* (2015; http://ua.edu.au/ccs/teaching/lsr).

102. D. Watson, L. A. Clark, G. Carey, Positive and Negative Affectivity and Their Relation to Anxiety and Depressive Disorders. *J. Abnorm. Psychol.* **97**, 346–353 (1988).

103. S. J. Bishop, C. Gagne, Anxiety, Depression, and Decision Making: A Computational Perspective. *Annu. Rev. Neurosci.* **41**, 371–388 (2018).

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, Trait anxiety affects decision-making differently in healthy men and women: Towards gender-specific endophenotypes of anxiety. *Neuropsychologia*. **48**, 1598–1606 (2010).

105. J. T. E. Richardson, Eta squared and partial eta squared as measures of effect size in educational research. *Educ. Res. Rev.* (2011), , doi:10.1016/j.edurev.2010.12.001.

106. C. M. Bishop, in *Information Science and Statistics* (2006).

107. M. Speekenbrink, E. Konstantinidis, Uncertainty and exploration in a restless bandit problem. *Top. Cogn. Sci.* (2015), doi:10.1111/tops.12145.

## Appendix 1

**Drug effect on response times**

There were no differences in response times (RT) between drug groups in the one-way ANOVA. Neither in the mean RT (ANOVA: $F(2, 54)=1.625$, $p=.206$, $\eta2=.057$) nor in its variability (standard deviation; $F(2, 54)=1.85$, $p=.16$, $\eta2=.064$).

**Bandit effect on response times**

There was no difference in response times between bandits in the repeated-measures ANOVA (bandit main effect: $F(1.78, 99.44)=1.634$, $p=.203$, $\eta2=.028$; Figure 3 – Figure supplement 1).

**Horizon effect on response times**

There were no differences in RT between horizon conditions in the repeated-measures ANOVA with the between-subject factor drug group, the within-subject factor horizon condition and the covariates WASI and PANAS negative score (horizon main effect: $F(1, 54)=1.443$, $p=.235$, $\eta2=.026$; drug main effect: $F(2, 54)=1.625$, $p=.206$, $\eta2=.057$; drug-by-horizon interaction: $F(2, 54)=.431$, $p=.652$, $\eta2=.016$. In the long horizon, the RT decreased with each sample (sample main effect: $F(1.36, 73.5)=13.626$, $p<.001$, $\eta2=0.201$; Pairwise comparisons: sample 1 vs 2: $t(59)=20.968$, $p<.001$, $d=2.73$; sample 2 vs 3: $t(59)=11.825$, $p<.001$, $d=1.539$; sample 3 vs 4: $t(59)=7.862$, $p<.001$, $d=1.024$; 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 – Figure supplement 1b).

**PANAS**

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

**Physiological effects**

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

**Task performance score**

The performance did not differ between drug groups (total score: drug main effect: $F(2, 5)=2.313$, $p=.109$, $\eta2=.079$) but it was increased in subjects with higher IQ scores (WASI main effect: $F(1, 54)=17.172$, $p<.001$, $\eta2=.241$).

In the long horizon, the score increased with each sample (sample main effect: $F(3.12, 174.97)=103.469$, $p<.001$, $\eta2=0.649$; Pairwise comparisons: sample 1 vs 2: $t(59)=-6.737$, $p<.001$, $d=0.877$; sample 2 vs 3: $t(59)=-3.69$, $p<.001$, $d=0.48$; sample 3 vs 4: $t(59)=-5.167$, $p<.001$, $d=0.673$; sample 4 vs 5: $t(59)=-2.832$, $p=.006$, $d=0.48$; sample 5 vs 6: $t(59)=-2.344$, $p=.022$, $d=0.673$; Figure 2 – Figure supplement 1a). The increase in reward was larger in trials where the first draw was exploratory (linear regression slope coefficient: mean=0.118, sd=0.038) compared to when it was exploitative (linear regression slope coefficient: mean=0.028, sd=0.041; t-tests for slope coefficients: $t(58)=-12.161$, $p<.001$, $d=-1.583$; Figure 2 - Figure supplement 1d), suggesting that exploration was used beneficially and subjects benefitted from their initial exploration.

45

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

The amisulpride group had a marginal tendency towards selecting the high-value bandit, meaning that they were disposed to exploit more overall (propranolol group excluded: horizon main effect: $F_{(1, 35)}=3.035$, p=.09, η2=.08; drug main effect: $F_{(1, 35)}=3.602$, p=.066, η2=.093; drug-by-horizon interaction: $F_{(1, 35)}=2.15$, p=.151, η2=.058). This trend effect was not observed when all 3 groups were included (horizon main effect: $F_{(1, 54)}=3.909$, p=.053, η2=.068; drug main effect: $F_{(2, 54)}=1.388$, p=.258, η2=.049; drug-by-horizon interaction: $F_{(2, 54)}=.834$, p=.44, η2=.03).

**Gender effects**

When adding gender as a between-subjects variable in the repeated-measures ANOVAs, none of the main results changed. Interestingly, we observed a drug-by-gender interaction in the prior variance $\sigma_0$ (drug-by-gender interaction: $F_{(2, 51)}=5.914$, p=.005, η2=.188; Figure 5 – Figure supplement 2), driven by the fact that, female subjects in the placebo group had a larger average $\sigma_0$ (across both horizon conditions) compared to males ($t(20)=2.836$, p=.011, d=1.268), whereas male subjects have a larger $\sigma_0$ compared to females in the amisulpride group, ($t(19)=-2.466$, p=.025, d=1.124; propranolol group: $t(20)=-0.04$, p=.969, d=.018). This suggests that in a placebo setting, females are on average more uncertain about an option's expected value, whereas in a dopamine blockade setting males are more uncertain. Besides this effect, we observed a trend-level significance in response times (RT), driven primarily by female subjects tending to have a faster RT in the long horizon compared to male subjects (gender main effect: $F_{(1, 51)}=3.54$, p=.066, η2=.065).

**Horizon and drug effects without covariate**

When analysing the results without correcting for IQ (WASI) and negative affect (PANAS), similar results are obtained. The high-value bandit is picked more in the short-horizon condition indicating exploitation ($F_{(1, 56)}=44.844$, p<.001, η2=.445), whereas the opposite phenomenon is observed in the low-value bandit ($F_{(1, 56)}=24.24$, p<.001, η2=.302) and the novel bandit (horizon main effect: $F_{(1, 56)}=30.867$, p<.001, η2=.355), indicating exploration. In line with these results, the model parameters for value-free random exploration ($\epsilon$: $F_{(1, 56)}=10.362$, p=.002, η2=.156) and novelty exploration ($\eta$: $F_{(1, 56)}=38.103$, p<.001, η2=.405) are larger in the long compared to the short horizon condition. Additionally, noradrenaline blockade reduces value-free random exploration as can be seen in the two behavioural signatures, frequency of picking the low-value bandit ($F_{(2, 56)}=2.523$, p=.089, η2=.083; Pairwise comparisons: placebo vs propranolol: $t(40)=2.923$, p=.005, d=.654; amisulpride vs placebo: $t(38)=-.587$, p=.559, d=.133; amisulpride vs propranolol: $t(38)=2.171$, p=.034, d=.496), and in the consistency ($F_{(2, 56)}=3.596$, p=.034, η2=.114; Pairwise comparisons: placebo vs propranolol: $t(40)=-3.525$, p=.001, d=.788; amisulpride vs placebo: $t(38)=1.107$, p=.272, d=.251; amisulpride vs propranolol: $t(38)=-2.267$, p=.026, d=.514), as well as in the model parameter for value-free random exploration ($\epsilon$: $F_{(2, 56)}=3.205$, p=.048, η2=.103; Pairwise comparisons: placebo vs propranolol: $t(40)=3.177$, p=.002, d=.71; amisulpride vs placebo: $t(38)=.251$, p=.802, d=.057; amisulpride vs propranolol: $t(38)=2.723$, p=.009, d=.626).

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

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

**Other model results**

When analysing the fitted parameter values of both the 2nd winning model (UCB + $\epsilon$ + $\eta$) and 3rd winning model (hybrid + $\epsilon$ + $\eta$), similar results pertain. Thus, a value-free random exploration parameter was reduced following noradrenaline blockade in the 2nd winning model ($\epsilon$: $F_{(2, 54)}=4.503$, p=.016, $\eta^2$=.143; Pairwise comparisons: placebo vs propranolol: $t(38)=2.185$, p=.033, d=.386; amisulpride vs propranolol: $t(40)=1.724$, p=.089, d=.501; amisulpride vs placebo: $t(40)=-.665$, p=.508, d=.151) and was affected at a trend-level significance in the 3rd

46

1179  winning model ($\epsilon$: F(2, 54)=3.04, p=.056, $\eta^2$=.101). These results highlight our finding that value-free random
1180  exploration is modulated by noradrenaline and additionally demonstrates this is independent of the complex
1181  exploration strategy used as well as the value function.
1182
1183  **Bandit combination effect**
1184  Behavioural results were analysed additionally for each bandit combination separately. The high-value bandit was
1185  picked more when there was no novel bandit (pairwise comparisons: [certain-standard, standard, low] vs [certain-
1186  standard, standard, novel]: t(59)=-15.122, p<.001, d=1.969 ; [certain-standard, standard, low] vs [certain-standard,
1187  novel, low]: t(59)=12.905, p<.001, d=1.68; [certain-standard, standard, low] vs [standard, novel, low]: t(59)=18.348,
1188  p<.001, d=2.389), and less when its value was less certain ([standard, novel, low] vs [certain-standard, standard,
1189  novel]: t(59)=6.986, p<.001, d=.909; [standard, novel, low] vs [certain-standard, novel, low] : t(59)=5.44, p<.001,
1190  d=.708; bandit combination main effect: F(1.81, 101.33)=237.051, p<.001, $\eta^2$=.809; [certain-standard, standard,
1191  novel] vs [certain-standard, novel, low]: t(59)=.364, p=.717, d=.047; Figure 3 – Figure supplement 2a). The novel
1192  bandit was picked the most when the high-value bandit was less certain, then when the high-value bandit was more
1193  certain and it was picked the least when both certain and certain standard bandits were present ([standard, novel,
1194  low] vs [certain-standard, novel, low]: t(59)=-5.001, p<.001, d=.651; [standard, novel, low] vs [certain-standard,
1195  standard, novel]: t(59)=-9.414, p<.001, d=1.226; [certain-standard, novel, low] vs [certain-standard, standard,
1196  novel]: t(59)=-4.146, p<.001, d=.54; bandit combination main effect: F(2, 112)=42.44, p<.001, $\eta^2$=.431; Figure 3 –
1197  Figure supplement 2b). The low-value bandit was picked less when the high-value bandit was more certain ([certain-
1198  standard, novel, low] vs [certain-standard, standard, low]: t(59)=2.731, p=.008, d=.356; [certain-standard, novel,
1199  low] vs [standard, novel, low]: t(59)=-1.958, p=.055, d=.255; bandit combination main effect: F(1.66, 92.74)=4.534,
1200  p=.019, η2=.075; [certain-standard, standard, low] vs [standard, novel, low]: t(59)=1.32, p=.192, d=.172; Figure 3 –
1201  Figure supplement 2c).
1202
1203  **Other effects on choice consistency**
1204  Our results demonstrate a drug-by-horizon interaction on choice consistency (F(2, 54)=3.352, p=.042, $\eta^2$=.110;
1205  Figure 3c), mainly driven by the fact that frequency of selecting the same option is increased in the long (compared
1206  to the short) horizon in the amisulpride group, while there is no significant horizon difference in the other two drug
1207  groups (pairwise comparison for horizon effect: amisulpride group: t(19)=2.482, p=.023, d=.569; propranolol group:
1208  t(20)=-1.91, p=.071, d=.427; placebo group: t(20)=.505, p=.619, d=.113). It is not entirely clear why catecholamines
1209  would increase the differentiation between the horizon conditions and this relatively weak effect should be
1210  replicated before interpreting.
1211
1212  **Stand-alone heuristic models**
1213  We also analysed stand-alone heuristic models, in which there is no value computation (value of each bandit $i$: $V_i =$
1214  0). The held-out data likelihood for such heuristic model combined with novelty exploration had a mean of
1215  m=0.367 (sd=0.005). The model in which we added value-free random exploration on top of novelty exploration
1216  had a mean of m=0.384 (sd=0.006). These models performed poorly, although better than chance level. Importantly,
1217  adding value-free random exploration improved performance. This highlights that subjects' combine complex and
1218  heuristic modules in exploration.

1219 **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) | $F(2,56)=.404, p=.669, \eta^2=.014$ |
| Intellectual abilities | 22.8 (1.85) | 22.6 (3.70) | 24.37 (2.45) | $F(2,56)=2.337, p=.106, \eta^2=.077$ |
| Positive affect | 24.55 (8.99) | 28.90 (7.56) | 29.58 (10.21) | $F(2,56)=1.832, p=.170, \eta^2=.061$ |
| Negative affect | 10.65 (.81) | 12.75 (3.63) | 11.16 (1.71) | $F(2,56)=4.259, p=.019, \eta^2=.132$ |

1220 **Appendix 2 Table 1.**

1221 Characteristics of drug groups. The drug groups did not differ in gender, age, nor in intellectual abilities (adapted
1222 WASI matrix test). Groups differed in negative affect (PANAS), driven by a higher score in the placebo group
1223 (pairwise comparisons: placebo vs propranolol: $t(56)=2.801, p=.007, d=.799$; amisulpride vs placebo: $t(56)=-2.096$,
1224 $p=.041, d=.557$; amisulpride vs propranolol: $t(56)=.669, p=.506, d=.383$). For more details cf. Appendix 1. Mean
1225 (SD).

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

1226 **Appendix 2 Table 2.**

1227 Physiological effects on drug groups. The drug groups also differed in post-experiment heart rate, driven by lower
1228 values in the propranolol group (pairwise comparisons: placebo vs propranolol: t(55)=3.5, p=.001, d=1.293;
1229 amisulpride vs placebo: t(55)= -.394, p=.695, d=.119 ; amisulpride vs propranolol: t(55)=3.013, p=.004, d=.921). For
1230 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, $\eta^2=.258$ | |
| | **long** | 6.221 (0.379) | | |
| Initial samples | **short** | 1.282 (0.247) | $F(1, 56)=58.78$, p<.001, $\eta^2=.512$ | |
| | **long** | 1.084 (0.329) | | |
| Score (1st sample) | **short** | 5.904 (0.192) | $F(1, 56)=58.78$, p<.001, $\eta^2=.512$ | |
| | **long** | 5.82 (0.182) | | |
| Score (average) | **short** | 5.904 (0.192) | $F(1, 56)=103.759$, p<.001, $\eta^2=.649$ | |
| | **long** | 6.098 (0.222) | | |

1231    **Appendix 2 Table 3.**

1232    Table of statistics and behavioural values of Figure 2. All of those measures were modulated by the horizon condition.

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

**Appendix 2 Table 4.**

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

| | | Thompson | | | | UCB | | | | Hybrid | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **Model** | | $+\epsilon$ | $+\eta$ | $+\epsilon +\eta$ | | $+\epsilon$ | $+\eta$ | $+\epsilon +\eta$ | | $+\epsilon$ | $+\eta$ | $+\epsilon +\eta$ |
| **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 | $\sigma_0$ | $\sigma_0, \epsilon$ | $\sigma_0, \eta$ | $\sigma_0, \epsilon, \eta$ | $\gamma, \beta$ | $\gamma, \beta, \epsilon$ | $\gamma, \beta, \eta$ | $\gamma, \beta, \epsilon, \eta$ | $\sigma_0, \gamma, \beta$ | $\sigma_0, \gamma, \beta, \epsilon$ | $\sigma_0, \gamma, \beta, \eta$ | $\sigma_0, \gamma, \beta, \epsilon, \eta$ |
| **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 |

**Appendix 2 Table 5.**

Table of parameters used for each model compared during model selection (Figure 4). Each of the 12 columns indicate a model. The three 'main models' studied were the Thompson model, the UCB model and a hybrid of both. Variants were then created by adding the $\epsilon$-greedy parameter, the novelty bonus and a combination of both. All the parameters besides $Q_0$ and w were fitted to each horizon separately. Parameters: $Q_0$=prior mean (initial estimate of a bandits mean); $\sigma_0$=prior variance (uncertainty about $Q_0$); $w$=contribution of UCB vs Thompson; $\gamma$ =information bonus; $\beta$=softmax inverse temperature; $\epsilon$=$\epsilon$-greedy parameter (stochasticity); $\eta$=novelty bonus. Model selection measures include the cross-validation held-out data likelihood averaged over subjects, mean (SD), as well as the subject count for which this model performed better over either 12 models or over the 3 best models.

| | | | Mean (sd) | | | Two-way repeated-measures ANOVA | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | Horizon | Amisulpride | Placebo | Propranolol | Main effect | | Interaction | |
| $\epsilon$-greedy parameter | | short | 0.10 (0.10) | 0.12 (0.08) | 0.07 (0.08) | D | $F(2, 54)=6.722$, $p=.002$, $\eta^2=.199$ | DH | $F(2, 54)=1.305$, $p=.280$, $\eta^2=.046$ |
| | | long | 0.17 (0.14) | 0.14 (0.10) | 0.08 (0.06) | H | $F(1, 54)=1.968$, $p=.166$, $\eta^2=.035$ | HW | $F(1, 54)=6.08$, $p=.017$, $\eta^2=.101$ |
| Novelty bonus $\eta$ | | short | 2.07 (0.98) | 2.26 (1.37) | 2.05 (1.16) | D | $F(2, 54)=.249$, $p=.780$, $\eta^2=.009$ | DH | $F(2, 54)=.03$, $p=.971$, $\eta^2=.001$ |
| | | long | 3.24 (1.19) | 3.12 (1.63) | 2.95 (1.70) | H | $F(1, 54)=1.839$, $p=.181$, $\eta^2=.033$ | HW | $F(1, 54)=8.416$, $p=.005$, $\eta^2=.135$ |
| Prior variance $\sigma_0$ | | short | 1.18 (0.20) | 1.12 (0.43) | 1.25 (0.34) | D | $F(2, 54)=.060$, $p=.942$, $\eta^2=.002$ | DH | $F(2, 54)=2.162$, $p=.125$, $\eta^2=.074$ |
| | | long | 1.41 (0.61) | 1.42 (0.59) | 1.21 (0.44) | H | $F(1, 54)=.129$, $p=.721$, $\eta^2=.002$ | HW | $F(1, 54)=.022$, $p=.882$, $\eta^2<.001$ |
| Prior mean $Q_0$ | | | 3.22 (1.05) | 3.20 (1.36) | 3.44 (1.05) | D | $F(2, 54)=.118$, $p=.889$, $\eta^2=.004$ | | |

1250 **Appendix 2 Table 6.**

1251 Table of statistics and fitted model parameters of Figure 5. The drug groups differed in $\epsilon$-greedy parameter value
1252 (pairwise comparisons: placebo vs propranolol: $t(40)=3.177$, $p=.002$, $d=.71$; amisulpride vs placebo: $t(38)=.251$,
1253 $p=.802$, $d=.057$; amisulpride vs propranolol: $t(38)=2.723$, $p=.009$, $d=.626$). The main effect is either of drug group (D)
1254 or of horizon (H). The interaction is either drug-by-horizon (DH) or horizon-by-WASI (measure of IQ; HW).

53

| | **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** | $\gamma = 0.7$ | $\gamma = 1.5$ | $\beta = 1.5, \epsilon = 0.2$ |

1255     **Appendix 2 Table 7**

1256     Parameter values used for simulations on Figure 1- Figure supplement 3-5. Parameter values for high and low
1257     exploration were selected empirically from pilot and task data. Value-free random exploration and novelty exploration
1258     were simulated with an argmax decision function, which always selects the value with the highest expected value. For
1259     simulating the long (versus short) horizon condition, we assumed that not only the key value but also the other
1260     exploration strategies increased, as found in our experimental data. For each simulation $Q_0 = 5$ and unless otherwise
1261     stated, $\sigma_0 = 1.5$.

1262



1263 **Figure 1 - Figure supplement 1**

1264 Visualisation of the 9 different sizes that the apples could take. The associated rewards went from 2

1265 (small apple on the left) to 10 (big apple on the right).

1266

1267

**Figure 1 - Figure supplement 2**

Comparison of value-based (softmax) and value-free ($\epsilon$-greedy) random exploration. (a) Changing the softmax inverse temperature affects the slope of the sigmoid while changing the $\epsilon$-greedy parameter (b) affects the compression of the sigmoid. Conceptually, in a softmax exploration mode, as each bandits' expected value is taken into account, (c) the $2^{nd}$ best bandit (medium-value bandit) will be favoured over one with a lower value (low-value bandit) when injecting noise. In contrast, in an $\epsilon$-greedy exploration mode, (d) bandits are explored equally often irrelevant of their expected value. Both simulations were performed on trials without novel bandit. When simulating on all trials we see that this also has a consequence on choice consistency, as (e) the $2^{nd}$ best option will most probably be explored (i.e. choice is still more consistent) in a softmax exploration mode versus (f) equal probability of exploring any of the 2 non-optimal options in an $\epsilon$-greedy exploration mode.
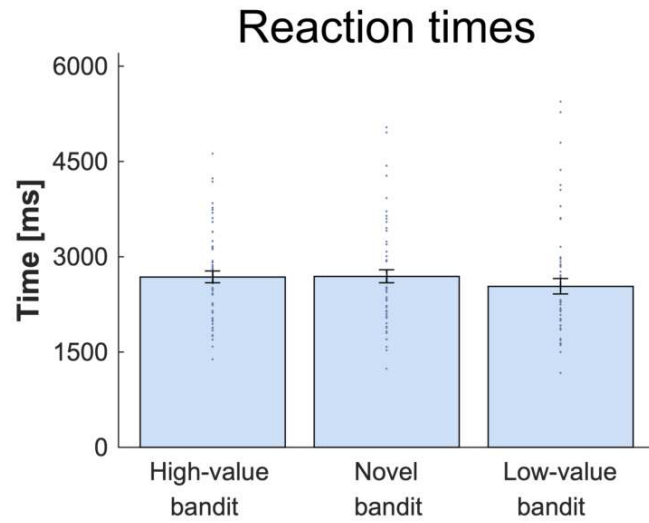
56

1279

**Figure 1 - Figure supplement 3**

Simulating the effect of the different exploration strategies on the frequency of picking the low-value bandit shows that (a) a higher value-free random exploration increases the selection of the low-value bandit, whereas neither (b) a higher novelty exploration, (c) a higher Thompson-sampling exploration nor (d) a higher UCB exploration affected this frequency. For simulating the long (versus short) horizon condition, we assumed that not only the key value but also the other exploration strategies increased, as found in our experimental data (cf. Appendix 2 Table 7 for parameter values).

1287

**Figure 1 - Figure supplement 4**

Simulating the effect of the different exploration strategies on choice consistency shows that (a) a higher value-free random exploration decreases the proportion of same choices, whereas neither (b) a higher novelty exploration, (c) a higher Thompson-sampling exploration nor (d) a higher UCB exploration affected this measure. For simulating the long (versus short) horizon condition, we assumed that not only the key value but also the other exploration strategies increased, as found in our experimental data (cf. Appendix 2 Table 7 for parameter values).

1294

**Figure 1 - Figure supplement 5**

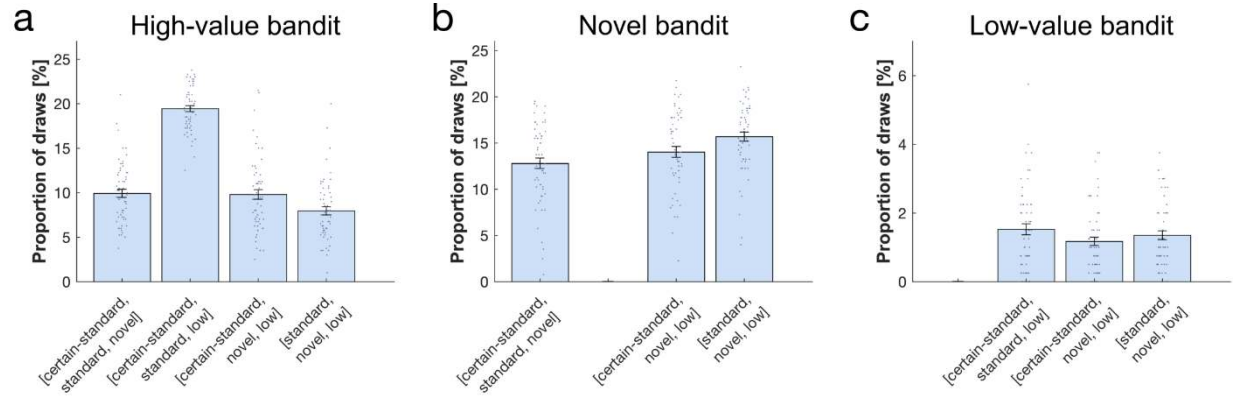Simulating the effect of the different exploration strategies on the frequency of picking the novel bandit shows that (a) a higher value-free random exploration has little effect on the selection of the novel bandit, whereas (b) a higher novelty exploration increases this frequency. (c) A higher Thompson-sampling exploration had little effect and (d) a higher UCB exploration affected this frequency but to a lower extend than novelty exploration. For simulating the long (versus short) horizon condition, we assumed that not only the key value but also the other exploration strategies increased, as found in our experimental data (cf. Appendix 2 Table 7 for parameter values).

1302

**Figure 2 - Figure supplement 1**

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

60

1314

1315 **Figure 3 - Figure supplement 1**

1316 Response time analysis per bandit. There was no difference in RT depending which bandit was chosen. For details
1317 and statistics cf. Appendix 1.

**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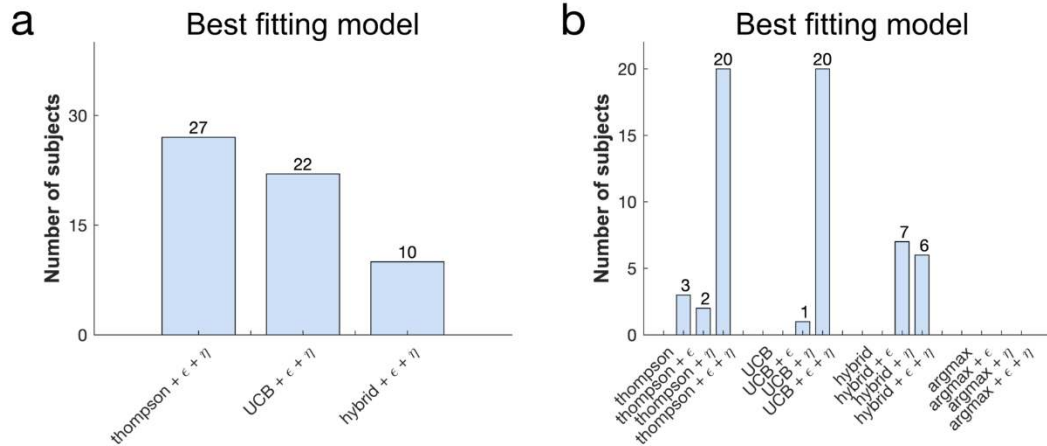 less certain. (b) The novel bandit was picked the most when the high-value bandit was less certain, then when the high-value bandit was more certain, and it was picked the least when both certain and certain standard bandits were present. (c) The low-value bandit was picked less when the high-value bandit was more certain. For statistics see Appendix 1.
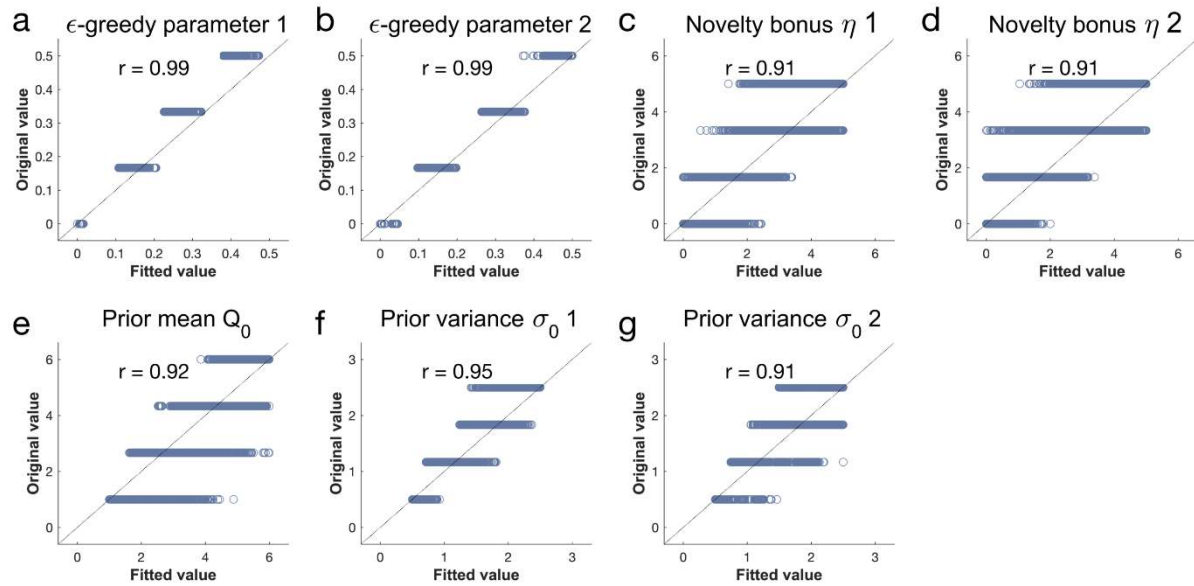
1325

**Figure 4 - Figure supplement 1**

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

1331

**Figure 4 - Figure supplement 2**
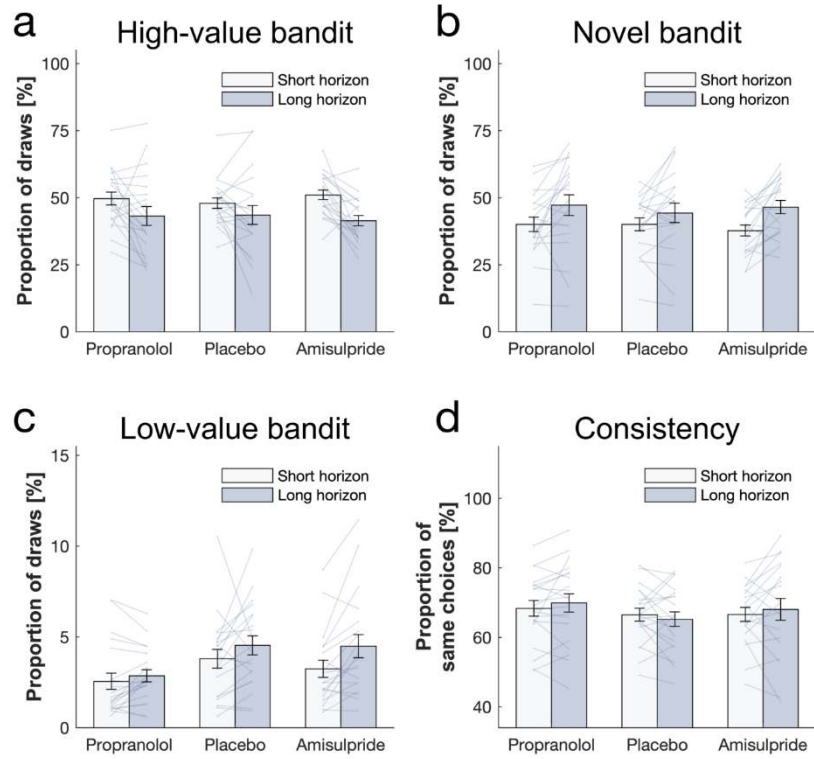
Correlations between model parameters and behaviour. The behavioural indicators of (a) value-free random exploration (left panel: draws from the low-value bandit; right panel: consistency) correlated with the $\epsilon$-greedy parameter values, and of (b) novelty exploration (draws from the novel bandit) correlated with the novelty bonus $\eta$.
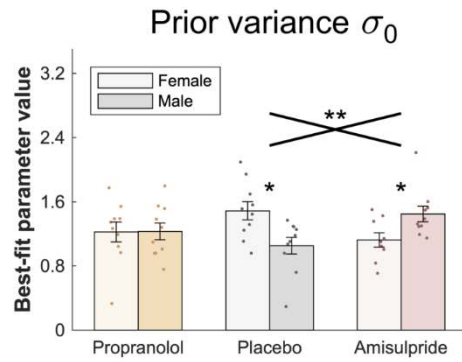
1336

**Figure 4 - Figure supplement 3**

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

1342

**Figure 5 - Figure supplement 1**

Simulated behaviour. We used each subjects' fitted parameters to simulate behaviour (blue diamonds; $N_{trials}$=4000) and superposed them to the real behaviour measures ($N_{trials}$=400) measures. Data are shown as mean ± SEM and each dot/line represent one agent.

1347

**Figure 5 – Figure supplement 2**

Gender effect on prior variance parameter. Mean values (across horizon conditions) of $\sigma_0$ were larger for female subjects, whereas in the amisulpride group, they were larger for male subjects. Data are shown as mean ± SEM and each dot represent one subject.